INLG 2023

The 16th International Natural Language Generation
Conference

Proceedings of the Conference

September 11 - 15, 2023

# Preface

We are excited to present the Proceedings of the 16th International Natural Language Generation Conference (INLG 2023). This year is the first time since the Covid-19 pandemic that the event will run mainly in-person again, from 11 to 15 September 2023 in Prague, Czech Republic. A novel aspect of this year's INLG is that, for the first time in its history, it is held jointly with the 24th Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDial 2023). INLG-SIGDIAL 2023 was locally organized by Charles University, thanks to the tireless efforts of the local chair Ondřej Dušek and his team.

The INLG conference is the main international venue for presentation of novel research and discussion of the computational task of Natural Language Generation (NLG) and its broad range of applications, including mainly data-to-text, text-to-text, and image-to-text approaches. Also this year, INLG consisted of several events.

The conference took place from 13 to 15 September. For the main track, we received a total of 98 conference submissions, 4 ARR submissions, and 4 demo paper submissions. After review by at least three program committee members and a meta review from the area chairs, 19 were accepted as long papers, 17 as short papers, and 4 as demo papers.

INLG, jointly with SIGDIAL, featured four keynote speakers, being:

- Barbara Di Eugenio, University of Illinois, Chicago, USA

- Emmanuel Dupoux, Ecole des Hautes Etudes en Sciences Sociales, France

- Ryan Lowe, OpenAI, USA

- Elena Simperl, King's College London, UK

The Generation Challenge, i.e., a set of shared tasks, was a track of the main conference also this year. It was chaired by Simon Mille. Details about the challenge and the proceedings will appear in a companion proceedings volume.

The main event was preceded by two days of workshops held jointly with SIGDIAL2023, of which two focussed on NLG, being the workshop on "Multimodal, Multilingual Natural Language Generation and Multilingual WebNLG Challenge" and a hackathon on practical "LLM-assisted data-to-text generation".

The event received sponsorship from: Liveperson and Luxai (Platinum), Apple (Gold), Furhat (Silver), and Bloomberg and Ax Semantics (Bronze).

It is also important to mention that the 16th INLG would not be possible without the help of the Area Chairs and Program Committee members for their reviewing contributions for whom we express our gratitude, and the expertise of SIGGEN representatives Raquel Hervás and Emiel van Miltenburg.

C. Maria Keet
Hung-yi Lee
Sina Zarrieß
INLG 2023 Program Chairs

# Organizing Committee

**Program Chairs**

    C. Maria Keet (University of Cape Town, South Africa)

    Hung-yi Lee (National Taiwan University, Taiwan)

    Sina Zarrieß (University of Bielefeld, Germany)

**Generation Challenge Chair**

    Simon Mille (ADAPT Research Centre, Dublin City University, Ireland)

**Local Organization Committee**

    Ondřej Dušek (Charles University, Czech Republic)

**SIGGEN Exec**

    Raquel Hervás (University Complutense of Madrid, Spain)

    Emiel van Miltenburg (Tilburg University, the Netherlands)

**Publication Chair**

    Chung-Chi Chen (National Institute of Advanced Industrial Science and Technology, Japan)

**Sponsor Chair**

    Ramesh Manuvinakurike (Intel Labs)

**Invited Speakers**

    Barbara Di Eugenio (University of Illinois Chicago, USA)

    Emmanuel Dupoux (Ecole des Hautes Etudes en Sciences Sociales, France)

    Ryan Lowe (OpenAI, USA)

    Elena Simperl (King's College London, UK)

# Organizing Committee

**Area Chairs**

    Suma Bhat (University of Illinois at Urbana-Champaign)

    Joan Byamugisha (IBM Research)

    Brian Davis (Dublin City University)

    Albert Gatt (Utrecht University)

    Yufang Hou (IBM Research)

    Wei-Yun Ma (Academia Sinica)

    Lara Martin (University of Maryland)

    Samira Shaikh (University of North Carolina)

    Kees van Deemter (Utrecht University)

    Chris van der Lee (Tilburg University)


**Program Committee**

    Manex Agirrezabal (University of Copenhagen)

    Mary-Jane Antia (University of Cape Town)

    Vinayshekhar Bannihatti Kumar (AWS AI)

    Anya Belz (ADAPT Research Centre, Dublin City University)

    Raffaella Bernardi (University of Trento)

    Jennifer Biggs (U.S. Naval Research Laboratory)

    Nadjet Bouayad-Agha (Universitat Oberta de Catalunya)

    Daniel Braun (University of Twente)

    Gordon Briggs (U.S. Naval Research Laboratory)

    Alberto Bugarín-Diz (Universidad Santiago de Compostela)

    Jan Buys (University of Cape Town)

    Michele Cafagna (University of Malta)

    Deng Cai (Tencent AI Lab)

    Eduardo Calò (Utrecht University)

    Thiago Castro Ferreira (Federal University of Minas Gerais)

    Khyathi Raghavi Chandu (Allen Institute of AI)

    Cheng-Han Chiang (National Taiwan University)

    Yagmur Gizem Cinar (Amazon)

    Elizabeth Clark (Google Research)

    Nina Dethlefs (University of Hull)

    Simon Dobnik (University of Gothenburg)

    Farhood Farahnak (Concordia University)

    Nicolas Garneau (Universite Laval)

    Pablo Gervás (Universidad Complutense de Madrid)

    Martijn Goudbeek (Tilburg University)

    Mika Hämäläinen (Rootroo Ltd)

    Ting Han (Smartnews Inc.)

    Aki Harma (Philips Research)

    Hiroaki Hayashi (Salesforce Research)

    Philipp Heinisch (Bielefeld University)

    Po-chun Hsu (National Taiwan University)

    Nikolai Ilinykh (University of Gothenburg)

    Takumi Ito (Tohoku University / Langsmith Inc. / Utrecht University)

    Mihir Kale (Google)

    Zdeněk Kasner (Charles University)

    Natthawut Kertkeidkachorn (Japan Advanced Institute of Science and Technology)

    Emiel Krahmer (Tilburg University)

Cyril Labbe (Université Grenoble Alpes)
Luc Lamontagne (Laval University)
Maurice Langner (Ruhr-Universität Bochum)
Yucheng Li (University of Surrey)
Yizhi Li (University of Sheffield)
Michela Lorandi (Dublin City University)
Saad Mahamood (Trivago N.V)
Zola Mahlaza (University of Cape Town)
Aleksandre Maskharashvili (Ohio State University)
Kathleen McCoy (University of Delaware)
David McDonald (Smart Information Flow Technologies)
Antonio Valerio Miceli Barone (The University of Edinburgh)
Fatemehsadat Mireshghallah (UC San Diego)
Ryo Nagata (Konan University)
Christina Niklaus (University of St. Gallen)
Avinesh P.V.S (Apple Inc.)
Patrizia Paggio (University of Copenhagen, University of Malta)
Daniel Paiva (Arria NLG)
Suraj Pandey (Open University UK)
Steffen Pauws (Philips Research)
Pablo Perez De Angelis (tuQuejaSuma.com)
Paul Piwek (Open University UK)
François Portet (Université Grenoble Alpes)
Toky Raboanary (University of Cape Town)
Philipp Sadler (University of Potsdam)
Daniel Sanchez (University of Granada)
Sashank Santhanam (University of North Carolina at Charlotte, Apple)
David Schlangen (University of Potsdam)
Simeon Schüz (Bielefeld University)
Balaji Vasan Srinivasan (Adobe Research, India)
Somayajulu Sripada (Arria NLG Plc, University of Aberdeen)
Symon Stevens-Guille (Ohio State University)
Kristina Striegnitz (Union College)
Hsuan Su (National Taiwan University)
Hiroya Takamura (The National Institute of Advanced Industrial Science and Technology)
Ece Takmaz (University of Amsterdam)
Xiangru Tang (Yale University)
Marc Tanti (University of Malta)
Mariët Theune (University of Twente)
Ross Turner (Arria NLG)
Henrik Voigt (Friedrich-Schiller-University)
Di Wang (ContextLogic Inc)
Qingyun Wang (University of Illinois at Urbana-Champaign)
Robert Weißgraeber (CTO @ AX Semantics)
Michael White (Ohio State University)
Yuan-Kuei Wu (National Taiwan University)
Juncheng Xie (National Taiwan University)
Xinnuo Xu (University of Edinburgh)
Bohao Yang (University of Sheffield)
Ziheng Zeng (University of Illinois at Urbana-Champaign)
Zhirui Zhang (Tencent AI Lab)

Zaixiang Zheng (ByteDance AI Lab)
Yinhe Zheng (miHoYo)
Wanzheng Zhu (Google)

# Table of Contents

x

# Guided Beam Search to Improve Generalization in Low-Resource Data-to-Text Generation

**Nicolas Garneau**
Department of Computer Science
University of Copenhagen, Denmark
and Université Laval, Canada

**Luc Lamontagne**
Department of Computer Science
Université Laval

## Abstract

In this paper, we introduce a new beam search algorithm that improves the generalization of neural generators to unseen examples, especially in low-resource data-to-text settings. Our algorithm aims to reduce the number of omissions and hallucinations during the decoding process. For this purpose, it relies on two regression models to explicitly characterize factual errors. We explain how to create a new dataset to train these models given an original training set of less than a thousand data points. We apply our approach in the low-resource, legal setting using the French *Plum2Text* dataset, as well as in English using *WebNLG*. We observe in our experiment that this combination improves the faithfulness of pre-trained neural text generators using both human and automatic evaluation. Moreover, our approach offers a level of interpretability by predicting the number of omissions and hallucinations present in a given generation with respect to the input data. Finally, we visualize our algorithm's exploration of the hypothesis space at different steps during the decoding process.

## 1 Introduction

Data-to-text generation is commonly referred to as the task of verbalizing a structured input also known as a table of values. The table may contain several types of values such as text, numbers, categories, etc. In our study, we are specifically interested in improving the faithfulness of neural data-to-text generators. The relevance of their generations can be evaluated with respect to the coverage of the input table, i.e. to what extent the model omits values from the table. Moreover, neural text generators unfortunately have the tendency to hallucinate facts from the training set. Hence, generations can also be evaluated based on the number of hallucinated facts produced by the model (Dušek et al., 2018; Ji et al., 2022). The tendency of neural data-to-text generators to omit values and/or hallucinate facts can be exacerbated in low-resource

settings. The models may overfit the training set, and bring generalization to unseen data points at stake.

In this paper, we propose to improve the faithfulness of data-to-text neural generators by reducing the number of hallucinations and omissions during the generation process, without having to re-train the generation models. This perspective has many incentives, especially since models are becoming larger and larger, thus harder to train (Brown et al., 2020; Hoffmann et al., 2022). To this end, we propose a modified version of the beam search algorithm specifically for the data-to-text setting. Meister et al. (2020) studied the behavior of the beam search algorithm under a regularized framework, showing that beam search enforces uniform information density. That is, "*it produces text with evenly distributed surprisal, a feature that human readers tend to prefer*". Inspired by this regularization framework, we introduce two characterization models that will guide the decoding algorithm by promoting generated beams containing fewer hallucinations and omissions.

The characterization of omissions and hallucinations is crucial in the legal setting. Hence, we apply our new decoding algorithm and analyze its benefits on the task of verbalizing criminal docket files using the *Plum2Text* dataset (Beauchemin et al., 2020; Garneau et al., 2021b). Using automatic and manual evaluation, we show that our algorithm improves generalization in a low-resource setting, especially on unseen data points. We also show that our approach generalizes to other datasets, such as WebNLG (Castro Ferreira et al., 2020). In the next section, we introduce related work regarding the mitigation of omissions and hallucinations for neural text generators. We then introduce the main contribution of this paper in Section 3, a new decoding algorithm for the data-to-text setting. We present the experiments and analysis in Section 4. We assess the generalization of our approach in Sec-

tion 5 by applying it on WebNLG (Castro Ferreira et al., 2020). We conclude with our observations in Section 6. The models, datasets, generations and human evaluations are made publicly available[1].

## 2 Related Work

In this section, we study current mitigation techniques of omissions and hallucinations to improve neural generators' performance, limiting ourselves to the data-to-text setting[2]. These techniques *may* require changing the architecture and are enforced either during training or during inference. We classify mitigation techniques as being "invasive" or "non-invasive" to the generator. Invasive techniques require fine-tuning, adding a new objective function, or modifying the inner architecture of the generator.

In this paper, we focus on non-invasive techniques, which consider the generator as a black box and act either on the input or during the decoding process. These techniques are appealing for the fact that they do not require re-training the generator on the original training dataset. For example, Shin et al. (2020) proposed AutoPrompt, a model that learns how to create prompts for various sets of tasks. They basically search for "trigger" tokens using the gradient from the downstream task. According to their results, AutoPrompt outperformed fine-tuning methods in cases where the training dataset is small (i.e. 100-1000 samples). This method, however, does not necessarily mitigate the omission and hallucinations in the data-to-text setting. Similarly, Prefix-Tuning (Li and Liang, 2021) proposed a lightweight alternative to fine-tuning for natural language generation tasks, which keeps language model parameters frozen, but optimizes a small continuous task-specific vector, called the prefix. Then again, their method does not specifically mitigate omissions and hallucinations. Ghazvininejad et al. (2017) proposed *Hafez*, a method weighing the current beam state based on a set of feature functions that take as input a target word, and sometimes the beam state (e.g. to check for repetitions). These feature functions, in our case, could be used to force the generation of proper charges, decisions, and pleading, for example. They added two terms to the standard Beam

Search algorithm, given a current beam state and a predicted word;

$$score(b_i, w) = score(b_{i-1}) + \log Gen(w) +$$
$$\sum_j \alpha_j * f_j(w); \forall w \in V_{suc} \quad (1)$$

where $score(b)$ is the score of the current beam state, $\log Gen(w)$ is the output logit of the generator, $f(*)$ are functions that scores word $w$ weighted by $\alpha_i$, and $V_{suc}$ is a predefined vocabulary. Similarly, *Mention Flags* (Wang et al., 2021) tries to identify the presence of tokens in the hypothesis given a set of flags. Both methods face the same problem since they operate on surface tokens.

Anderson et al. (2017) also proposed to constrain the beam search algorithm operating at the lexical level using a finite-state machine that enforces the use of a specific vocabulary in the image captioning setting. However, their method does not scale well when the input is composed of sentences, since we don't know apriori the vocabulary we want to constrain. Balakrishnan et al. (2019) proposed a constrained decoding technique that leverages tree-structured meaning representations to control the semantic correctness of the generated text. While not explicitly characterizing omissions and hallucinations, their approach improved the faithfulness of the generative models. The prior work closest to ours is RANKGEN (Krishna et al., 2022), a ranking model that can be incorporated into the beam search scoring function during the decoding process. However, their method is designed for open-ended generation and does not yet scale to methods having a constrained output such as data-to-text, summarization, and machine translation.

Guerreiro et al. (2023) introduced DEHALLU-CINATOR, a model that flags hypotheses once they are fully generated so that they can be overwritten. Our model differs from their approach since we are guiding the exploration of the tree during decoding. Finally, (Vijayakumar et al., 2016) introduced Diverse Beam Search, an algorithm that promotes diverse generations amongst groups of beams but does not strictly reward or penalize beams for specific properties. To the best of our knowledge, no method in the literature proposes a way that can be adapted without major changes to handle both omissions and hallucinations at the semantic level during the decoding step. Moreover, none of the methods can explicitly estimate the number of hallucinations and omissions in the hypotheses. We

---

[1] https://drive.google.com/drive/folders/1gdgWXr7GXoww2dC83xcH3TXgwyB3ov4W?usp=sharing

[2] One can see the mitigation as the optimization of the precision (reduce hallucination) and the recall (reduce omissions) of a given data-to-text generator.

thus wish to fill this gap by proposing a guided beam search algorithm to create more faithful neural data-to-text generations.

## 3 Guided Decoding by Predicting Omissions and Hallucinations

In this section, we introduce a new decoding algorithm that is designed to mitigate and explicitly characterize omissions and hallucinations for data-to-text generation. To this end, we create two predictive models: one predicting the number of omitted values from the table, $m_o$, and the other predicting the number of hallucinations, $m_h$. These models will thus weigh the current beam score to promote generated sequences with few, or hopefully no omissions or hallucinations, enforcing semantically accurate generations.

### 3.1 Characterization Models

The proposed models are designed to take as input the table's values, as well as the current generated sequence, and output a real value as the following;

$$o_i = m_o(V_i, s_i) \quad (2)$$
$$h_i = m_h(V_i, s_i) \quad (3)$$

where $o_i$ is the predicted number of omissions, $h_i$ is the predicted number of hallucinations. $m_o$ is the omission model, $m_h$ is the hallucination model, $V_i$ is the set of table of values, and $s_i$ is the current generated sequence. To obtain these models, we need to train them using a dataset that has as input the table, the generated sequence as well as their true labels, i.e. the number of omissions and hallucinations in the sequence. We further detail in the next section how we obtain such datasets from the original training set using *Plum2Text* as an example (Plum2Text's training set contains around 1K examples).

### 3.2 Training Data

We hereby propose to build one training dataset for each model, $\mathcal{O}$ and $\mathcal{H}$, based on the overlapping table values across the original training examples. It is important to note that each actual training example is used in both $\mathcal{O}$ and $\mathcal{H}$, labeled with zero omission and zero hallucination respectively. We create the other training examples as follows;

1. We randomly select two training instances $(V_i, r_i)$, and $(V_j, r_j)$ where $r_i$ and $r_j$ are reference texts of both examples

2. The set of omitted values $O_i$ for $r_j$ with respect to $V_i$ correspond to the set difference between $V_i$ and $V_j$

3. Similarly, the set of hallucinated values $H_i$ for $r_i$ with respect to $V_j$ correspond to the set difference between $V_j$ and $V_i$.

We formally describe the dataset creation in Algorithm 1 and we illustrate in Figure 1 the construction of a training example, created from two original examples taken from the *Plum2Text* dataset.

---

**Algorithm 1** Creating Datasets $\mathcal{O}$ and $\mathcal{H}$

$\mathcal{O} \leftarrow \{\}$                ▷ set of omissions
$\mathcal{H} \leftarrow \{\}$             ▷ set of hallucinations
**for** $(V_i, r_i), (V_j, r_j)$ in the training set **do**
    $O_i \leftarrow V_i \setminus V_j$     ▷ set diff. between $V_i$ and $V_j$
    $H_i \leftarrow V_j \setminus V_i$     ▷ set diff. between $V_j$ and $V_i$
    $\mathcal{O} \leftarrow \mathcal{O} \cup \{(V_i, r_j), |O_i|\}$
    $\mathcal{H} \leftarrow \mathcal{H} \cup \{(V_i, r_j), |H_i|\}$
**end for**
**return** $\mathcal{O}, \mathcal{H}$

---

Using *Plum2Text*, the omissions dataset $\mathcal{O}$ consists of 12,460 examples using an 80%–20% split resulting in train and test sets of 9,968 and 2,492 examples respectively. The hallucination dataset $\mathcal{H}$ consists of 30,473 examples also using an 80%–20% split resulting in train and test sets of 24,378 and 6,095 examples respectively. With respect to the training architecture, we used the multilingual version of BERT (Devlin et al., 2019) of 178M parameters available in the HuggingFace library[3]. We used the mean squared error loss and AdamW (Loshchilov and Hutter, 2019) as the optimizer with a learning rate of 0.001. We used a batch size of 10 on a GeForce 2080Ti Nvidia graphic card. To automatically evaluate the architectures, we considered several metrics: mean squared error (MSE), root mean square error (RMSE), mean average error, $\mathcal{R}^2$, and accuracy defined as follows;

$$a = \begin{cases} 1 \text{ if } p - t < 0.5 \\ 0 \text{ otherwise} \end{cases} \quad (4)$$

where $p$ is the prediction and $t$ is the true value. As we can see in Table 1, both models achieve high performance across all metrics on the test set.

---

[3]We used the multilingual BERT (Devlin et al., 2019) because it provides a version with a pre-trained classification head, whereas CamemBERT (Martin et al., 2020) does not.

**Table 1**

| Accusation: Provision 320.14 (1) a) |
|---|
| Every person commits an offence who :<br>(a) operates a conveyance while his or her ability to drive is impaired to any degree by the effect of alcohol or a drug or by the combined effect of alcohol and a drug; |
| Plea |
| Pleaded not guilty |
| Decision |
| Declared guilty |

**Table 2**

| Accusation: Provision 265 (1) a) |
|---|
| A person commits an assault when :<br>(a) without the consent of another person, he applies force intentionally to that other person, directly or indirectly; |
| Plea |
| – |
| Decision |
| Declared guilty |

**Reference 1**

*"PER pleaded not guilty on a count of impaired driving and was declared guilty."*

**Omitted:**
1. Provision 265 (1) a)
**Hallucinated:**
1. Provision 320.14 (1) a)
2. Guilty plea

**Reference 2**

*"PER is accused on a count of assaulting another person by applying force intentionally and was declared guilty."*

**Omitted:**
1. Provision 320.14 (1) a)
2. Guilty plea
**Hallucinated:**
1. Provision 265 (1) a)

Figure 1: Given two training instances from *Plum2Text*, each with their respective table and reference, we pair the table from the first example with the reference of the second one and vice versa. This creates in total four training instances, two in each dataset $\mathcal{O}$ and $\mathcal{H}$. The first two "omitted" training instances are $V_j$, $r_i$ paired with the omitted value "Provision 265 (1) a)", and $V_i$, $r_j$ paired with the 2 omitted values "Provision 320.14 (1) a)" and "Guilty plea". The same procedure applies for the creation of the hallucinated training instances.

We also show the distribution of predicted vs actual values in Figure 2 using confusion matrices. The regression model on the omissions tends to underestimate the number of omissions in a given generation. The regression model on the hallucinations seems more balanced except for the cases where there are one or two hallucinations, underestimating them.

| | **Models** | |
|---|---|---|
| **Metric** | **Omission** | **Hallucination** |
| **MSE** | 0.05 | 0.05 |
| **RMSE** | 0.23 | 0.22 |
| **MAE** | 0.10 | 0.08 |
| $\mathcal{R}^2$ | 0.99 | 0.99 |
| **Accuracy** | 0.96 | 0.97 |

Table 1: Performance of both omission and hallucination models on *Plum2Text* w.r.t the mean squared error (MSE), the root mean squared error (RMSE), the mean average error (MAE), $\mathcal{R}^2$, and accuracy.

### 3.3 Guided Decoding for Omission and Hallucination Mitigation

In order to mitigate omissions and hallucinations, we propose the following weighted beam search score $b_i$;

$$b_i = score(b_{i-1}) + \log(Gen(w_i)) + \phi_i \quad (5)$$

where $score(b_{i-1})$ is the previous beam's score, $\log(Gen(w_i))$ is the score for word $w_i$ provided by the generator, and $\phi$ is the following function based on the omission and hallucination scores $o_i$ and $h_i$ obtained from the characterization models:

$$\phi_i = \omega \cdot (v_i - o_i) - \gamma \cdot h_i \quad (6)$$

where $\omega$ and $\gamma$ are parameters to weigh the omissions and hallucinations respectively. While the hallucinations $h_i$ are treated as a penalty on a beam score, the omissions are treated as a reward: $v_i$ corresponds to the actual number of values in the table, whereas $o_i$ is the number of detected omissions. If zero omissions are detected, the current beam will get a reward of $\omega \cdot v_i$.

In our experiments, $\omega$ and $\gamma$ are initialized to 1 and we perform a grid search over a set of values between 0.0 and 5.0 to find the optimal ones depending on the use case. The number of beams $b$

(a) Predicted omissions.



(b) Predicted hallucinations.

Figure 2: Confusion matrices of the predicted omissions and hallucinations by the regression models on the *Plum2Text* test set.

parameterizes the original beam search algorithm. During the generation process, omission rewards and hallucination penalties are cumulated at each step. Regardless of the values of $\omega$ and $\gamma$, we apply a final processing step to fully reevaluate the ranking of the candidates w.r.t the generator's final log-likelihood and the omission/hallucination models using values of 1 for both $\omega$ and $\gamma$. This is motivated by the fact that the models, trained on full sentences, may provide more accurate predictions and thus result in a better candidate ranking.

## 4 Experiments

In our experiment, we use *CriminelBART*, a generative model introduced by Garneau et al. (2021a). We only analyze the vanilla and guided versions of *CriminelBART* since other methods proposed in the literature do not explicitly mitigate omissions and hallucinations. We trained *CriminelBART* on the train set of *Plum2Text*, and we begin by automatically evaluating different versions of the weighted beam search using a grid search over the hyperparameters previously introduced. We then manually evaluate the performance of our new algorithm in Section 4.2. To assess the generalization performance of our algorithm, we added examples with 37 new provisions from the Criminal Code of Canada having no or very few occurrences in the original training set. Furthermore, we qualitatively analyze the behavior of our algorithm in Section 4.3.

### 4.1 Guided Decoding

In order to find the best generation model using the weighted decoding algorithm aforementioned, we performed a grid search exploration with the following hyper-parameters:

1. $\omega$, the weight for omission detection.

2. $\gamma$, the weight for hallucination detection.

3. $\beta$, the number of beams.

| Parameters | Values |
|---|---|
| Omission $-\omega$ | 0.0, 0.1, 0.2, 0.5, 1.0, 2.0, 5.0 |
| Hallucination $-\gamma$ | 0.0, 0.1, 0.2, 0.5, 1.0, 2.0, 5.0 |
| Beam size $-\beta$ | 5, 10, 15 |

Table 2: Hyper-parameters search on the omission and hallucination weight ($\omega$ and $\gamma$ respectively) and the beam size $\beta$.

Table 2 provides the values tried for each hyper-parameter. Among the 147 combinations, the best model uses weights of $0.2$ for both omissions and hallucinations and a beam size of $15$.

Evaluation results are presented in Table 3 for both the best-performing model using guided decoding and the original version of *Criminel-BART*. We considered BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), METEOR (Banerjee and Lavie, 2005), BertScore (Zhang et al., 2020) dubbed as BScore, NLI (Dušek and Kasner, 2020),

5

and RANK (Garneau and Lamontagne, 2021). It has been shown that RANK highly correlates with human judgment (CITE), so we used this metric in the cross-validation step to select the best hyper-parameters, $\omega$, $\gamma$, and $\beta$, for each model. The results indicate that guided *CriminelBART* outperforms the original *CriminelBART* on 6 automatic evaluation metrics out of 9. The guided version of *CriminelBART* obtains similar performance with respect to BLEU-1, METEOR, and NLI.

The guided version of *CriminelBART* using the post-processing step described in the previous section obtains similar performance but we observe an interesting two-point gain on the RANK metric, improving from 0.76 to 0.78, over the original version limited to 0.72. It is important to note that RANK tends to have the highest correlation score with respect to human evaluation. Overall, we can conclude that the guided version of *CriminelBART* obtains better performance than the original one by up to 6 points with respect to the RANK metric. We also note that the number of predicted hallucinations and omissions also considerably decrease, going from 0.28 and 0.24 to 0.11 and 0.11 respectively. In the next section, we manually evaluate the generations.

## 4.2 Human Evaluation

In this section, we further analyze the generalization performance of both models by considering 45 table values that are either not in the training set or appear rarely. We hired three annotators that followed the same evaluation procedure introduced by Garneau et al. (2022) to manually assess the performance of both models. For our application, these table values correspond to legal provisions from the Criminal Code of Canada (CCC). From these 45 provisions, we added 37 new ones that we selected by skimming through the whole CCC. These are listed in Appendix A. We list down in Table 4 the whole set of provisions considered in this manual evaluation. We decided to not manually evaluate examples where other provisions were found often in the training set because both models are having a similar performance for these frequent cases.

We manually evaluated the generations of both the original version of *CriminelBART* and the model using guided beam search. We recruited three evaluators from a Faculty of Law that assigned a score between 1 to 10, 1 corresponding

to a generation completely off-track, and 10 being a perfect generation. We used Krippendorff's alpha coefficient (Krippendorff, 2004) to analyze the inter-annotator agreement which is 0.69. We can see from Table 4 that the guided version of *CriminelBART* achieves better generalization performance on unseen provisions with an average score of 7.4, compared to the original version with a score of 3.9. That is, guided *CriminelBART* produces generations that verbalize the good provision with some hallucinations and/or omissions, whereas the original version mostly generates on-theme or off-track descriptions. It seems like the hallucination and omission models enable better exploration of the generation tree than regular beam search using maximum log-likelihood estimation. This can lead to better generations when using a higher number of beams (Meister et al., 2020). We specifically discuss and illustrate this phenomenon in Section 4.3.

Comparing *CriminelBART* and Guided *CriminelBART*, we found that for 10 out of 45 generations, the original version of *CriminelBART* generated commonly seen provisions such as 320.14 (driving under the influence), 266 (assault and battery), or 151 (sexual interference). We provide an example in Table 5 where the guided *CriminelBART* generated the good provision, but the original version generated unrelated content with respect to the input. There is one particular case where the original version produced a better generation which is on provision 345, "Stopping mail with intent" (see Table 4). Indeed, the guided version of *CriminelBART* produced a generation not capturing the act of stealing **mail**, while the original version did. In every other case, the original version attempted at generating meaningful content as being "on-theme", but most of the time the guided version was able to generate the right provision, with some factual errors, having a score above 5.

## 4.3 Beam Search Analysis

To better understand the behavior of our approach, we analyze the different steps in the beam search algorithm of both models for one generation involving table value "provision 431": *Attack on-premises, residence or transport of internationally protected person.* We illustrate in Figure 3 the paths taken by the two versions of beam search. The starting point, where the algorithms respectively branched on their own, is illustrated in blue.

6

| | | | BLEU | | | | ROUGE | METEOR | BScore | NLI | RANK | Rates | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\omega$ | $\gamma$ | $\beta$ | 1 | 2 | 3 | 4 | | | | | | Hal. | Om. |
| 0.0 | 0.0 | 5 | 0.73 | 0.58 | 0.47 | 0.41 | 0.42 | 0.38 | 0.78 | 0.34 | 0.72 | 0.28 | 0.24 |
| 0.2 | 0.2 | 15 | 0.73 | **0.59** | **0.48** | **0.43** | **0.44** | 0.38 | **0.79** | 0.34 | **0.76** | 0.13 | 0.11 |
| Post processing | | | 0.73 | **0.58** | **0.48** | **0.42** | **0.43** | 0.37 | **0.79** | 0.34 | **0.78** | **0.11** | **0.11** |

Table 3: Automatic evaluation results of the best performing original *CriminelBART* ($\omega = 0.0$, $\gamma = 0.0$, $\beta = 5$), the best-performing model using the weighted beam search algorithm ($\omega = 0.2$, $\gamma = 0.2$, $\beta = 15$), and that same model using the post-processing finalization step.



Figure 3: Analysis of *CriminelBART* using the original (orange) and the guided (green) beam search algorithms on the generation of provision 431 (translated in English): "*Attack on-premises, residence or transport of internationally protected person.*". The predicted number of omissions ($o_i$) and hallucinations ($h_i$) are presented at each timestep.

The original beam search algorithm is illustrated in yellow, while the guided beam search is illustrated in green. Each time step in the figures may be an aggregation of several generation steps, for easier understanding. Each time step is associated with the predicted number of omissions and hallucinations, $o_i$ and $h_i$ respectively, with respect to their associated models.

In this particular example, the models start with omitted and hallucinated values of one regarding the initial generation "*the or on the DATE, at LOC, LOC, exercised*"[4]. The decoding algorithm branch out on the next token, generating "assault" for the original version and "an" for the guided one. It is only with two generation steps that the guided beam search obtains lower predicted values in terms of omissions and hallucinations (0.64 and 0.28 respectively). The final generation obtains

scores of omissions and hallucinations of 0.00 and 0.04. The original version of the beam search on the other hand clearly omits to generate the proper provision, and hallucinates the provision of "assault", ending with both omission and hallucinations scores of 1.0. Finally, the original version of *CriminelBART* obtains a human evaluation score of 1.0, compared to the guided version having 8.33.

## 5 Generalization of the Approach

We illustrate the generalization of our proposed approach to improve the beam search algorithm to other data-to-text settings by using WebNLG, a well-known dataset in the NLP community. To this end, we used the same methodology described in Section 3.3:

- From the original WebNLG training dataset, we build two datasets, $\mathcal{O}$ and $\mathcal{H}$. Training instances of these datasets consist of a set of triplets each containing a table, a generation,

---

[4]DATE, PER and LOC are special tokens from the *Plum2Text¨* dataset where dates, persons, and locations have been anonymized.

| Provision | *CriminelBART* | Guided *CriminelBART* |
|---|---|---|
| 46 | 1.00 | 8.00 |
| 57 | 3.00 | 8.00 |
| 58 | 2.33 | 7.00 |
| 83.04 | 2.67 | 8.00 |
| 83.08 | 3.00 | 8.00 |
| 83.21 | 5.33 | 8.00 |
| 83.181 | 1.00 | 8.00 |
| 123 | 1.00 | 8.00 |
| 148 | 7.67 | 8.67 |
| 150 | 3.67 | 8.33 |
| 170 | 2.33 | 5.00 |
| 173 | 2.33 | 8.33 |
| 202 | 1.00 | 4.67 |
| 218 | 1.00 | 5.67 |
| 243 | 4.33 | 6.67 |
| 245 | 2.00 | 7.33 |
| 253 | 6.00 | 8.00 |
| 267 | 6.33 | 8.00 |
| 270.1 | 3.33 | 8.67 |
| 318 | 7.00 | 8.33 |
| 342 | 8.67 | 9.00 |
| 342.1 | 2.33 | 9.67 |
| 344 | 4.00 | 8.67 |
| 345 | 7.67 | 1.00 |
| 347 | 1.00 | 6.00 |
| 351 | 7.00 | 9.00 |
| 354 | 3.00 | 8.00 |
| 355 | 5.00 | 7.67 |
| 356 | 1.00 | 7.67 |
| 364 | 1.00 | 8.67 |
| 368 | 7.33 | 9.00 |
| 374 | 4.67 | 5.00 |
| 382.1 | 8.33 | 4.00 |
| 398 | 8.00 | 6.00 |
| 402.2 | 8.00 | 8.33 |
| 406 | 3.33 | 8.00 |
| 431 | 1.00 | 8.33 |
| 432 | 5.00 | 4.33 |
| 437 | 1.00 | 4.33 |
| 438 | 5.67 | 8.33 |
| 439 | 2.33 | 8.33 |
| 445.1 | 3.00 | 9.00 |
| 446 | 2.33 | 8.67 |
| 467.111 | 8.33 | 8.67 |
| 810.2 | 2.33 | 5.67 |
| **Average** | **3.9** | **7.4** |

Table 4: Human evaluation of the original version of *CriminelBART* and the one using guided beam search on the 45 unseen provisions.

and the associated number of omissions or hallucinations.

- Using the previously created datasets, we train two models to predict the number of omissions and hallucinations given the input table and its corresponding generation.

- We use the trained models to predict, during the decoding process, the number of omissions and hallucinations and weigh the beams accordingly.

- We apply the finalization step to select the best hypothesis.

The omission dataset $\mathcal{O}$ of WebNLG consists of 20,448 examples resulting in train and test sets of 16,358 and 4,090 examples respectively using an 80%–20% split. The hallucination dataset $\mathcal{H}$ consists of 20,600 examples resulting in train and test sets of 16,480 and 4,120 examples respectively also using an 80%–20% split. Similar to the *Plum2Text* setting, we trained the English version of BERT (Devlin et al., 2019) available in the HuggingFace library using the same hyperparameters previously selected. For the actual data-to-text generation task we trained BART (Lewis et al., 2020) on the training set of WebNLG and evaluated it on the test set after performing a grid-search over the guided beam search hyperparameters. We used the same automatic evaluation metrics as with *Plum2Text*. We can see from Table 6 that the guided version of BART on WebNLG improves the performance on 6 metrics out of 9. Considerable improvements are made regarding NLI and RANK, two metrics that were proven to be correlated with manual evaluation of WebNLG test instances (Garneau et al., 2022). Similar to the *Plum2Text* case, using the post-processing step to re-rank candidates improved the generations for almost all metrics. While the hallucination and omission rates are already low, using the guided version of BART slightly improves them.

## 6  Conclusion

In this paper, we introduced a new guiding mechanism for the beam search algorithm in the data-to-text generation setting. We presented how to train two models, one to predict the number of omissions, and the other for the number of hallucinations. These predictors are used not only to mitigate the number of omissions and hallucinations

| Input Data | **Provision:** 123: Influencing a municipal officer; is guilty of an indictable offence and liable to imprisonment for a term not exceeding five years or is guilty of an offence punishable on summary conviction if he influences or attempts to influence a municipal officer to do any of the things referred to in paragraphs (1)(a) to (d): by threats or deception. |
|---|---|
| *Criminel*BART<br>Score: 1.0 | On the LABELD, PER is arrested for driving with more than 80 mg of alcohol per 100 ml of blood. |
| **Guided** *Criminel*BART<br>Score: 8.0 | On or about LABELD, at LOC, LOC, willfully attempted to influence a municipal official, thereby committing the criminal act under section 120(a) of the Criminal Code. |

Table 5: Comparing the generation of the original *CriminelBART* and the model using guided beam search on provision 123: "Influencing a municipal officer". The original version of *CriminelBART* generates one of the most common provisions, 320.14 (driving under the influence) resulting in a score of 1.0, whereas the model using guided beam search generates a description about the right provision resulting in a score of 8.0.

| | | | BLEU | | | | | | | | | Rates | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\omega$ | $\gamma$ | $\beta$ | 1 | 2 | 3 | 4 | ROUGE | METEOR | BScore | NLI | RANK | Hal. | Om. |
| 0.0 | 0.0 | 5 | 0.81 | 0.71 | 0.64 | 0.58 | 0.55 | 0.54 | 0.94 | 0.63 | 0.64 | 0.11 | 0.00 |
| 0.2 | 0.2 | 15 | **0.83** | **0.73** | **0.65** | **0.59** | 0.53 | 0.54 | 0.94 | **0.68** | **0.65** | 0.10 | 0.00 |
| Post processing | | | **0.84** | **0.74** | **0.66** | **0.60** | 0.54 | 0.54 | 0.94 | **0.68** | **0.66** | 0.10 | 0.00 |

Table 6: Automatic evaluation results of the best performing BART model on WebNLG ($\omega = 0.0$, $\gamma = 0.0$, $\beta = 5$) and the best-performing BART model using the weighted beam search algorithm ($\omega = 0.2$, $\gamma = 0.5$, $\beta = 10$).

but also to favor the exploration of the possible generation space. This new mechanism improves the generation quality with respect to automatic evaluation metrics and shows significant generalization improvement regarding unseen data points during human evaluation. Moreover, our mechanism offers a new degree of a posteriori interpretability given a list of potential hypotheses, since the characterization models provide estimates of the number of omissions and hallucinations. Finally, we showed that our approach generalizes not only to *Plum2Text*, a challenging low-resource dataset but also to a well-known dataset such as WebNLG. In future works, it would be interesting to investigate the identification of omitted values and hallucinated tokens. The identification of omitted values is easier to perform since we already provide a way to build such a dataset and train a model accordingly. However, identifying the hallucinated tokens requires a sequence-to-sequence tagger and its respective training set, which most likely can only be obtained with manual annotations.

## Ethics Statement

The scope of this work is to improve the faithfulness of neural data-to-text generators. Faithfulness is extremely important in the legal field since we do not want to generate false accusations about litigants. There is a potential risk to using neural data-to-text generators in production, and we provided not only improve their performance but also analyzed their behavior. In the end, the purpose of this work is largely motivated by the ethical use of neural text generators and a better understanding of their implications.

## Acknowledgements

## References

Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. 2017. Guided open vocabulary image captioning with constrained beam search. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 936–

945, Copenhagen, Denmark. Association for Computational Linguistics.

Anusha Balakrishnan, Jinfeng Rao, Kartikeya Upasani, Michael White, and Rajen Subba. 2019. Constrained decoding for neural NLG from compositional representations in task-oriented dialogue. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 831–844, Florence, Italy. Association for Computational Linguistics.

Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.

David Beauchemin, Nicolas Garneau, Eve Gaumond, Pierre-Luc Déziel, Richard Khoury, and Luc Lamontagne. 2020. Generating intelligible plumitifs descriptions: Use case application with ethical considerations. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 15–21, Dublin, Ireland. Association for Computational Linguistics.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Thiago Castro Ferreira, Claire Gardent, Nikolai Ilinykh, Chris van der Lee, Simon Mille, Diego Moussallem, and Anastasia Shimorina, editors. 2020. *Proceedings of the 3rd International Workshop on Natural Language Generation from the Semantic Web (WebNLG+)*. Association for Computational Linguistics, Dublin, Ireland (Virtual).

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Ondřej Dušek and Zdeněk Kasner. 2020. Evaluating semantic accuracy of data-to-text generation with natural language inference. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 131–137, Dublin, Ireland. Association for Computational Linguistics.

Ondřej Dušek, Jekaterina Novikova, and Verena Rieser. 2018. Findings of the E2E NLG challenge. In *Proceedings of the 11th International Conference on Natural Language Generation*, pages 322–328, Tilburg University, The Netherlands. Association for Computational Linguistics.

Nicolas Garneau, Eve Gaumond, Luc Lamontagne, and Pierre-Luc Déziel. 2021a. Criminelbart: A french canadian legal language model specialized in criminal law. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Law*, ICAIL '21, page 256–257, New York, NY, USA. Association for Computing Machinery.

Nicolas Garneau, Eve Gaumond, Luc Lamontagne, and Pierre-Luc Déziel. 2021b. Plum2text: A french plumitifs–descriptions data-to-text dataset for natural language generation. In *Proceedings of the 18th International Conference on Artificial Intelligence and Law*, Sao Paulo, Brazil. International Association for Artificial Intelligence and Law.

Nicolas Garneau, Eve Gaumond, Luc Lamontagne, and Pierre-Luc Déziel. 2022. Evaluating legal accuracy of neural generators on the generation of criminal court dockets description. In *Proceedings of the 15th International Conference on Natural Language Generation*, pages 73–99, Waterville, Maine, USA and virtual meeting. Association for Computational Linguistics.

Nicolas Garneau and Luc Lamontagne. 2021. Trainable ranking models to evaluate the semantic accuracy of data-to-text neural generator. In *Proceedings of the 2nd Workshop on Evaluation and Comparison of NLP Systems*, pages 51–61, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Marjan Ghazvininejad, Xing Shi, Jay Priyadarshi, and Kevin Knight. 2017. Hafez: an interactive poetry generation system. In *Proceedings of ACL 2017, System Demonstrations*, pages 43–48, Vancouver, Canada. Association for Computational Linguistics.

Nuno M. Guerreiro, Elena Voita, and André Martins. 2023. Looking for a needle in a haystack: A comprehensive study of hallucinations in neural machine translation. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1059–1075, Dubrovnik, Croatia. Association for Computational Linguistics.

Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. 2022. Training compute-optimal large language models.

Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2022. Survey of hallucination in natural language generation. *ACM Comput. Surv.* Just Accepted.

Klaus Krippendorff. 2004. *Content Analysis: An Introduction to Its Methodology (second edition)*. Sage Publications.

Kalpesh Krishna, Yapei Chang, John Wieting, and Mohit Iyyer. 2022. Rankgen: Improving text generation with large ranking models. *ArXiv*, abs/2205.09726.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4582–4597, Online. Association for Computational Linguistics.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In *International Conference on Learning Representations*.

Louis Martin, Benjamin Muller, Pedro Javier Ortiz Suárez, Yoann Dupont, Laurent Romary, Éric de la Clergerie, Djamé Seddah, and Benoît Sagot. 2020. CamemBERT: a tasty French language model. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7203–7219, Online. Association for Computational Linguistics.

Clara Meister, Ryan Cotterell, and Tim Vieira. 2020. If beam search is the answer, what was the question? In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2173–2185, Online. Association for Computational Linguistics.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. 2020. AutoPrompt: Eliciting Knowledge from Language Models with Automatically Generated Prompts. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4222–4235, Online. Association for Computational Linguistics.

Ashwin K Vijayakumar, Michael Cogswell, Ramprasath R Selvaraju, Qing Sun, Stefan Lee, David Crandall, and Dhruv Batra. 2016. Diverse beam search: Decoding diverse solutions from neural sequence models. *arXiv preprint arXiv:1610.02424*.

Yufei Wang, Ian Wood, Stephen Wan, Mark Dras, and Mark Johnson. 2021. Mention flags (MF): Constraining transformer-based text generators. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 103–113, Online. Association for Computational Linguistics.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

## A New Unseen Provisions

List of the 37 new unseen provisions and their associated texts.

- **46 (1) a)**: "*High treason. Every person commits high treason who, in Canada, wages war against Canada or does any act preparatory thereto;*"

- **57 (2)**: "*Misrepresentation in relation to a passport. Every person who, in Canada or elsewhere, for the purpose of obtaining a passport for himself or herself or for another person or for the purpose of obtaining a material alteration or addition to such a passport, makes a written or oral statement that he or she knows to be false or misleading is guilty.*"

- **58 (1) a)**: "*Fraudulent use of citizenship certificate. Every person who, while in Canada or outside Canada, as the case may be, uses a certificate of citizenship or a certificate of naturalization for a fraudulent purpose is guilty of an indictable offence and liable to imprisonment for a term not exceeding two years or is guilty of an offence punishable on summary conviction;*"

- **83.04 a)**: "*Using or possessing property for terrorist purposes. Any person who: directly or indirectly uses property, in whole or in part, for or to facilitate a terrorist activity is guilty of an indictable offense punishable by imprisonment for not more than ten years;*"

- **83.08 (1) a)**: "*Freezing of property. No person in Canada and no Canadian outside Canada shall: knowingly deal with property owned or controlled, directly or indirectly, by a terrorist group;*"

- **83.21 (1)**: "*Instructing a person to carry out an activity for a terrorist group. Every person who knowingly directs, directly or indirectly, any person to carry out any activity for the benefit of, at the direction of, or in association with a terrorist group for the purpose of enhancing the ability of any terrorist group to facilitate or carry out a terrorist activity is guilty of an indictable offence and liable to imprisonment for life.*"

- **123 (2)**: "*Influencing a municipal officer. Every person who influences or attempts to influence a municipal officer to do anything mentioned in paragraphs (1)(a) to (d) is guilty of an indictable offence and liable to imprisonment for a term not exceeding five years or is guilty of an offence punishable on summary conviction:*"

- **148 a)**: "*Assisting prisoner of war to escape. Every one who knowingly: aids a prisoner of war in Canada to escape from a place of confinement is guilty of an indictable offence and liable to imprisonment for a term not exceeding five years or is guilty of an offence punishable on summary conviction;*"

- **170**: "*Father, mother or guardian who procures. A parent or guardian of a person under the age of eighteen years who causes that person to engage in sexual acts prohibited by this Act with a third party is guilty of an indictable offence and liable to imprisonment for a term not exceeding fourteen years and to a minimum punishment of one year.*"

- **173 (2)**: "*Exhibitionism. Any person who, in any place whatsoever, for sexual purposes, exhibits his or her genitals in front of a person under the age of sixteen years is guilty of:*"

- **202 (1) a)**: "*Gambles, bookmaking, etc. Every person commits an offence who: uses or knowingly permits to be used any premises under his control for the purpose of registering or recording bets or selling a pool bet;*"

- **218**: "*Abandonment of child. Whoever unlawfully abandons or exposes a child under the age of ten years, so that the life of such child is actually endangered or exposed to be endangered, or the health of such child is actually permanently endangered or exposed to be endangered, is guilty of:*"

- **243**: "*Suppression of part. Whoever in any way causes the corpse of a child to disappear with the intention of concealing the fact that its mother gave birth to it, whether the child died before, during or after birth, is guilty:*"

- **245 (1)**: "*Administering deleterious substance. Whoever administers or causes to be administered to any person any poison or other destructive or deleterious substance, shall be guilty of:*"

- **270.1 (1)**: "*Disarming a peace officer. Every person commits an offence who takes or attempts to take a weapon from the possession of a peace officer acting in the performance of his or her duties, without the consent of the peace officer.*"

- **318 (1)**: "*Advocacy of genocide. Anyone who advocates or foments genocide is guilty of an indictable offence and liable to imprisonment for a term not exceeding five years.*"

- **342 (3)**: "*Unauthorized use of credit card data. Any person who fraudulently and without the appearance of right has in his possession or uses data, whether genuine or not, relating to a credit card, including a personal authenticator, which would enable the use of the same or the obtaining of services connected with its use, traffics in such data or allows another person to use the same, shall be guilty:*"

- **342.1 (1) a)**: "*Unauthorized use of computer. Every person who fraudulently and without colour of right, directly or indirectly, obtains computer services is guilty of an indictable offence and liable to imprisonment for a term not exceeding ten years or is guilty of an offence punishable on summary conviction;*"

- **345**: "*Stopping the mail with intent to rob. Anyone who stops a mail transport with the intention of stealing or searching it is guilty of a criminal act and liable to life imprisonment.*"

- **347 (1)**: "*Criminal rate of interest. Notwithstanding any other federal law, any person who enters into an agreement or arrangement to charge interest at a criminal rate or charges interest, even partially, at a criminal rate is guilty of:*"

- **351 (1)**: "*Possession of burglary tools. Whoever, without lawful excuse, has in his possession any instrument which may be used to break into any place, motor vehicle, vault or safe, knowing that the instrument has been used or is intended to be used for such purpose, is guilty of:*"

- **354 (2)**: "*Possession of motor vehicle with identification number obliterated. In proceedings under subsection (1), evidence that a person is in possession of a motor vehicle, or any part thereof, the identification number of which has been wholly or partly removed or obliterated is, in the absence of any evidence to the contrary, proof that it was obtained by the commission in Canada of an offence punishable on indictment;*"

- **356 (1) a)**: "*Theft of mail. Every person commits an offence who: steals anything sent by mail after it has been deposited in a post office and before it is delivered, or after it has been delivered but before it is in the possession of the addressee or any person who may reasonably be regarded as authorized by the addressee to receive the mail;*"

- **364 (1)**: "*Fraudulent obtaining of food and lodging. Any person who fraudulently obtains food, drink, or other commodities in any establishment dealing in them is guilty of a summary conviction offense.*"

- **368 (1) a)**: "*Using, possessing or trafficking in a forged document. Every person commits an offence who, knowing or believing that a document is counterfeit, as the case may be: uses, treats or acts with respect to it as if it were genuine;*"

- **374 (a)**: "*Unauthorized drafting of document. Any person who, with intent to defraud and without lawful authority, makes, subscribes, draws, signs, accepts or endorses a document in the name of or on behalf of another person, by proxy or otherwise, is guilty of an indictable offence and liable to imprisonment for a term not exceeding fourteen years;*"

- **382.1 (1) a)**: "*Insider trading. Every person who knowingly sells or buys securities, even indirectly, using confidential information that he or she holds as a shareholder of the issuer of the securities in question is guilty of an indictable offence and liable to imprisonment for a term not exceeding ten years or is guilty of an offence punishable on summary conviction;*"

- **398**: "*Falsifying record of employment. Every person who, with intent to mislead, falsifies a record of employment by any means, including*"

*the punching of a time clock, is guilty of a summary conviction offence."*

- **402.2**: *"Identity theft. Every person commits an offense who obtains or has in his or her possession identifying information about another person with the intent to use that information to commit an indictable offence, one of the elements of which is fraud, deceit or falsehood."*

- **406 a)**: *"Infringement of Trade-mark. For the purposes of this Part, a person who, without the consent of the owner of the trade-mark, makes or reproduces in any manner that trademark or a mark so nearly resembling it as to be likely to mislead;"*

- **431**: *"Attack on the official premises, private dwelling or means of transport of an internationally protected person. Any person who makes an attack accompanied by violence on the official premises, private dwelling or means of transportation of an internationally protected person in such a manner as to be likely to endanger the life or liberty of that person shall be guilty of an indictable offence punishable by imprisonment for a term not exceeding fourteen years."*

- **432 (1)**: *"Unauthorized recording of a motion picture. Whoever, without the consent of the manager of the cinema, records a cinematographic work - as that term is defined in section 2 of the Copyright Act - that is shown in a cinema, or its soundtrack, is guilty of:"*

- **437**: *"False alarm. Any person who willfully, without reasonable cause, by shouting, ringing bells, using a fire alarm, telephone or telegraph, or in any other manner, sounds or spreads or causes to be sounded or spread a fire alarm, is guilty."*

- **438 (2)**: *"Obstructing salvage of wreck. Every person who wilfully prevents or hinders, or wilfully seeks to prevent or hinder, the salvage of a wreck is guilty of an offence punishable on summary conviction."*

- **439**: *"Disturbance of marine signals. Every person who moors a ship or boat to a signal, buoy or other landmark used for navigation is guilty of an offence punishable on summary conviction."*

- **467.111**: *"Recruitment of members by criminal organization. Whoever recruits a person to be a member of a criminal organization-or invites, encourages, coerces, or solicits a person to be a member of a criminal organization-for the purpose of increasing the ability of the organization to facilitate or commit a criminal act under this or any other federal law is guilty of an indictable offense and liable:"*

# XF2T: Cross-lingual Fact-to-Text Generation for Low-Resource Languages

**Shivprasad Sagare, Tushar Abhishek, Bhavyajeet Singh, Anubhav Sharma**
**Manish Gupta and Vasudeva Varma**
IIIT-Hyderabad, India
{shivprasad.sagare,tushar.abhishek}@research.iiit.ac.in
{bhavyajeet.singh,anubhav.sharma}@research.iiit.ac.in
{manish.gupta,vv}@iiit.ac.in

## Abstract

Multiple business scenarios require an automated generation of descriptive human-readable text from structured input data. This has resulted into substantial work on fact-to-text generation systems recently. Unfortunately, previous work on fact-to-text (F2T) generation has focused primarily on *English* mainly due to the high availability of relevant datasets. Only recently, the problem of cross-lingual fact-to-text (XF2T) was proposed for generation across multiple languages along-with a dataset, XALIGN for eight languages. However, there has been no rigorous work on the actual XF2T generation problem. We extend XALIGN dataset with annotated data for four more languages: Punjabi, Malayalam, Assamese and Oriya. We conduct an extensive study using popular Transformer-based text generation models on our extended multilingual dataset, which we call XALIGNV2. Further, we investigate the performance of different text generation strategies: multiple variations of pretraining, fact-aware embeddings and structure-aware input encoding. Our extensive experiments show that a multi-lingual mT5 model which uses fact-aware embeddings with structure-aware input encoding leads to best results (30.90 BLEU, 55.12 METEOR and 59.17 chrF++) across the twelve languages. We make our code and dataset publicly available[1], and hope that this will help advance further research in this critical area.

## 1 Introduction

Fact-to-text (F2T) is a natural language generation (NLG) task where input is structured data (like facts[2]) and output is its natural language description. F2T systems have been shown to be effective in many applications like automated dialog

---

[1] https://github.com/blitzprecision/
XAlignV2

[2] A fact is a triple composed of subject, relation and object.



Figure 1: XF2T example from XALIGNV2: Generating English, Hindi, Bengali, Gujarati, Tamil and Punjabi sentences to capture semantics from English facts.

systems (Wen et al., 2016), domain-specific chatbots (Novikova et al., 2017), open domain question answering (Chen et al., 2020), authoring sports reports (Chen and Mooney, 2008), financial reports (Plachouras et al., 2016), news reports (Leppänen et al., 2017), etc. Recently, several English F2T systems have been proposed, but lack of training data in low-resource languages (LRLs) implies that there are hardly any such systems for LRLs.

Across many business domains, there is abundance of facts (or key-value stores) in English, and consumers want to access that information in their own regional languages. For example, users want product descriptions, weather report, match report, financial report in various LRLs. Another related problem is to automatically populate first sentence for LRL Wikipedia pages using facts from English Wikidata. If such facts were in LRLs and there were models to do F2T in those LRLs, we could leverage those. However, neither exist. Even LRL facts on Wikidata are very sparse. Another approach could be to do F2T in English and then translate the output to LRLs. But our experiments show that this leads to poor quality primarily due to lack of robust translation systems for LRLs.

Specifically, we focus on the F2T problem of

15

| Dataset | Languages | A/M | \|I\| | F/I | \|P\| | \|T\| | X-Lingual |
|---|---|---|---|---|---|---|---|
| WikiBio | en | A | 728K | 19.70 | 1740 | 26.1 | No |
| E2E | en | M | 50K | 5.43 | 945 | 20.1 | No |
| WebNLG 2017 | en | M | 25K | 2.95 | 373 | 22.7 | No |
| fr-de Bio | fr, de | A | 170K, 50K | 8.60, 12.6 | 1331, 1267 | 29.5, 26.4 | No |
| TREX | en | A | 6.4M | 1.77 | 642 | 79.8 | No |
| WebNLG 2020 | en, ru | M | 40K, 17K | 2.68, 2.55 | 372, 226 | 23.7 | Yes |
| KELM | en | A | 8M | 2.02 | 663 | 21.2 | No |
| WITA | en | A | 55K | 3.00 | 640 | 18.8 | No |
| WikiTableT | en | A | 1.5M | 51.90 | 3K | 115.9 | No |
| GenWiki | en | A | 1.3M | 1.95 | 290 | 21.5 | No |
| XALIGN | en + 7 LR | A | 0.45M | 2.02 | 367 | 19.8 | Yes |
| XALIGNV2 | en + 11 LR | A | 0.55M | 1.98 | 374 | 19.7 | Yes |

Table 1: Statistics of popular Fact-to-Text datasets: WikiBio (Lebret et al., 2016), E2E (Novikova et al., 2017), WebNLG 2017 (Gardent et al., 2017), WebNLG 2020 (Ferreira et al., 2020), fr-de Bio (Nema et al., 2018), KELM (Agarwal et al., 2021), WITA (Fu et al., 2020), WikiTableT (Chen et al., 2021), GenWiki (Jin et al., 2020), TREX (Elsahar et al., 2018), XAlign (Abhishek et al., 2022), and XALIGNV2 (ours). Alignment method could be A (automatic) or M (manual). |I|=number of instances. F/I=number of facts per instance. |P|=number of unique relations. |T|=average number of tokens per instance.

generating LRL person biographies (like a sentence on Wikipedia page) from English Wikidata facts. While millions of English person entities exist on Wikidata, there are a total of only 168K (non-unique) person Wikidata entries across 11 LRLs of our interest. As an extreme, Assamese has only 1.7K person entries! Even worse, average number of facts per entity on Wikidata in LRLs (10.39) is less than half of that of English (22.8). Monolingual F2T for LRLs suffers from lack of training data. Translating English output (using English F2T) to LRLs leads to poor results. This necessitates us to build *cross-lingual F2T generation (XF2T)* systems, wherein the input is a set of English facts and output is a sentence capturing the fact-semantics in the specified LR language, as introduced in our previous work (Abhishek et al., 2022).

In (Abhishek et al., 2022), we proposed transfer learning and distance supervision based methods for cross-lingual alignment for aligning English Wikidata facts with equivalent text from LRL Wikipedia pages. In that paper, we used such alignment methods to contribute the XALIGN dataset which consists of sentences from LR language Wikipedia aligned with English fact triples from Wikidata. It contains data for the following eight languages: Hindi (hi), Telugu (te), Bengali (bn), Gujarati (gu), Marathi (mr), Kannada (kn), Tamil (ta) and English (en). In that paper, we focused on dataset creation and not much on the XF2T task. In this paper, we extend this dataset to four more LR languages: Punjabi (pa), Malayalam (ml), Assamese (as) and Oriya (or). Fig. 1 shows an XF2T example from our extended dataset, XALIGNV2. Further, we rigorously investigate models for the

XF2T problem. First, we experiment with standard existing Transformer-based multi-lingual encoder-decoder models like the vanilla Transformer, IndicBART and mT5. Next, we explore performance across various training setups: bi-lingual, translate-output, translate-input and multi-lingual. Further, we systematically explore various strategies for improving XF2T generation like multi-lingual data-to-text pre-training, fact-aware embeddings, and structure-aware encoding. Overall, we make the following contributions in this work.

- We extend the XALIGN dataset with annotated XF2T data corresponding to four more LR languages, leading to a new dataset, XALIGNV2.

- We rigorously experiment with multiple encoder-decoder models, training setups, pre-training methods, and input representations toward building a robust XF2T system.

- We find that a multi-lingual mT5 model which uses fact-aware embeddings along with structure-aware input encoding leads to best results. Our best small-scale model achieves an average BLEU of 29.27, METEOR of 53.64, and chrF++ of 57.30 for XF2T across 12 languages. We make the code and dataset publicly available[1].

## 2   Related Work

**Multi-lingual and Cross-lingual NLG**: Recently there has been a lot of work on multi-lingual and cross-lingual NLG tasks like machine translation (Chi et al., 2021; Liu et al., 2020), question generation (Chi et al., 2020; Mitra et al., 2021),

news title generation (Liang et al., 2020), and summarization (Zhu et al., 2019; Taunk et al., 2023) thanks to models like XNLG (Chi et al., 2020), mBART (Liu et al., 2020), mT5 (Xue et al., 2021), etc. In this work, we investigate effectiveness of multiple modeling techniques for the *XF2T* task. Further, from a knowledge graph (KG) and text linking perspective, our work is related to tasks like entity linking (link mention in a sentence to a KG entity) (Botha et al., 2020) and fact linking (linking sentence to a set of facts) (Kolluru et al., 2021). As against this, XF2T is the problem of generating a sentence given a set of facts. XF2T is also related to graph-to-text (Ribeiro et al., 2021) where our fact triples about an entity can be mapped to a star-like graph, but no cross-lingual graph-to-text methods exist unfortunately.

**F2T Datasets**: Several F2T datasets have been proposed in the literature: WikiBio (Lebret et al., 2016), E2E (Novikova et al., 2017), WebNLG 2017 (Gardent et al., 2017), WebNLG 2020 (Ferreira et al., 2020), fr-de Bio (Nema et al., 2018), KELM (Agarwal et al., 2021), WITA (Fu et al., 2020), WikiTableT (Chen et al., 2021), GenWiki (Jin et al., 2020), TREX (Elsahar et al., 2018) and XAlign (Abhishek et al., 2022). These datasets contain text from various domains like people, sports, restaurants, airports, politicians, artists, etc. Also, these datasets vary widely in terms of statistics like the number of instances, number of facts per instance, number of unique relations and average number of tokens per instance. All of these are English only except fr-de Bio (which has French and German), WebNLG 2020 (which has English and Russian) and XAlign (which has English and 7 other LR languages). Both fr-de Bio and WebNLG 2020 propose multi-lingual but not cross-lingual F2T tasks. Unlike other datasets, XALIGN and our dataset, XALIGNV2 are cross-lingual. Our proposed dataset, XALIGNV2, contains 12 languages, has 0.55M instances, 374 unique relations, avg 19.7 tokens/instance and avg 1.98 facts/instance. Table 1 shows basic statistics of popular F2T datasets.

**F2T Generation**: Training F2T models requires aligned data with adequate content overlap. Some previous studies like WebNLG (Gardent et al., 2017) collected aligned data by crowdsourcing, while others have performed automatic alignment by heuristics like TF-IDF. In (Abhishek et al., 2022), we explored two unsupervised methods to perform a cross-lingual alignment. We leverage the "transfer learning from Natural Language Inference task" based method for this work.

Initial F2T methods were template-based and were therefore proposed on domain-specific data like medical (Bontcheva and Wilks, 2004), cooking (Cimiano et al., 2013), person (Duma and Klein, 2013), etc. They align entities in RDF triples with entities mentioned in sentences, extract templates from the aligned sentences, and use templates to generate sentences given facts for new entities. Template-based methods are brittle and do not generalize well. Recently, Seq-2-seq neural methods (Lebret et al., 2016; Mei et al., 2016) have become popular for F2T. These include vanilla LSTMs (Vougiouklis et al., 2018), LSTM encoder-decoder model with copy mechanism (Shahidi et al., 2020), LSTMs with hierarchical attentive encoder (Nema et al., 2018), pretrained Transformer based models (Ribeiro et al., 2021) like BART (Lewis et al., 2020) and T5 (Raffel et al., 2020). Vougiouklis et al. (2018) proposed a method which uses feedforward neural networks to encode RDF triples and concatenate them as the input of the LSTM decoder. Variations of LSTM encoder-decoder model with copy mechanism (Shahidi et al., 2020) or with hierarchical attentive encoder (Nema et al., 2018) have also been proposed. Recently, pretrained Transformer based models like BART (Lewis et al., 2020) and T5 (Raffel et al., 2020) have been applied for monolingual English Fact-to-Text (Ribeiro et al., 2021).

Richer encoding of the input triples has also been investigated using a combination of graph convolutional networks and Transformers (Zhao et al., 2020), triple hierarchical attention networks (Chen et al., 2020), or Transformer networks with special fact-aware input embeddings (Chen et al., 2020). Some recent work also explores specific F2T settings like plan generation when the order of occurrence of facts in text is available (Zhao et al., 2020) or partially aligned F2T when the text covers more facts than those mentioned in the input (Fu et al., 2020). However, all of these methods focus on English fact to text only. Only recently, we proposed the XF2T problem in our previous paper Abhishek et al. (2022) but in that work, our focus was on problem formulation and dataset contribution. In this paper, we extensively evaluate multiple methods for the XF2T generation task.

| | |V| | Train+Validation | | | Manually Labeled Test | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | |I| | |T| | |F| | κ | |A| | |I| | |T| | |F| |
| hi | 75K | 57K | 25.3/5/99 | 2.0 | 0.81 | 4 | 842 | 11.1/5/24 | 2.1 |
| mr | 50K | 19K | 20.4/5/94 | 2.2 | 0.61 | 4 | 736 | 12.7/6/40 | 2.1 |
| te | 61K | 24K | 15.6/5/97 | 1.7 | 0.56 | 2 | 734 | 9.7/5/30 | 2.2 |
| ta | 121K | 57K | 16.7/5/97 | 1.8 | 0.76 | 2 | 656 | 9.5/5/24 | 1.9 |
| en | 104K | 133K | 20.2/4/86 | 2.2 | 0.74 | 4 | 470 | 17.5/8/61 | 2.7 |
| gu | 35K | 9K | 23.4/5/99 | 1.8 | 0.50 | 3 | 530 | 12.7/6/31 | 2.1 |
| bn | 131K | 121K | 19.3/5/99 | 2.0 | 0.64 | 2 | 792 | 8.7/5/24 | 1.6 |
| kn | 88K | 25K | 19.3/5/99 | 1.9 | 0.54 | 4 | 642 | 10.4/6/45 | 2.2 |
| pa | 59K | 30K | 32.1/5/99 | 2.1 | 0.54 | 3 | 529 | 13.4/5/45 | 2.4 |
| as | 27K | 9K | 19.23/5/99 | 1.6 | - | 1 | 637 | 16.22/5/72 | 2.2 |
| or | 28K | 14K | 16.88/5/99 | 1.7 | - | 2 | 242 | 13.45/7/30 | 2.6 |
| ml | 146K | 55K | 15.7/5/98 | 1.9 | 0.52 | 2 | 615 | 9.2/6/24 | 1.8 |

Table 2: Basic Statistics of XALIGNV2. |I|=# instances, |T|=avg/min/max word count, |F|=avg #facts, |V|=Vocab. size, $\kappa$=Kappa score, |A|=#annotators. For Train+Validation, min and max fact count is 1 and 10 resp across languages.[4]

## 3 XALIGNV2: Data Collection, Pre-processing and Alignment

**Data Collection and Pre-processing**: We start by gathering a list of ∼95K person entities from Wikidata each of which has a link to a corresponding Wikipedia page in at least one of our 11 LR languages. This leads to a dataset $D$ where every instance $d_i$ is a tuple ⟨entityID, English Wikidata facts, LRL, LRL Wikipedia URL for the entityID⟩. We extract facts (in English) from the 20201221 WikiData dump for each entity in $D$ using the Wiki-Data API[3]. We gathered facts corresponding to only the following Wikidata property (or relation) types that capture most useful factual information for person entities: WikibaseItem, Time, Quantity, and Monolingualtext. We retain any supporting information associated with the fact triple as a fact qualifier. This leads to overall ∼0.55M data instances across all the 12 languages. Also, for each language, we gather sentences (along with section information) from 20210520 Wikipedia XML dump using same pre-processing steps as described in (Abhishek et al., 2022).

**Fact-to-Text Alignment**: For every (entity $e$, language $l$) pair, the pre-processed dataset has a set $F_{el}$ of English Wikidata facts and a set of Wikipedia sentences $S_{el}$ in that language. Next, we use a two-stage automatic aligner as proposed in (Abhishek et al., 2022) to associate a sentence in $S_{el}$ with a subset of facts from $F_{el}$. We run this aligner for the new four LR languages to obtain the corresponding Train+Validation part of XALIGNV2.

---

**Manual Annotations for Ground-Truth Data**: We need manually annotated data for evaluation of our XF2T generation. Again, we follow the same procedure as outlined in (Abhishek et al., 2022) to get annotations for the new four languages in XALIGNV2. Detailed annotation guidelines are also mentioned here[1]. Our annotator pool is selected from the National Register of Translators[5]. Annotators were in age range 25 to 40 years; 46% females and 54% males; occupations varied as linguists, editors, translators, freelancers; qualifications varied as BA, MA, MSc, LLB, PhD. We report details of this test part of our XALIGNV2 dataset in Table 2. On average, a sentence can be verbalized using ∼2 fact triples.

**XALIGNV2 Dataset Analysis**: Table 2 shows the dataset statistics. Figs. 2 and 3 show fact count distribution. We observe that a large percent of sentences contain more than one fact across languages. Also, the distribution is similar across languages and data subsets. Finally, Table 3 shows top 10 frequent fact relations across all the languages.

## 4 XF2T Approaches

In this section, we first discuss our input representation. Next, we discuss various Transformer-based methods, different training setups, multiple pretraining methods, and discussion on fact-aware embeddings.

**Structure-aware Input encoding**: Each input instance consists of multiple facts $F = \{f_1, f_2, \ldots, f_n\}$ and a section title $t$. A fact $f_i$ is a tuple composed of subject $s_i$, relation $r_i$, object $o_i$ and $m$ qualifiers $Q = q_1, q_2, \ldots, q_m$. Each qualifier provides more information about the fact. Each of the qualifiers $\{q_j\}_{j=1}^m$ can be linked to the fact using a fact-level property which we call as qualifier relation $qr_j$. For example, consider the sentence: "Narendra Modi was the Chief Minister of Gujarat from 7 October 2001 to 22 May 2014, preceded by Keshubhai Patel and succeeded by Anandiben Patel." This can be represented by a fact where subject is "Narendra Modi", relation is "position held", object is "Chief Minister of Gujarat" and there are 4 qualifiers each with their qualifier relations as follows: (1) $q_1$="7 October 2001", $qr_1$="start time", (2) $q_2$="22 May 2014", $qr_2$="end time", (3) $q_3$="Keshubhai Patel", $qr_3$="replaces", and (4) $q_4$="Anandiben Patel", $qr_4$="replaced by".
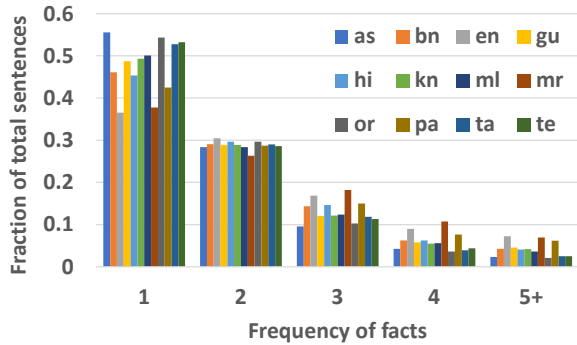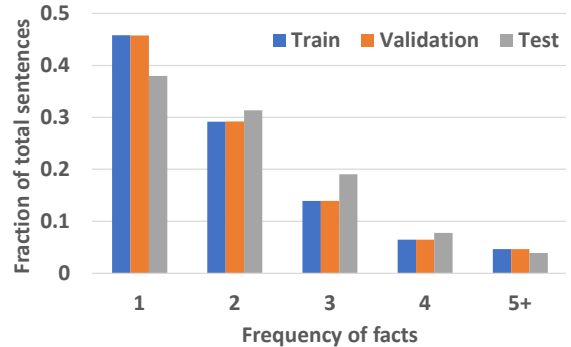
---

Figure 2: Fact Count Distribution across languages



Figure 3: Fact Count Distribution across data subsets

| | |
|---|---|
| hi | occupation, date of birth, position held, cast member, country of citizenship, award received, place of birth, date of death, educated at, languages spoken written or signed |
| mr | occupation, date of birth, position held, date of death, country of citizenship, place of birth, member of sports team, member of political party, cast member, award received |
| te | occupation, date of birth, position held, cast member, date of death, place of birth, award received, member of political party, country of citizenship, educated at |
| ta | occupation, position held, date of birth, cast member, country of citizenship, educated at, place of birth, date of death, award received, member of political party |
| en | occupation, date of birth, position held, country of citizenship, educated at, date of death, award received, place of birth, member of sports team, member of political party |
| gu | occupation, date of birth, cast member, position held, award received, date of death, languages spoken written or signed, place of birth, author, country of citizenship |
| bn | occupation, date of birth, country of citizenship, cast member, member of sports team, date of death, educated at, place of birth, position held, award received |
| kn | occupation, cast member, date of birth, award received, position held, date of death, performer, place of birth, author, educated at |
| pa | occupation, date of birth, place of birth, date of death, cast member, country of citizenship, educated at, award received, languages spoken, written or signed, position held |
| as | occupation, date of birth, cast member, position held, date of death, place of birth, country of citizenship, educated at, award received, member of political party |
| or | occupation, date of birth, position held, cast member, member of political party, place of birth, date of death, award received, languages spoken, written or signed, educated at |
| ml | occupation, cast member, position held, date of birth, educated at, award received, date of death, place of birth, author, employer |

Table 3: Top-10 frequent fact relations across languages.

Each fact $f_i$ is encoded as a string and the overall input consists of a concatenation of such strings across all facts in $F$. The string representation for a fact $f_i$ is "$\langle S \rangle s_i \langle R \rangle r_i \langle O \rangle o_i \langle R \rangle qr_{i_1} \langle O \rangle q_{i_1} \langle R \rangle qr_{i_2} \langle O \rangle q_{i_2} \ldots \langle R \rangle qr_{i_m} \langle O \rangle q_{i_m}$" where $\langle S \rangle$, $\langle R \rangle$, $\langle O \rangle$ are special tokens. Finally, the overall input with $n$ facts is obtained as follows: "generate [language] $f_1$ $f_2$ $\ldots f_n \langle T \rangle [t]$" where "[language]" is one of our 12 languages, $\langle T \rangle$ is the section title delimiter token, and $t$ is the section title.

**Standard Transformer-based Models**: For XF2T generation, we train multiple popular multilingual text generation models on Train+Validation part of our XALIGN dataset. We use a basic Transformer model, mT5-small, and the IndicBART (Dabre et al., 2021) for the XF2T task. We do not experiment with mBART (Liu et al., 2020) and Muril (Khanuja et al., 2021) since their small sized model checkpoints are not publicly available. We train these models in a multi-lingual cross-lingual manner. Thus, we train a single model using training data across languages without

any need for translation.

**Bi-lingual, Multi-lingual & Translation models**: Next, we experiment with different training setups. We first build bilingual models, where input is in English and output could be in any of the 12 languages. A drawback with this approach is the need to maintain one model per language which is cumbersome.

Further, we also train two translation based models. In the "translate-output" setting, we train a single English-only model which consumes English facts and generates English text. The English output is translated to desired language at test time using IndicTrans (Ramesh et al., 2021). In the "translate-input" setting, English facts are translated to LR language and fed as input to train a single multi-lingual model across all languages. While translating if mapped strings for entities were present in Wikidata they were directly used. A drawback with these approaches is the need for translation at test time.

**Pretraining approaches**: Pretraining has been a standard method to obtain very effective models

19

even with small amounts of labeled data across several tasks in natural language processing (NLP). Domain and task specific pretraining has been shown to provide further gains (Gururangan et al., 2020). We experiment with the following four pretraining strategies on top of the already pretrained encoder-decoder model before finetuning it on XALIGNV2 dataset. (1) Multi-lingual pretraining: Wang et al. (2021) provide a noisy, but larger corpus (542192 data pairs across 15 categories) crawled from Wikipedia for English F2T task. The dataset is obtained by coupling noisy English Wikipedia data with Wikidata triples. We translate English sentences from the Wikipedia-based Wang et al. (2021)'s data to our LR languages. Thus, the multi-lingual pretraining data contains ∼6.5M data pairs. For translating sentences, we use IndicTrans (Ramesh et al., 2021). (2) Translation-based pretraining: Translation is a preliminary task for effective cross-lingual NLP. Thus, in this method, we pretrain mT5 on translation data corresponding to English to other language pairs with ∼0.25M data instances per language. (3) Two-stage pretraining: This combines the above two methods. In the first stage, we do translation-based pretraining. In the second stage, we perform multi-lingual pretraining. (4) Multi-task pretraining: This method also involves training for both translation as well as XF2T tasks. Unlike the two-stage method where pretraining is first done for translation and then for XF2T (multi-lingual pretraining), in this method we perform the two tasks jointly in a multi-task learning setup.

**Fact-aware embeddings**: The input to mT5 consists of token embeddings as well as position embeddings. For XF2T, the input is a bunch of facts. Facts contain semantically separate units each of which play a different role: subject, relation, object. We extend the standard mT5 input with specific (fact-aware) role embeddings. Specifically, we use four role IDs: ROL1 for subject, ROL2 for relation and qualifier relation, ROL3 for object and qualifier tokens, and ROL0 for everything else, as shown in Fig. 4. These are randomly initialized and learned while training. We hope that this explicit indication of the role played by each token in the input facts, will help the model for improved XF2T generation.

We also experimented with (1) separate role embeddings for qualifier relation and qualifier, and (2) adding fact id embeddings, i.e., if the input contains $K$ facts, we have $K$ fact IDs, and all tokens corresponding to a fact gets the same fact ID embedding. However, these did not lead to better results and thus we do not report those results.

## 5   Experiments

**Implementation Details for Reproducibility**:
We closely follow Abhishek et al. (2022)'s data-collection and XF2T alignment method for the creation of cross-lingual fact-to-text dataset for four additional languages. All XF2T generation approaches were run on a machine equipped with four 32GB V100 GPUs. For all experiments, we use IndicNLP (Kakwani et al., 2020) to convert the low-resource languages of XALIGNV2 to the unified Devanagari script. All Transformer models have 6 encoder and 6 decoder layers. For Vanilla Transformer, we follow the standard architecture and hyper-parameters suggested by Vaswani et al. (2017). For other methods, we optimize cross entropy loss using AdamW with constant learning rate of 3e-5 with L2-norm weight decay of 0.001, batch size of 20 and dropout of 0.1. We closely follow (Dabre et al., 2021) for finetuning IndicBart.

When applicable, we pretrain for 7 epochs. For multi-lingual pretraining, we use full validation set. In two-stage pretraining, we save best checkpoint of first stage (translation task) on validation set of translation task and use it to initialize model parameters for second stage. For multi-task pretraining, we create new validation set by combining validation set of translation task and XF2T task. We finetune for 30 epochs and use beam search with width of 4.

**Evaluation Metrics**: We use overall BLEU scores (Ramesh et al., 2021) for evaluating the multi-lingual models for English-Indic fact-sentence pairs. Following previous work, we also use METEOR (Banerjee and Lavie, 2005) and chrF++ (Popović, 2017). PARENT (Dhingra et al., 2019) relies on the word overlap between input and the prediction text. Since the input and prediction in XF2T are in different languages, we cannot compute PARENT scores.

## 6   Results and Analysis

Since XF2T is a very recently proposed task, there are not many baseline methods to compare with. In this section, we will present results using methods described in Section 4. Due to lack of space, we show per language results only for our best model,

Figure 4: English facts being passed as input to mT5's encoder with token, position and (fact-aware) role embeddings.

but present language-wise results for other models in the Appendix. For other comparisons and analysis, we show average across all languages while pointing out any interesting per-language insights.

| | BLEU | METEOR | chrF++ |
|---|---|---|---|
| Vanilla Transformer | 21.93 | 50.21 | 50.89 |
| IndicBART | 23.78 | 50.80 | 53.88 |
| mT5 | **28.13** | **53.54** | **57.27** |

Table 4: XF2T scores on XALIGNV2 test set using standard Transformer-based encoder-decoder models. The best results are highlighted.

| | BLEU | METEOR | chrF++ |
|---|---|---|---|
| Bi-lingual mT5 (12 models) | 25.88 | 50.91 | 52.88 |
| Translate-Output mT5 (1 model) | 18.91 | 42.83 | 49.10 |
| Translate-Input mT5 (1 model) | 26.53 | 52.24 | 55.32 |
| Multi-lingual mT5 (1 model) | **28.13** | **53.54** | **57.27** |

Table 5: XF2T scores on XALIGNV2 test set using bi-lingual, multi-lingual and translation-based variants of mT5 model. Best results are highlighted.

**Standard Transformer-based Models**: Table 4 shows BLEU results across different (model, metric) combinations using three standard Transformer-based encoder-decoder models. Across the 12 languages, on average for each metric, mT5 performs better than IndicBART, which is better than vanilla Transformer. We observed that IndicBART performed exceptionally well for Bengali but is exceptionally poor on English. Given that mT5 is better on average amongst the three, we perform further experiments using mT5.

| No. | Method | BLEU | METEOR | chrF++ |
|---|---|---|---|---|
| 1 | No pretraining and no fact-aware embeddings | 28.13 | 53.54 | 57.27 |
| 2 | Two-stage Pretraining | 27.70 | 51.87 | 55.32 |
| 3 | Multi-task Pretraining | 28.45 | 51.87 | 55.20 |
| 4 | Translation-based Pretraining | 27.53 | 50.67 | 53.71 |
| 5 | Multi-lingual Pretraining | 28.71 | **53.83** | **57.58** |
| 6 | Fact-aware embeddings | **29.27** | 53.64 | 57.30 |

Table 6: XF2T scores on XALIGNV2 test set using different pretraining strategies and fact-aware embeddings for the mT5 model. Best results are highlighted. Row 1 is same as last row from Table 5.

**Bi-lingual, Multi-lingual & Translation models**: Table 5 shows results when mT5 model is trained using various bi-lingual, multi-lingual and translation-based settings. We observe that across all settings, the initial setting of training a single multi-lingual cross-lingual model is the best on average across all metrics. That said, for Bengali, a bi-lingual model, i.e., a model specifically trained for en→bn, is much better[6]. Translate-output and translate-input settings lead to slightly improved models for English and Tamil respectively. On average, translate-output setting performs the worst while the multi-lingual setting performs the best. Although we use the state-of-the-art translation method, we believe low accuracy for translate-output setting is mainly due to poor translation quality.

**Pretraining approaches**: Table 6 (lines 1 to 5) shows results using different pretraining strategies. We observe that multi-lingual pretraining leads to improvements compared to no XF2T specific pretraining across 2 of the 3 metrics. Two-stage pretraining is slightly better than translation-based pretraining but not as good as multi-lingual pretraining. Finally, multi-task performs better than two-stage. For English and Bengali, we found that two-stage pretraining provided best results. However, multi-lingual pretraining is the best on average across languages, with biggest wins for Malayalam and Oriya.

**Fact-aware embeddings**: Table 6 (line 6) shows that fact-aware embeddings lead to improvements over the vanilla mT5 method without fact-aware embeddings (line 1).

In summary, we note that both the proposed methods (multi-lingual pretraining, fact-aware embedding) lead to improvements over the vanilla mT5. We also experimented with combinations

---

[6]Even later we observe that translation-only pretraining helps improve Bengali performance. We hypothesize this is because of huge influence English has had over Bengali historically.

| | Vanilla mT5 | | | Multi-lingual Pretraining | | | Fact-aware embeddings | | |
|---|---|---|---|---|---|---|---|---|---|
| | BLEU | METEOR | chrF++ | BLEU | METEOR | chrF++ | BLEU | METEOR | chrF++ |
| hi | 44.65 | 68.58 | 68.49 | 43.32 | 68.19 | 68.21 | 42.72 | 67.49 | 68.03 |
| mr | 26.47 | 56.85 | 59.17 | 27.64 | 56.34 | 57.74 | 29.06 | 55.40 | 57.97 |
| te | 14.46 | 43.45 | 52.58 | 15.94 | 42.71 | 52.40 | 16.21 | 42.14 | 51.25 |
| ta | 18.37 | 46.15 | 57.42 | 16.68 | 42.32 | 54.88 | 19.07 | 43.65 | 56.01 |
| en | 46.94 | 70.60 | 65.20 | 46.61 | 70.45 | 65.33 | 48.29 | 70.75 | 65.42 |
| gu | 22.69 | 50.31 | 51.36 | 21.39 | 47.98 | 50.14 | 23.27 | 50.00 | 50.64 |
| bn | 40.38 | 61.71 | 68.71 | 50.89 | 75.62 | 77.43 | 49.48 | 73.03 | 76.19 |
| kn | 10.66 | 32.58 | 46.92 | 11.61 | 33.00 | 47.18 | 11.57 | 33.44 | 46.66 |
| ml | 26.22 | 56.71 | 57.01 | 27.38 | 56.63 | 57.35 | 29.04 | 57.15 | 57.60 |
| pa | 26.96 | 54.82 | 52.33 | 26.04 | 54.17 | 52.50 | 28.65 | 55.19 | 53.38 |
| or | 47.17 | 67.82 | 71.20 | 44.97 | 66.49 | 70.64 | 41.75 | 63.77 | 67.96 |
| as | 12.61 | 32.93 | 36.91 | 12.00 | 32.04 | 37.15 | 12.16 | 31.61 | 36.44 |
| Avg | 28.13 | 53.54 | 57.27 | 28.71 | 53.83 | 57.58 | 29.27 | 53.64 | 57.30 |

Table 7: XF2T scores on XALIGNV2 test set using vanilla mT5, multi-lingual pretrained mT5 and mT5 with fact-aware embedding models.
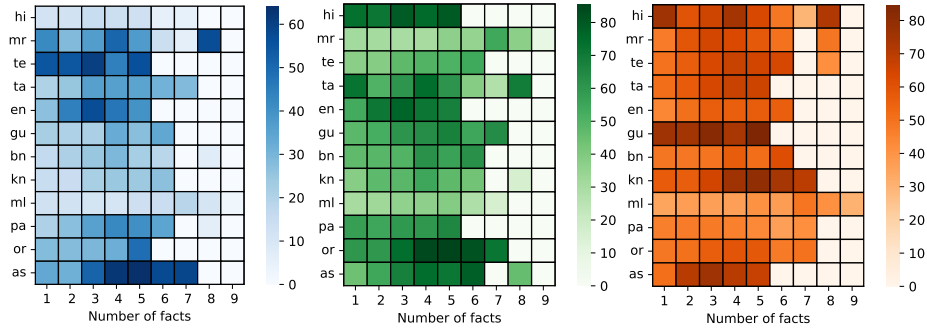


Figure 5: BLEU (left), METEOR (middle) and chrF++ (right) scores for the best model across languages for test instances with a specific number of facts. White cells indicate absence of instances in that cell.

| Lang. | Input | Reference Text | Generated Text |
|---|---|---|---|
| hi | generate Hindi <S> Asha Nautiyal <R> member of political party <O> Bharatiya Janata Party <R> date of birth <O> 25 June 1969 <R> occupation <O> politician <R> country of citizenship <O> India <T> introduction | आशा नौटियाल ( जन्म 25 जून, १९६९) एक भारतीय राजनीतिज्ञ हैं जो भारतीय जनता पार्टी से हैं। | आशा नौटियाल ( जन्म 25 जून, १९६९) एक भारतीय राजनीतिज्ञ और भारतीय जनता पार्टी की सदस्य हैं। |
| en | generate English <S> Kedarnath Singh <R> date of death <O> 19 March 2018 <R> date of birth <O> 07 July 1934 <R> occupation <O> poet <R> languages spoken, written or signed <O> Hindi <R> country of citizenship <O> India <T> introduction | Kedarnath Singh ( 7 July 1934 - 19 March 2018 ) was an Indian poet who wrote in Hindi. | Kedarnath Singh ( 7 July 1934 - 19 March 2018 ) was a Hindi poet from Uttar Pradesh, India. |
| mr | generate Marathi <S> Théodore de Banville <R> date of death <O> 13 March 1891 <R> date of birth <O> 14 March 1823 <R> occupation <O> writer <R> country of citizenship <O> France <T> introduction | थेओदोर दि बॅनव्हिल ( मार्च १४, इ. स. १८२३ – मार्च १३, इ. स. १८९१ ) हा फ्रेंच साहित्यिक होता. | थॉडेर द बॅनव्हिल ( मार्च १४, इ. स. १८२३ – मार्च १३, इ. स. १८९१ ) हा फ्रेंच लेखक होता. |
| te | generate Telugu <S> Sushmita Sen <R> date of birth <O> 19 November 1975 <R> place of birth <O> Hyderabad <T> introduction | ఈమె 1975 నవంబరు 19 న హైదరాబాదులో జన్మించింది. | సుష్మితా సేన్ 1975 నవంబరు 19 న హైదరాబాదులో జన్మించింది. |
| ta | generate Tamil <S> Kirti Kumari <R> member of political party <O> Bharatiya Janata Party <R> date of birth <O> 13 August 1967 <R> date of death <O> 28 August 2017 <R> occupation <O> politician <R> country of citizenship <O> India <T> introduction | கீர்த்தி குமாரி ( 13 ஆகத்து 1967 - 28 ஆகத்து 2017 ) பாரதீய ஜனதா கட்சியின் இந்திய அரசியல்வாதி ஆவார். | கீர்த்தி குமாரி ( 13 ஆகத்து 1967 - 28 ஆகத்து 2017 ) ஓர் இந்திய அரசியல்வாதியும், பாரதீய ஜனதா கட்சியின் முன்னாள் சட்டமன்ற உறுப்பினரும் ஆவார். |
| kn | generate Kannada <S> Barry C. Barish <R> award received <O> Henry Draper Medal <R> point in time <O> 2017 <T> awards and honors | ಮತ್ತು ಬ್ಯಾರಿಶ್ ಅವರಿಗೆ ೨೦೦೭ ರ ಹೆನ್ರಿ ಡ್ರೇಪರ್ ಪದಕವನ್ನು ನೀಡಲಾಯಿತು. | ೨೦೦೭ ರಲ್ಲಿ ಅವರು ಹೆನ್ರಿ ಡ್ರೇಪರ್ ಪದಕವನ್ನು ಪಡೆದರು. |
| bn | generate Bengali <S> Jim Pothecary <R> member of sports team <O> South Africa national cricket team <R> occupation <O> cricketer <T> introduction | দক্ষিণ আফ্রিকা ক্রিকেট দলের অন্যতম সদস্য ছিলেন তিনি। | দক্ষিণ আফ্রিকা ক্রিকেট দলের অন্যতম সদস্য ছিলেন তিনি। |
| gu | generate Gujarati <S> Krishnalal Shridharani <R> date of birth <O> 16 September 1911 <R> date of death <O> 23 July 1960 <R> occupation <O> poet <R> occupation <O> playwright <R> languages spoken, written or signed <O> Gujarati <T> introduction | કૃષ્ણલાલ શ્રીધરાણી ( ૧૬ સપ્ટેમ્બર ૧૯૧૧ – ૨૩ જુલાઈ ૧૯૬૦ ) ગુજરાતી ભાષાના કવિ અને નાટ્યકાર હતા. | કૃષ્ણલાલ શ્રીધરાણી ( ૧૬ સપ્ટેમ્બર ૧૯૧૧ – ૨૩ જુલાઈ ૧૯૬૦ ) ગુજરાતી કવિ, નાટ્યકાર અને નાટ્યકાર હતા. |
| pa | generate Punjabi <S> Orhan Pamuk <R>award received <O> Nobel Prize in Literature <R> point in time <O> 2006 <R> date of birth <O> 07 June 1952 <R> occupation <O> novelist <R> languages spoken, written or signed <O> Turkish <T> introduction | ਓਰਹਾਨ ਪਾਮੁਕ ( ਜਨਮ 7 ਜੂਨ 1952 ) ਇੱਕ ਤੁਰਕੀ ਨਾਵਲਕਾਰ ਹੈ ਜਿਸ ਨੇ 2006 ਵਿੱਚ ਸਾਹਿਤ ਲਈ ਨੋਬਲ ਇਨਾਮ ਹਾਸਿਲ ਕੀਤਾ। | ਓਰਹਾਨ ਪਾਮੁਕ ( ਜਨਮ 7 ਜੂਨ 1952 ) ਇੱਕ ਤੁਰਕੀ ਨਾਵਲਕਾਰ ਹੈ ਜਿਸ ਨੂੰ 2006 ਵਿੱਚ ਸਾਹਿਤ ਲਈ ਨੋਬਲ ਪੁਰਸਕਾਰ ਨਾਲ ਸਨਮਾਨਿਤ ਕੀਤਾ ਗਿਆ । |
| ml | generate Malayalam <S> Naomi Scott <R> date of birth <O> 06 May 1993 <R> place of birth <O> London <R> country of citizenship <O> United Kingdom <T> introduction | 1993 മേയ് 6 ന് ഇംഗ്ലണ്ടിലെ ലണ്ടനിലാണ് സ്കോട്ട് ജനിച്ചത്| | 1993 മേയ് 6 ന് ഇംഗ്ലണ്ടിലെ ലണ്ടനിലാണ് സ്കോട്ട് ജനിച്ചത്| |
| or | generate Odia <S> Ajay Swain <R> award received <O> Odisha Sahitya Akademi Award <R> point in time <O> 2012 <T> introduction | ସେ ୨୦୧୨ ମସିହାରେ ଓଡ଼ିଶା ସାହିତ୍ୟ ଏକାଡେମୀ ପୁରସ୍କାର ଲାଭ କରିଥିଲେ। | ୨୦୧୨ ମସିହାରେ ଓଡ଼ିଶା ସାହିତ୍ୟ ଏକାଡେମୀ ପୁରସ୍କାର ଲାଭ କରିଥିଲେ। |
| te | generate Telugu <S> Sushmita Sen <R> date of birth <O> 19 November 1975 <R> place of birth <O> Hyderabad <T> introduction | ఈమె 1975 నవంబరు 19 న హైదరాబాదులో జన్మించింది. | సుష్మితా సేన్ 1975 నవంబరు 19 న హైదరాబాదులో జన్మించింది. |
| as | generate Assamese <S> Harishankar Parsai <R> date of death <O> 10 August 1995 <R> date of birth <O> 22 August 1922 <R> occupation <O> writer <R> country of citizenship <O> British India <R> country of citizenship <O> Dominion of India <R> occupation <O> author <T> introduction | হৰিশংকৰ পৰসাই ( ২২ আগষ্ট, ১৯২২ - ১০ আগষ্ট, ১৯৯৫ ) আছিল হিন্দী সাহিত্যৰ এগৰাকী প্ৰসিদ্ধ লেখক আৰু ব্যংগকাৰ । | হৰিশংকৰ পৰসাই ( ২২ আগষ্ট, ১৯২২ - ১০ আগষ্ট, ১৯৯৫ ) এজন ভাৰতীয় লেখক । |

Table 8: Test examples with reference text and predictions from our fact-aware embedding model.

22

| Input | Reference Text | Generated Text | Reason |
|---|---|---|---|
| <S> rajendra pratap singh <R> date_of_birth <O> 20 october 1954 <R> place_of_birth <O> pratapgarh uttar pradesh <T> early life and education | इनका जन्म 20 अक्टूबर 1954 को प्रतापगढ़ जिले में हुआ । | राजेन्द्र प्रताप सिंह का जन्म २० अक्टूबर १९५४ को उत्तर प्रदेश के प्रतापगढ़ में हुआ था । | The generated output covers more facts from input |
| <S> imam-ul-haq <R> date_of_birth <O> 12 december 1995 <R> place_of_birth <O> lahore <T> personal life | उनका जन्म 12 दिसंबर 1995 को लाहौर में हुआ था । | इमाम – उल – हक का जन्म १२ दिसम्बर १९९५ को लाहौर में हुआ था । | Better transliteration of dates |
| <S> ryan mclaren <R> country_of_citizenship <O> south africa <R> date_of_birth <O> 09 february 1983 <R> occupation <O> cricketer <T> introduction | रियान मैकलारेन ( जन्म: 9 फरवरी 1983 ) एक दक्षिण अफ़्रीकी क्रिकेट खिलाड़ी हैं । | रायन मैकलेरन ( जन्म 9 फरवरी 1983 ) एक दक्षिण अफ़्रीकी क्रिकेटर हैं । | Multiple possible transliterations of English word in phonetic languages like Hindi |
| <S> mirza salaamat ali dabeer <R> languages_spoken,_written_or_signed <O> urdu <R> occupation <O> poet <T> introduction | मिर्ज़ा सलामत अली दबीर उर्दू के एक कवि थे । | मिर्ज़ा सालमत अली दारबीर उर्दू भाषा के विख्यात साहित्यकार हैं । | Model outputs synonyms/related terms |

Table 9: Some examples of Hindi sentences generated by our best model which got low BLEU scores but are actually better than ground truth.

of these approaches but did not observe better results. Amongst these, multi-lingual pretraining performs the best on two of the metrics (METEOR and chrF++) while fact-aware embeddings perform best on BLEU. Hence, we present language-wise detailed comparison across these three models in Table 7. As expected, since we have relatively more training data for bn, en and hi, overall the models perform well on these languages. We also observe that the models do not perform very well for languages of the Dravidian family (te, ta and kn) even though ta has as many training examples as hi. Oriya is an exception – models perform almost as good as on Hindi, even with almost one third amount of training data. We hypothesize that this is limited variety in terms of types of person entities in Oriya compared to that in Hindi.

Fig. 5 shows BLEU, METEOR and chrF++ scores for the best model across languages for test instances with a specific number of facts. Number of facts per instance range from 1 to 9. We observe that the model performs best on instances with 2–4 facts across languages and across all metrics.

Table 8 shows XF2T prediction examples for our fact-aware embedding model. In general, across examples, we observe that the generated text is fluent and correct. Most of the input facts are covered by the generated sentence. Sometimes, though, the model hallucinates and brings in extra information in the output, e.g., for English, "Uttar Pradesh" is not mentioned as part of input facts.

**Scaling study**: So far we presented results using small-scale models. For the fact-aware embedding model, we also train a large scale checkpoint with 12 encoder and 12 decoder layers. We observe that it leads to a BLEU of 30.90, METEOR of 55.12 and chrF++ of 59.17 which is significantly better compared to the small model as expected.

**Human Evaluation Results**: Finally, we obtain human annotations to evaluate the perceived qual-

ity of the generated text. Table 10 shows results for our best model across three metrics: fluency, coverage and hallucination in the generated output. Higher the better. The evaluation has been done on 100 samples for 7 languages on a 5-point Likert scale per metric. The table shows values averaged across judgments from three annotators. Fluency checks for coherence and grammar correctness of generated output. Coverage verifies if most facts are captured in the sentence correctly. Absence of extra information verifies if the model does not generate any hallucinated information. Fluency, coverage and hallucination are 4.71, 4.31, 4.37 on average for our best model respectively.

Further, we observed that even though our models generate reasonable results, sometimes they are wrongly penalized using automated metrics for multiple reasons as shown in Table 9.

| | Fluency | Coverage | Hallucination |
|---|---|---|---|
| hi | 4.89 | 4.75 | 4.37 |
| ml | 4.87 | 4.42 | 4.73 |
| ta | 4.45 | 4.07 | 4.36 |
| te | 4.65 | 4.18 | 4.14 |
| pa | 4.69 | 4.23 | 4.29 |
| mr | 4.70 | 4.35 | 4.44 |
| en | 4.69 | 4.17 | 4.29 |

Table 10: Human Evaluation Results for our best model

## 7 Conclusion

In this paper, we worked on the XF2T problem. We contributed the XALIGNV2 dataset which has instances with English facts aligned to 12 languages. We investigated several multi-lingual Transformer methods with different training setups, pretraining setups and input representations. We obtained models with best metrics of 30.90 BLEU, 55.12 METEOR and 59.17 chrF++ for XF2T. We make our code and dataset[1] publicly available to empower future research in this critical area.

## 8 Ethical Concerns

We do not foresee any harmful uses of this technology. In fact, F2T generation systems are vital in many downstream Natural Language Processing (NLP) applications like automated dialog systems (Wen et al., 2016), domain-specific chatbots (Novikova et al., 2017), open domain question answering (Chen et al., 2020), authoring sports reports (Chen and Mooney, 2008), etc. We believe that these systems will be useful for powering business applications like Wikipedia text generation given English Infoboxes, automated generation of non-English product descriptions using English product attributes, etc.

As part of this work, we collected labeled data as discussed in Section 3. The dataset does not involve collection or storage of any personally identifiable information or offensive information at any stage. Human annotators were paid appropriately while performing data collection according to the standard wages set by National Translation Mission (`https://www.ntm.org.in/`) and mutually agreed upon. The data is publicly released under MIT Open-Source License. The annotation exercise was approved by the Institutional Review Board of our institute.

Usage of XALIGN dataset: Our usage was consistent with its intended use. The dataset was made available to us by the authors under MIT Open-Source License.

## References

Tushar Abhishek, Shivprasad Sagare, Bhavyajeet Singh, Anubhav Sharma, Manish Gupta, and Vasudeva Varma. 2022. Xalign: Cross-lingual fact-to-text alignment and generation for low-resource languages. In *The World Wide Web Conference*, pages 171–175.

O Agarwal, H Ge, S Shakeri, and R Al-Rfou. 2021. Knowledge graph based synthetic corpus generation for knowledge-enhanced language model pre-training. In *NAACL-HLT*, pages 3554–3565.

Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.

K Bontcheva and Y Wilks. 2004. Automatic report generation from ontologies: the miakt approach. In *Conf. on application of natural language to info. systems*, pages 324–335.

J A Botha, Z Shan, and D Gillick. 2020. Entity linking in 100 languages. In *EMNLP*, pages 7833–7845.

D L Chen and R J Mooney. 2008. Learning to sportscast: a test of grounded language acquisition. In *ICML*, pages 128–135.

M Chen, S Wiseman, and K Gimpel. 2021. Wikitablet: A large-scale data-to-text dataset for generating wikipedia article sections. In *ACL-IJCNLP Findings*, pages 193–209.

W Chen, Y Su, X Yan, and W Y Wang. 2020. Kgpt: Knowledge-grounded pre-training for data-to-text generation. *arXiv:2010.02307*.

Z Chi, L Dong, S Ma, S Huang, X-L Mao, H Huang, and F Wei. 2021. Mt6: Multilingual pretrained text-to-text transformer with translation pairs.

Z Chi, L Dong, F Wei, W Wang, X-L Mao, and H Huang. 2020. Cross-lingual natural language generation via pre-training. In *AAAI*, volume 34, pages 7570–7577.

P Cimiano, J Lüker, D Nagel, and C Unger. 2013. Exploiting ontology lexica for generating natural language texts from rdf data. In *European Workshop on NLG*, pages 10–19.

Raj Dabre, Himani Shrotriya, Anoop Kunchukuttan, Ratish Puduppully, Mitesh M Khapra, and Pratyush Kumar. 2021. Indicbart: A pre-trained model for natural language generation of indic languages. *arXiv preprint arXiv:2109.02903*.

Bhuwan Dhingra, Manaal Faruqui, Ankur Parikh, Ming-Wei Chang, Dipanjan Das, and William Cohen. 2019. Handling divergent reference texts when evaluating table-to-text generation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4884–4895.

D Duma and E Klein. 2013. Generating natural language from linked data: Unsupervised template extraction. In *IWCS*, pages 83–94.

H Elsahar, P Vougiouklis, A Remaci, C Gravier, J Hare, F Laforest, and E Simperl. 2018. T-rex: A large scale alignment of natural language with knowledge base triples. In *LREC*.

Thiago Ferreira, Claire Gardent, Nikolai Ilinykh, Chris Van Der Lee, Simon Mille, Diego Moussallem, and Anastasia Shimorina. 2020. The 2020 bilingual, bi-directional webnlg+ shared task overview and evaluation results (webnlg+ 2020). In *Proceedings of the 3rd International Workshop on Natural Language Generation from the Semantic Web (WebNLG+)*.

Z Fu, B Shi, W Lam, L Bing, and Z Liu. 2020. Partially-aligned data-to-text generation with distant supervision. *arXiv:2010.01268*.

C Gardent, A Shimorina, S Narayan, and L Perez-Beltrachini. 2017. The webnlg challenge: Generating text from rdf data. In *INLG*, pages 124–133.

Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A Smith. 2020. Don't stop pretraining: Adapt language models to domains and tasks. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360.

Z Jin, Q Guo, X Qiu, and Z Zhang. 2020. Genwiki: A dataset of 1.3 million content-sharing text and graphs for unsupervised graph-to-text generation. In *COLING*, pages 2398–2409.

Divyanshu Kakwani, Anoop Kunchukuttan, Satish Golla, NC Gokul, Avik Bhattacharyya, Mitesh M Khapra, and Pratyush Kumar. 2020. Indicnlpsuite: Monolingual corpora, evaluation benchmarks and pre-trained multilingual language models for indian languages. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4948–4961.

S Khanuja, D Bansal, S Mehtani, S Khosla, A Dey, B Gopalan, D K Margam, P Aggarwal, R T Nagipogu, S Dave, et al. 2021. Muril: Multilingual representations for indian languages. *arXiv:2103.10730*.

K Kolluru, M Rezk, P Verga, W W Cohen, and P Talukdar. 2021. Multilingual fact linking. In *AKBC*.

R Lebret, D Grangier, and M Auli. 2016. Neural text generation from structured data with application to the biography domain. In *EMNLP*, pages 1203–1213.

L Leppänen, M Munezero, M Granroth-Wilding, and H Toivonen. 2017. Data-driven news generation for automated journalism. In *INLG*, pages 188–197.

M Lewis, Y Liu, N Goyal, M Ghazvininejad, A Mohamed, O Levy, V Stoyanov, and L Zettlemoyer. 2020. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *ACL*, pages 7871–7880.

Y Liang, N Duan, Y Gong, N Wu, F Guo, W Qi, M Gong, L Shou, D Jiang, G Cao, et al. 2020. Xglue: A new benchmark dataset for cross-lingual pre-training, understanding and generation. *arXiv:2004.01401*.

Y Liu, J Gu, N Goyal, X Li, S Edunov, M Ghazvininejad, M Lewis, and L Zettlemoyer. 2020. Multilingual denoising pre-training for neural machine translation. *TACL*, 8:726–742.

H Mei, M Bansal, and M R Walter. 2016. What to talk about and how? selective gen. using lstms with coarse-to-fine alignment. In *NAACL-HLT*, pages 720–730.

Rajarshee Mitra, Rhea Jain, Aditya Srikanth Veerubhotla, and Manish Gupta. 2021. Zero-shot multilingual interrogative question generation for" people also ask" at bing. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 3414–3422.

P Nema, S Shetty, P Jain, A Laha, K Sankaranarayanan, and M M Khapra. 2018. Generating descriptions from structured data using a bifocal attention mechanism and gated orthogonalization. In *NAACL-HLT*, pages 1539–1550.

J Novikova, O Dušek, and V Rieser. 2017. The e2e dataset: New challenges for end-to-end generation. *arXiv:1706.09254*.

V Plachouras, C Smiley, H Bretz, O Taylor, J L Leidner, D Song, and F Schilder. 2016. Interacting with financial data using natural language. In *SIGIR*, pages 1121–1124.

Maja Popović. 2017. chrf++: words helping character n-grams. In *Proceedings of the second conference on machine translation*, pages 612–618.

C Raffel, N Shazeer, A Roberts, K Lee, S Narang, M Matena, Y Zhou, W Li, and P J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *JMLR*, 21:1–67.

G Ramesh, S Doddapaneni, A Bheemaraj, M Jobanputra, Raghavan AK, A Sharma, et al. 2021. Samanantar: The largest publicly available parallel corpora collection for 11 indic languages. *arXiv:2104.05596*.

Leonardo FR Ribeiro, Martin Schmitt, Hinrich Schütze, and Iryna Gurevych. 2021. Investigating pretrained language models for graph-to-text generation. In *Proceedings of the 3rd Workshop on Natural Language Processing for Conversational AI*, pages 211–227.

H Shahidi, M Li, and J Lin. 2020. Two birds, one stone: A simple, unified model for text generation from structured and unstructured data. In *ACL*, pages 3864–3870.

Dhaval Taunk, Shivprasad Sagare, Anupam Patil, Shivansh Subramanian, Manish Gupta, and Vasudeva Varma. 2023. Xwikigen: Cross-lingual summarization for encyclopedic text generation in low resource languages. In *Proceedings of the ACM Web Conference 2023*, pages 1703–1713.

A Vaswani, N Shazeer, N Parmar, J Uszkoreit, L Jones, A N Gomez, Ł Kaiser, and I Polosukhin. 2017. Attention is all you need. In *NIPS*, pages 5998–6008.

P Vougiouklis, H Elsahar, L-A Kaffee, C Gravier, F Laforest, J Hare, and E Simperl. 2018. Neural wikipedian: Generating textual summaries from knowledge base triples. *J. Web Semantics*, 52:1–15.

Qingyun Wang, Semih Yavuz, Xi Victoria Lin, Heng Ji, and Naznen Rajani. 2021. Stage-wise fine-tuning for graph-to-text generation. In *Proceedings of the*

*59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: Student Research Workshop*, pages 16–22.

T-H Wen, M Gasic, N Mrksic, L M Rojas-Barahona, P-H Su, D Vandyke, and S Young. 2016. Multi-domain neural network language generation for spoken dialogue systems. *arXiv:1603.01232*.

L Xue, N Constant, A Roberts, M Kale, R Al-Rfou, A Siddhant, A Barua, and C Raffel. 2021. mt5: A massively multilingual pre-trained text-to-text transformer. In *NAACL-HLT*, pages 483–498.

C Zhao, M Walker, and S Chaturvedi. 2020. Bridging the structural gap between encoding and decoding for data-to-text generation. In *ACL*, pages 2481–2491.

J Zhu, Q Wang, Y Wang, Y Zhou, J Zhang, S Wang, and C Zong. 2019. Ncls: Neural cross-lingual summarization. In *EMNLP-IJCNLP*, pages 3054–3064.

Table 12 shows detailed results per language using various bi-lingual, multi-lingual and translation-based settings.

## A   Limitations

In this work, we performed XF2T for a total of 12 languages. Clearly, the work can be extended to include many more low resource languages. Further, the amount of training data per language varies significantly. Gathering more labeled data across languages is difficult but should help improve accuracy of the trained models.

For some languages, finding qualified annotators was very difficult. For Assamese, we could obtain only one annotator. For Oriya, we found two annotators but due to their limited bandwidth, we did not get overlapping samples annotated by them and hence cannot compute inter-annotator agreement. While our annotation guidelines are clear, and inter-annotator agreement is high on most languages, we acknowledge that the annotation quality may have suffered for Assamese and Oriya.

The best automatic evaluation results from our models as well as human evaluation results show that there is a lot of scope for further work in this area.

## B   Detailed results

Table 11 shows detailed results per language. We observe that IndicBART performed exceptionally well for Bengali but is exceptionally poor on English.

|  | Vanilla Transformer | | | IndicBART | | | mT5 | | |
|---|---|---|---|---|---|---|---|---|---|
|  | BLEU | METEOR | chrF++ | BLEU | METEOR | chrF++ | BLEU | METEOR | chrF++ |
| hi | 35.04 | 63.46 | 60.85 | 40.44 | 66.41 | 66.27 | **44.65** | **68.58** | **68.49** |
| mr | 18.28 | 50.66 | 49.87 | **28.08** | 55.35 | 57.73 | 26.47 | **56.85** | **59.17** |
| te | 6.95 | 36.17 | 41.70 | **15.67** | 41.52 | 50.40 | 14.46 | **43.45** | **52.58** |
| ta | 14.67 | 44.64 | 53.03 | **19.37** | 45.78 | 56.63 | 18.37 | **46.15** | **57.42** |
| en | 37.12 | 65.32 | 59.69 | 10.47 | 42.35 | 34.35 | **46.94** | **70.60** | **65.20** |
| gu | 15.66 | 47.70 | 46.29 | 19.16 | 47.92 | 49.30 | **22.69** | **50.31** | **51.36** |
| bn | 48.55 | 74.18 | 75.68 | **55.90** | **79.29** | **80.51** | 40.38 | 61.71 | 68.71 |
| kn | 4.78 | 28.96 | 37.60 | 10.30 | **33.55** | 46.65 | **10.66** | 32.58 | **46.92** |
| ml | 16.29 | 50.84 | 47.26 | **27.41** | 56.27 | 56.80 | 26.22 | **56.71** | **57.01** |
| pa | 17.76 | 50.27 | 44.73 | 22.32 | 53.20 | 50.74 | **26.96** | **54.82** | **52.33** |
| or | 39.94 | 61.09 | 62.79 | 22.16 | 53.76 | 58.30 | **47.17** | **67.82** | **71.20** |
| as | 8.08 | 29.27 | 31.24 | **14.07** | **34.25** | **38.87** | 12.61 | 32.93 | 36.91 |
| Avg | 21.93 | 50.21 | 50.89 | 23.78 | 50.80 | 53.88 | **28.13** | **53.54** | **57.27** |

Table 11: XF2T scores on XALIGNV2 test set using standard Transformer-based encoder-decoder models. Best results for a (metric, language) combination are highlighted.

|  | Bi-lingual (12 models) | | | Translate-Output (1 model) | | | Translate-Input (1 model) | | | Multi-lingual (1 model) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | BLEU | METEOR | chrF++ | BLEU | METEOR | chrF++ | BLEU | METEOR | chrF++ | BLEU | METEOR | chrF++ |
| hi | 41.07 | 66.15 | 65.57 | 24.88 | 55.91 | 54.48 | 41.98 | 66.14 | 66.47 | **44.65** | **68.58** | **68.49** |
| mr | 16.74 | 49.36 | 48.40 | 20.62 | 46.87 | 52.23 | 24.90 | 54.56 | 57.25 | **26.47** | **56.85** | **59.17** |
| te | 12.23 | 37.85 | 44.94 | 14.13 | 38.69 | 50.36 | 13.11 | 40.83 | 49.64 | **14.46** | **43.45** | **52.58** |
| ta | 18.37 | **46.57** | 57.10 | 8.36 | 30.41 | 46.35 | **19.23** | 45.68 | **57.54** | 18.37 | 46.15 | 57.42 |
| en | 45.79 | 69.90 | 63.79 | **50.81** | 70.47 | **65.43** | 45.12 | 69.88 | 64.11 | 46.94 | **70.60** | 65.20 |
| gu | 12.49 | 38.73 | 37.01 | 18.23 | 42.25 | 46.27 | 20.84 | 48.71 | 49.30 | **22.69** | **50.31** | **51.36** |
| bn | **53.61** | **75.42** | **78.12** | 20.57 | 46.58 | 56.60 | 40.56 | 67.75 | 71.36 | 40.38 | 61.71 | 68.71 |
| kn | 8.71 | 31.02 | 41.16 | 7.93 | 27.58 | 44.47 | 7.75 | 30.82 | 41.44 | **10.66** | **32.58** | **46.92** |
| ml | 24.28 | 55.37 | 55.49 | 18.60 | 47.39 | 51.47 | 26.16 | 56.49 | **57.22** | **26.22** | **56.71** | 57.01 |
| pa | 21.92 | 51.10 | 47.82 | 26.24 | 53.18 | 51.57 | 24.42 | 51.64 | 49.28 | **26.96** | **54.82** | **52.33** |
| or | 45.53 | 62.91 | 65.30 | 9.37 | 29.40 | 37.80 | 43.43 | 64.12 | 65.20 | **47.17** | **67.82** | **71.20** |
| as | 9.76 | 26.48 | 29.80 | 7.15 | 25.25 | 32.19 | 10.89 | 30.27 | 35.00 | **12.61** | **32.93** | **36.91** |
| Avg | 25.88 | 50.91 | 52.88 | 18.91 | 42.83 | 49.10 | 26.53 | 52.24 | 55.32 | **28.13** | **53.54** | **57.27** |

Table 12: XF2T scores on XALIGNV2 test set using bi-lingual, multi-lingual and translation-based variants of mT5 model. Best results for a (metric, language) combination are highlighted.

# Preventing Generation of Verbatim Memorization in Language Models Gives a False Sense of Privacy

**Daphne Ippolito**[1]     **Florian Tramèr**[*2]     **Milad Nasr**[*1]

**Chiyuan Zhang**[*1]     **Matthew Jagielski**[*1]     **Katherine Lee**[*1,3]

**Christopher A. Choquette-Choo**[*1]     **Nicholas Carlini**[1]

[1] *Google Research*     [2] *ETH Zurich*     [3] *Cornell University*

## Abstract

Studying data memorization in neural language models helps us understand the risks (e.g., to privacy or copyright) associated with models regurgitating training data and aids in the development of countermeasures. Many prior works—and some recently deployed defenses—focus on "verbatim memorization", defined as a model generation that exactly matches a substring from the training set. We argue that verbatim memorization definitions are too restrictive and fail to capture more subtle forms of memorization. Specifically, we design and implement an efficient defense that *perfectly* prevents all verbatim memorization. And yet, we demonstrate that this "perfect" filter does not prevent the leakage of training data. Indeed, it is easily circumvented by plausible and minimally modified "style-transfer" prompts—and in some cases even the non-modified original prompts—to extract memorized information. We conclude by discussing potential alternative definitions and why defining memorization is a difficult yet crucial open question for neural language models.

## 1 Introduction

The ability of neural language models to memorize their training data has been studied extensively (Kandpal et al., 2022; Lee et al., 2021; Carlini et al., 2022; Zhang et al., 2021; Thakkar et al., 2021; Ramaswamy et al., 2020). When language models, especially ones used in production systems, are susceptible to *data extraction* attacks, it can lead to practical problems ranging from privacy risks to copyright concerns. For example, Carlini et al. (2021) showed that the GPT-2 language model could output personally identifying information of individuals contained in the training dataset.



Figure 1: Illustration of Memorization-free Decoding, a defense which can eliminate verbatim memorization in the generations from a large neural language model, but does not prevent approximate memorization.

One natural way to avoid this risk is to filter out any generations which copy long strings verbatim from the training set. GitHub's Copilot, a language-model-based code assistant, deploys this defense by giving users the option to "block suggestions matching public code" (GitHub, 2022).

In this work, we ask the question: "*Do language models emit paraphrased memorized content?*" This scenario can happen maliciously (e.g., adversaries trying to extract private user data) or through honest interactions (e.g., users prompting in real-world scenarios). Indeed, we find that Copilot's filtering system is easy to circumvent by applying plausible "style transfers" to the prompt. For example, by translating variable names from English to French the model outputs completely memorized examples, but post-processed with the en-fr style transfer. We further show that GPT-3 (Brown et al., 2020), a model trained on natural language, is also vulnerable to extraction attacks.

Unfortunately, Copilot's training set and precise algorithm for their defense are non-public. Therefore, to investigate this phenomenon systematically, we develop MEMFREE decoding (Figure 1), an efficient defense that is guaranteed to prevent all verbatim memorization, and which scales to training sets consisting of hundreds of gigabytes of text. In

---

*Remaining authors ordered by Algorithm 18 in Appendix H; briefly, we require Daphne be listed first, and Nicholas listed last, and we search for the first permutation of authors' first names which satisfies these constraints, where permutations order names by their salted MD5 hash.
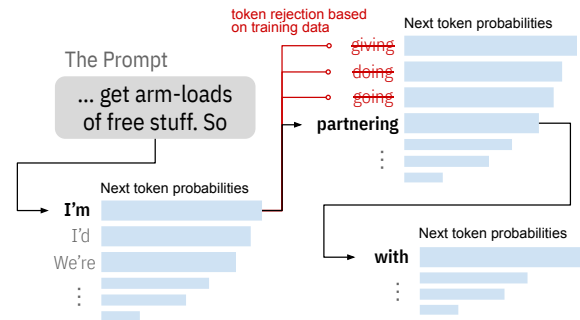
MEMFREE decoding, at each step of generation we check whether the model's chosen next token would create an $n$-gram found in the training set. If it does, an alternative next token is selected (without a computationally expensive regeneration) by sampling from the model's token posterior. The check for membership in the training set is performed efficiently using a Bloom filter containing all common $n$-grams from the training set.

We use MEMFREE to study Copilot's verbatim-filtering defense on other state-of-the-art large language models such as GPT-Neo (Gao et al., 2020). We first confirm that even honestly designed prompts often bypass verbatim memorization checks. Then, we observe another interesting phenomenon: language models succeed at emitting *approximate memorization* that bypass our filter all by themselves. Indeed, when prevented from generating exact $n$-grams from the training set, models are capable of "cheating" the filter by producing close paraphrases–for example, inserting spelling errors, adjusting punctuation or whitespace, or using synonyms (e.g., swapping 'and' with '&'). These changes lead to generated text a human would perceive as nearly identical, even if it is not verbatim memorization.

Clearly, defenses which prevent verbatim copying are *necessary but not sufficient* to protect against training data leakage. As a result of these failure modes, we argue that a broader definition of memorization is necessary when reasoning about training set memorization in language models. Such a definition should not only capture verbatim notions of memorization, but also notions based on high "semantic similarity" between model outputs and training data. We conclude our work by comparing approximate and verbatim memorization, discussing their relation to other domains of literature, and the challenges surrounding the ambiguity of approximate memorizations. Future work that aims to faithfully measure or prevent memorization in language models will need to take this ambiguity into account—for example, our analysis suggests that the fraction of datasets that large language models is likely far larger than the fraction as reported in prior work (Carlini et al., 2022).

## 2 Background

**Language Models.** We consider auto-regressive language models that operate over sequences of text and, given a prefix $p$, output a probability dis-

tribution for the next token in the sequence. To generate text for a prompt $p$, the language model starts with an empty suffix $s$, and repeatedly samples the next token from its prediction on $p + s$, and then appends this token to $s$. The success of neural language models has, in large part, been driven by the transformer architecture introduced of Vaswani et al. (2017), which allowed models to scale from millions to hundreds of billions of parameters over the past half-decade (Brown et al., 2020; Chowdhery et al., 2022; Zhang et al., 2022). This increase in model sizes has likewise driven increases in dataset sizes, with most of this data coming from internet crawls (Lee et al., 2021; Raffel et al., 2020; Gao et al., 2020).[1]

Prior work has shown that large language models can memorize and regurgitate potentially private information, like phone numbers and addresses, as well as memorize long sequences from their training sets (Carlini et al., 2019, 2021; Lee et al., 2021; Carlini et al., 2022; Zhang et al., 2021; Thakkar et al., 2021; Ramaswamy et al., 2020; Kandpal et al., 2022). Our work focuses on large language models trained to generate English text or code, and our work does not distinguish between problematic memorization (e.g. exposure of private information) and non-problematic memorization (e.g. quoting perfectly from a presidential speech).

**Measuring Memorization.** Many studies of memorization stem from a concern of privacy leakage: if a model memorizes sensitive training data and can generate it, then interactions with a model can lead to the leakage of that sensitive data. Nearly all of this literature is focused on measuring *verbatim cases of memorization*.

*Eidetic memorization* (Carlini et al., 2021) defines a string $s$ as memorized if there exists a prompt $p$ so that $\text{LM}(p) = s$ and $s$ is contained in the training dataset. This definition and variations of it have been used widely in the literature (Kandpal et al., 2022; Lee et al., 2021; Carlini et al., 2022). For example, Tirumala et al. (2022) study a similar per-token definition called *exact memorization* and Kandpal et al. (2022) a document-level definition called *perfect memorization*.

There is also a newly emerging line of works exploring *differential-privacy (DP)-based definitions* (Zhao et al., 2022; Stock et al., 2022), which relate to document-level DP guarantees in language

---

[1]A common source for datasets is the Common Crawl dataset found at: https://commoncrawl.org/

modelling (Yu et al., 2021). These works differ from the above in that they define a probabilistic leakage measure. However, this is based on the probability of generating—verbatim—a canary sentence $s$, depending on whether $s$ was contained in the training set or not. There are different probabilistic definitions, also based on verbatim sequences, such as the *counterfactual memorization* proposed by Zhang et al. (2021).

In the domain of language model memorization, the most similar work to ours is Lee et al. (2021) who also argue for a more relaxed definition of memorization. Lee *et al.* say any model output for a prompt $p$ is memorized if it is within some chosen edit distance of the prompt's true continuation in the training set. As we will discuss, a small edit distance may not capture all forms of approximate memorization either—such as our examples of "style-transfer" applied to memorized content.

**Preventing Memorization.** Differentially private training, e.g., using DP stochastic gradient descent (Abadi et al., 2016), is the gold standard for training models which provably do not memorize individual training examples. However, in practice, these techniques result in worse generative models (Anil et al., 2021)—thus, no state-of-the-art, large, language models are trained with DP. Instead, data deduplication has arisen as a pragmatic countermeasure against data memorization (Lee et al., 2021; Kandpal et al., 2022; Carlini et al., 2022). The core idea is to remove any duplicated content—e.g., repeated documents—because duplicated content is much more likely to be memorized. However, deduplication does not *guarantee* that a model will not still memorize individual (deduplicated) examples, necessitating defenses that operate at inference-time.

## 3 Preventing Models from Emitting Verbatim Training Data

In this paper, we consider inference-time defenses that eliminate the generation of memorized content from the training set. The most immediate way to do this is simply to filter all model outputs using some fixed definition of memorization. For example, in Carlini et al. (2022), a continuation $s = \text{LM}(p)$ of a $k$-length prompt $p$ is said to be memorized if the string $s$ exists verbatim in the training dataset. A straightforward implementation checks each generation $s$ against the training set and rejects any matches. We call the approach of

re-running a language model, possibly many times with different seeds, until a qualifying generation is produced, **retroactive censoring**.

The problem with retroactive censoring is that it effectively prevents the model from emitting any output when the model's confidence in a memorized string is too high. To encourage a model to generate novel outputs, we could also adopt a more granular filtering approach: rather than censoring memorized content solely at the level of an entire sequence $s$, we could instead check and mark each $n$-gram within $s$ individually. Filtering for memorization at the $n$-gram-level rather than at the sequence level allows substrings of a generation which may be novel to be kept, and only the pieces that are verbatim memorized to be modified. We call this approach **MEMFREE decoding**, as the defense is applied at decoding time.

Both retroactive censoring and MEMFREE decoding explicitly prohibit the model from emitting a sequence if it is contained (entirely or partially) in the training dataset. However, in retroactive censoring, if a generation starts off with memorized text, but then veers off track from the true continuation (a common occurrence), this would not be marked as memorization, even though a portion of the output sequence is clearly memorized. The MEMFREE decoding approach performs a more fine-grained and aggressive check by filtering out all memorized subsequences of a given length. In this work we use the MEMFREE decoding approach to show that even when a model is restricted from emitting any output with snippets of verbatim memorization, the model can still leak training data.

### 3.1 MEMFREE Decoding Details

In order to implement MEMFREE decoding, we alter the model's generation in an online manner by restricting the production of tokens which would result in an $n$-gram memorization. Let $p$ be the current working prefix and $t$ be the next proposed token when running the model forward.

Our algorithm first checks if any $n$-gram in the concatenated sequence $p\|t$ is contained in the training dataset $D$. If it is, we suppress this generated token and re-sample from the model. To avoid potentially expensive resamplings, we equivalently express this as altering the model's output probability distribution by removing the probability mass from token $t$. In this way, we guarantee that prior to sampling the probability of outputting a mem-

30

orization will be 0. Appendix B.1 gives a formal procedure for this method.

Altering the token posterior allows any sampling strategy to be used on top of memorization-free decoding. For example, if one uses top-$k$ sampling, tokens that result in memorization are disqualified before the probability distribution is truncated to the $k$ next most likely tokens. This procedure is guaranteed to generate non-memorized text.

## 3.2 Querying the Training Set Efficiently

Our MEMFREE defense has assumed that it is easy to perform the query $s \in D$ to test if any given string is contained in the training dataset. Because the defense works at inference-time, it is necessary that this query is computationally efficient to maintain utility of the language model. Given that training sets may contain terabytes of data (Brown et al., 2020), it is infeasible to maintain an entire copy of the training dataset in an efficiently accessible storage. Thus, we explore three optimizations to speed up the process of memorization checking.

**First**, as a direct result of our $n$-gram memorization definition, we can equivalently check only the $n$-gram ending in the current predicted token $t$; we can thus avoid many $n$-gram queries for each token. Further, and in addition to preventing subsequence memorization, this allows us to avoid queries into a large set of all prefixes and continuations.

**Second**, we only check against sequences that have a reasonable probability of being memorized by the model. In theory, this could be easily determined by running each $n$-gram $s \in \mathcal{D}$ through the model and then filtering out all sequences with high loss (thus unlikely to be memorized). However, this is a computationally expensive procedure as it requires re-processing every substring of the training dataset. Instead, a computationally- and storage-efficient procedure could be to only store $n$-grams which occur more than once in the training set—prior work has shown duplicate text is the most likely to be memorized (Lee et al., 2021; Kandpal et al., 2022).

**Third**, by being willing to tolerate some false positives (labeling an $n$-gram as memorized when it is in fact not), we can take advantage of probabilistic data structures such as Bloom filters (Bloom, 1970), which admits no false negatives but trades off time and space with the false positive rate (which can be computed exactly). Thus, by using a Bloom Filter, we guarantee that no mem-

orized $n$-gram will ever be released (i.e., a false negative) but we may (rarely) prevent the emission of non-memorized content (i.e., a false positive).

Integrating a Bloom Filter into our defense is straightforward. Let $\mathcal{F}_{fp}(\mathcal{D}_n)$ represent the Bloom Filter of dataset $\mathcal{D}$, generated by adding each $n$-gram of the dataset $s \in \mathcal{D}_n$ to the Bloom filter, with false positive rate $fp$. Then, any memorization check $s \in \mathcal{D}_n$ in Algorithm 1 can be replaced with $s \in \mathcal{F}_{fp}(\mathcal{D}_n)$. The Bloom filter can be generated with a single pass over the model's training set, which could be performed in parallel with one epoch of model training.

**Additional Parameters.** We must choose an appropriate false positive rate based on memory constraints and the chosen $n$-gram length. Choosing $n$ has two major impacts: on the population size (i.e., the number of unique $n$-grams) and thus the size of the filter, and on the effectiveness of memorization mitigation. If $n$ is set too low, then we will certainly prevent all memorized sequences but might also prevent too many common phrases. But if we set $n$ too high, we might not prevent actually memorized sequences from being emitted by the model. We discuss these tensions in Appendix B, along with two additional takeaways: (1) that MEMFREE *does not impact downstream model performance* (which may result from false positives), and (2) that our chosen optimizations maintain a suitably low false negative rate (we observed a 3000x improvement). These optimizations led to a filter of size 1.6 gigabytes (or, 40.5 gigabytes if all, even non-duplicated, 10-grams were stored) when run over the 800GB Pile dataset.

## 3.3 Measuring Approximate Memorization

To show that defenses against verbatim memorization still allow approximate memorization, we need a definition for approximate memorization. We consider two definitions. First, drawing from standard NLP evaluation techniques, we measure the BLEU score between the generated and ground-truth continuations. Second, we measure the length-normalized character-level Levenshtein similarity between the generated and ground-truth continuations. Appendix C.1 gives implementation details. In Section 5, we investigates how these two similarity metrics decrease with MEMFREE decoding.

For situations requiring a binary label of whether approximate memorization has occurred, we use the following definition: a suffix $s$ for prefix $p$ is

**Standard prompting with original prefix and format**

```
float Q_rsqrt( float number )
{
  long i;
  float x2, y;
  const float threehalfs = 1.5F;

  x2 = number * 0.5F;
  y  = number;
  i  = * ( long * ) &y;
Copilot no longer generates continuations
```

**Prompt with Python-style comment**

```
# float Q_rsqrt( float number )
# {
#   long i;
#   float x2, y;
#   const float threehalfs = 1.5F;
#
#   x2 = number * 0.5F;
#   y  = number;
#   i  = * ( long * ) &y;
#   i  = 0x5f3759df - ( i >> 1 );
#   y  = * ( float * ) &i;
#   y  = y * (threehalfs - (x2*y*y));
#
#   return y;
# }
```

**Prompt with French translation (alternate naming convention)**

```
float Q_sqrt( float nombre )
{
  long i;
  float x2, y;
  const float trois_moitie = 1.5F;

  x2 = nombre * 0.5F;
  y  = nombre;
  i  = * ( long * ) &y;
  i  = 0x5f3759df - ( i >> 1 )
  y  = * ( float * ) &i;
  y  = y * ( trois_moitie - (x2*y*y));
  //y = y * ( trois_moitie - (x2*y*y));

  return nombre * y;
}
```

Figure 2: **Honest "style-transfer" prompts evade verbatim memorization filters.** Trivially modifying prompts causes GitHub's Copilot language model to emit memorized, but not verbatim, content. Prompts highlighted in blue. Model evaluated with the option "block suggestions matching public code" enabled. For brevity, we removed comments from model outputs.

labeled as memorized if for generation $g = f(p)$, $\text{BLEU}(g, s) > 0.75$. This threshold was chosen by qualitatively inspecting examples. Several example generations that are close to this threshold are shown in Table A12.

When we repeat the prefix-extraction experiment from (Carlini et al., 2022) to measure incidents of generations that could be considered memorized, but using this approximate definition instead of a verbatim one, we find that hat prior literature has significantly underestimated memorization leakage. In Figure 3, the shaded region represents the fraction of memorized samples that would have bypassed a verbatim memorization filter: in the worst case, there is a factor-of-two increase.



Figure 3: **Significantly more examples are approximately memorized (BLEU > 0.75) than are found to be exactly memorized by Carlini et al. (2022).** This is for undefended generation.

However, we caution that this definition of approximate memorization is inaccurate, potentially both over and under counting approximate memorization. While our choice of a $0.75$ BLEU score threshold shows a significant increase in approximate vs. verbatim memorization, it is not clear that all identified cases of memorization would be perceptually tagged as such by a human judge. This is one reason why simply switching to this definition for defenses may not be ideal—it could introduce significant false positives.

## 4 Evading Verbatim Memorization Defenses

In this section, we show how retroactive censoring of verbatim memorization can be evaded, even in settings where models are used honestly. We first present a case study with Copilot, which has implemented retroactive censoring in production. We then show how a large English language models like GPT-3 and PaLM are susceptible to the same vulnerability, should a defense similar to Copilot's be deployed. In short, protecting against verbatim memorization can lead to a false sense of privacy.

### 4.1 Evading Copilot's Memorization Filter

Copilot is a code auto-complete service which is trained on GitHub code. Copilot is built using the Codex language model designed by OpenAI (Chen et al., 2021). To prevent generating memorized code, Copilot uses a filtering mechanism that blocks model outputs from being suggested if they overlap significantly (approximately 150 characters) with a training example. This is a practical example of a filter that aims at preventing perfect verbatim memorization, presumably by using a procedure similar to Algorithm 1 (the

exact mechanism used by GitHub is not public). However, we find that the filter fails to prevent the leakage of training data in many settings.

**Style-transfer prompting.** In Figure 2, we show that Copilot's filter can easily be bypassed by prompts that apply various forms of "style-transfer" to model outputs, thereby causing the model to produce memorized (but not verbatim) outputs.

As a concrete example, we demonstrate how to extract the public code for Quake's "Fast Inverse Square Root". If we naively prompt the model with the function definition "`float Q_rsqrt ( float number )`", Copilot correctly aborts generation of the full function ("standard prompting").

However, we find that simple style-transfers applied to the prompt allow us to easily bypass Copilot's restrictions. First, via prompting with "Python-style comments" we begin our prompt with Python's comment character "#". Even though this is syntactically invalid C code, Copilot outputs the entire verbatim fast inverse square root algorithm, but commented out. Second, in prompting with "French translations" we change the naming convention to French. As a result, the generations follow the new naming convention and are no longer flagged as a verbatim match. Other naming conventions, such as pre-pending "_" to the variable or changing the language to Spanish, also work.

These strategies work because the Copilot model is sufficiently powerful: it can both follow the style-transfer prompt (by e.g., renaming variables) while simultaneously regurgitating memorized training data. We provide more examples in Appendix F.

**Copilot evades its own filter.** Not only do *actively* style-transferred prompts evade the verbatim memorization filter, but even passively prompting Copilot with highly duplicated text from the Pile dataset can too. We find several examples where Copilot evades its own filter to output memorized text, some of which we show in Figure 5. We see that Copilot evades the filter by (1) changing capitalization, (2) making small non-stylistic errors, and (3) changing whitespaces. The latter evasion (changing whitespaces) is surprising, as Copilot's documentation reports ignoring whitespace in its filtering mechanism (Appendix A). However, we hypothesize that this can be explained by the model replacing tabs with space characters. We can verify this by adding tabs to the beginning of each line of the Q_sqrt function, as an application of our



Figure 4: Fraction of prompts which discover approximate memorization, grouped by domain (left) and by style transfer applied (right). We tested two versions of GPT-3 DaVinci and two sizes of PaLM. Full plot in Appendix D.

style-transfer strategy.

## 4.2 English Language Models

Following our analysis of Copilot, we ask whether this vulnerability is pervasive in other language models too. We use API access to four large (English) language models—GPT-3 Davinci Original and V2 and PaLM 62B and 540B–to test whether they would be susceptible to style transfer of the prompt. We assume that the training sets for these models are unknown and prompt with documents we believe are likely to have been memorized: open-source licenses, famous speeches and monologues, novel openings, and song lyrics. For each document, we prompt the model with 100 words of either (1) the original document ("base"), (2) the document with all spaces doubled ("spaces"), (3) the document in all lowercase ("lower"), and (4) the document in all uppercase ("caps"). We report approximate memorization results of this experiment in Figure 4, with additional figures in Appendix D.

We see that even when prompting with style-transfered prompts, GPT-3 and PaLM are still often able to generate memorized continuations. Defenses for verbatim memorization are therefore incomplete. Among the three techniques, uppercasing was the least likely to lead to memorized generations. For the two PaLM models, the larger one is much more capable of memorization than the smaller one, which validates prior work (Carlini et al., 2022). The two version GPT-3 are purportedly the same size model but have quite different memorization tendencies. For example, V2 is much

**1) Misspelling and changed capitalization**

```
This program is free software" ; you can redistribute redistribut it and / or *
modify it under the terms of the GNU General Public License
* asAS published by the Free Software Foundation ; either
version 2 * of the License , orOR ( at your option ) any later
version "
```

**2) Small non-stylistic errors**

```
@aws-sdk/protocol-http" ;
import { Command as $Command } from "@aws-sdk/smithy-client";
import {FinalizeHandlerArguments,Handler,HandlerExecutionContext,
MiddlewareStack,HttpHandlerOptions} as __HttpHandlerOptions, MetadataBearer
as __MetadataBearerMetadataBear ,
```

**3) Changed whitespace**

```
[...]IPV6_2292PKTINFO(2 ws)\t\t\t\t\t = 0x2\nIPV6_2292PKTOPTIONS \t\t\t\t = 0x6\n
IPV6_2292RTHDR [20 spaces][9 spaces]= 0x5 \n
```

Figure 5: **CoPilot can "cheat" and emit nearly verbatim memorized content.** Here, we show prompts from the training set, where the model makes slight errors causing the continuations to pass the filter. Prompts are in cyan, followed by CoPilot's continuation where errors are highlighted as **model's generation in orange** with the **correct characters in green**.

more susceptible to the "double spaces" style transfer than the Original Davinci. This emphasizes the importance of models' training set compositions and training methods on memorization tendencies.

# 5 MEMFREE Decoding Experiments

In this section, we study the effectiveness of our proposed MEMFREE decoding defense from Section 3.3, and the appropriateness of our proposed definition of approximate memorization.

## 5.1 Experimental Design

It is not possible to apply MEMFREE to the models from Section 4 since their training sets are non-public. Instead, we turn to the GPT-Neo langue model family (Black et al., 2021). These models are trained on the Pile, a publicly available 825GB dataset (Gao et al., 2020). We build a Bloom filter over all 10-grams occur 10 or more times.[2] In all experiments, we generate text using $\arg\max$ decoding as the sampling method. We investigate four model sizes: 125M–6B parameters.

We evaluate using substrings of the Pile released by Carlini et al. (2022). The dataset includes 30k strings of length 150 tokens taken from the training set. These are divided into 30 buckets of 1k strings, sampled such that the strings in bucket $i$ occur in

---

[2]Note that the choice of $n$=10 for the $n$-gram size is very conservative, and common phrases that happen to be composed of 10+ tokens will get filtered out by this check. We discuss why we chose these particular values in Appendix B.



(a) BLEU (word-level)    (b) Edit similarity (char-level)

Figure 6: **MEMFREE reduces similarity when the continuation would have been highly similar to the ground-truth**, and has little impact otherwise. For 5,000 prompts, we plot the similarity of the groundtruth continuation with the generation from MEMFREE (y-axis) and with the undefended generation (x-axis). Generations on the diagonal were not memorized.

the Pile between $2^{i/4}$ and $2^{(i+1)/4}$ times. For each string, we use the first 50 tokens as a prompt $p$ and generate a 50-token long continuation.

## 5.2 Reduction in Memorization

MEMFREE significantly reduces the similarity of generations to the groundtruth, compared to performing undefended generation (Figure 6). We also observe that when undefended generation already results in low similarity with the groundtruth, MEMFREE does not significantly alter the generations, as desired.

Previous work shows that increasing model size increases discoverable memorization (Carlini et al., 2022; Kandpal et al., 2022). We again find a clear trend that generations from larger models have, on average, a much higher similarity with the original continuation (Figure 8). Despite this, MEMFREE remains effective at all model sizes (BLEU remains near-flat around 0.6). Even when a sequence has many duplicates in the train set (a strong indicator of memorization), MEMFREE significantly decreases similarity with the groundtruth at all model sizes (Figure 7).

## 5.3 Failures in Preventing Memorization

A defense against memorization fails when it allows a sequence to be generated which a human would perceive as substantially copied from the true continuation—even if it is not verbatim memorized. This failure case can be seen as the points where the MEMFREE generation is still a close match to the ground-truth continuation (Figure 6). It occurs because the defense only adjusted a few

Figure 7: **MEMFREE decreases the BLEU score of generations more for highly duplicated examples**.



Figure 8: **MEMFREE remains effective at reducing similarity between the generated and groundtruth continuations even as models grow larger.**,

tokens (e.g., 1 after every sequence of 10). When looking at these examples, many, but not all, are lists of numbers. Some examples are included in Table A17. There is also a second failure-case: when a full (50 token) generation is made more similar with the ground-truth by MEMFREE (on 10-grams) than without. This may happen depending on the model's token posterior's after removing all tokens that fail the MEMFREE check. Almost all of these cases had a trivial increase in similarity. However, $0.16\%$ of samples had a similarity increase above $0.1$. We found qualitatively that many of these cases did have significant overlap with the true continuation.

## 6 Discussion

**Defining memorization in language models.** While verbatim definitions have helped discover significant memorization in large language models, they are insufficient to capture more subtle forms of memorization. Our work highlights two such situations: "style-transfer" prompting, where defenses for verbatim memorization can be actively subverted, and when models "cheat" by outputting similar, but not verbat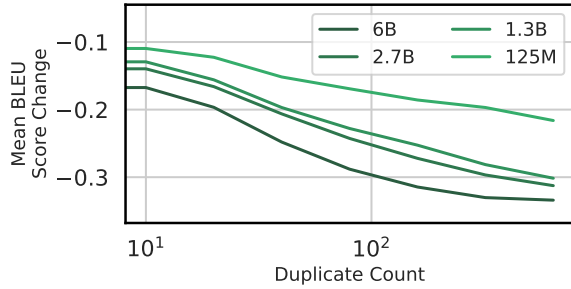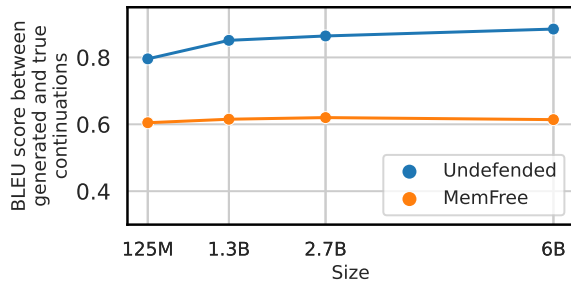im, continuations. As a result, our work suggests that memorization prevention must capture these types of paraphrased memo-

rizations in addition to the previously considered verbatim definitions. However, exhaustively anticipating styles to incorporate into defenses is an innumerable problem that will become harder as models become more powerful.

This emphasizes two major challenges in defining approximate memorization. First, since new approximate cases must be discoverable by the definition, this can result in a cat-and-mouse game. Second, the definition of memorization is domain-dependent. For example, our paper focuses on language models trained to output English and code, which each have different standards for what it means to memorize. Other languages will require different considerations when defining memorization.

There are a few areas of research which may help in improving memorization definitions. The field of **image generation memorization** is already comfortable with measuring *fuzzy* (in our terms, *approximate*) memorization, where generated items may be perceptually similar to training set examples, despite having high distance according to standard metrics. For example, Fredrikson et al. (2015) consider "model inversion", where an image is successfully recovered from the model if it is identifiable to a human worker. In Zhang et al. (2020), model inversion success is measured based on pixel similarity and feature space similarity to training images. These works also recover "representative" images from different classes, rather than specific training examples. Recent work on reconstructing training images have used feature similarity (Haim et al., 2022) and pixel similarity (Balle et al., 2022). In each of these papers, "fuzzy" reconstructions are allowed by the evaluation metrics and, indeed, are common in their reconstructions.

The inherent limitations of verbatim definitions of text regurgitation have also been well documented in the literature on **plagiarism detection**—both for text and code. Existing plagiarism tools, and their evaluations, go far beyond verbatim matches and consider fuzzy data "clones" ranging from simple transformations (e.g., word variations or shuffles) to arbitrary semantics-preserving paraphrasing (Roy et al., 2009; Potthast et al., 2010). Re-purposing techniques from the plagiarism detection literature to minimize generation of memorized data in LLMs is an interesting direction toward achieving better approximate memorization definitions in machine learning.

35

**Consequences for machine learning research.** In relaxing definitions of memorization, our paper acknowledges the blurred line between memorization (e.g., of personal information) and knowledge (e.g., of common facts). Because we use a 10-gram overlap, our MEMFREE decoding algorithm should not significantly impact utility, however studying this interplay is an important area of future work. However, still, identifying which data is considered "memorized" cannot be done only by looking for verbatim reproductions of the training set. This may make the task of understanding memorization and generalization more difficult.

We do not believe that our work requires abandoning all research directions which rely on prior verbatim definitions. These definitions are still useful as an efficient way to test for obvious and undeniable memorization. However it will be necessary to continue studying further relaxations of memorization definitions to adequately capture and measure the space of privacy concerns for language models.

## 7 Ethics & Broader Impact

Improving the privacy of neural language models—and especially those trained on user data—is an important and timely research problem. In this paper we hope to help both researchers and practitioners develop a more nuanced understanding of what constitutes memorization in language models. In particular, just because a sequence does not appear verbatim in a training dataset does not mean the example is a novel generation: as we have shown, models today are sufficiently powerful to minimally transform memorized data to make it appear superficially different even if the underlying content remains memorized.

Our observation will complicate the privacy evaluation of future machine learning models. It should no longer be deemed sufficient to check for (verbatim) matches between generated output and a training example. Practitioners in the future will need to be aware of this potential failure mode when applying output post-processing defenses to mitigate memorization. To the best of our knowledge, the only deployed system affected by our analysis is GitHub's Copilot. In order to mitigate harm here we shared a copy of our paper with the relevant researchers at both GitHub and OpenAI prior to paper submission.

In this paper we focus our efforts entirely on *public* datasets that other researchers have extensively studied (Gao et al., 2020) to minimize any harm caused by demonstrating extraction results. However, just because the data that we study is public does not mean there are no privacy concerns. As Brown et al. (2022) argue, there are many other considerations when discussing the privacy of large models trained on "public" datasets.

## Contributions

- Daphne Ippolito posed the idea of memory-free decoding using a bloom filter as a solution to memorization, worked on the MEMFREE implementation, ran experiments with GPT-3 and PaLM, and contributed to paper writing.
- Christopher Choquette analyzed how MEM-FREE used the bloom filter, created figures, and contributed to paper writing.
- Matthew Jagielski qualitatively analyzed generations from MEMFREE, created figures, and contributed to paper writing.
- Katherine Lee led figure-making, contributed to paper writing, and resolved TODOs.
- Milad Nasr ran experiments with Copilot and contributed to paper writing.
- Florian Tramèr came up with the idea of style transferring prompts and contributed to paper writing.
- Chiyuan Zhang generated figures, prepared qualitative examples, and contributed to paper writing.
- Nicholas Carlini identified the weaknesses in memory-free decoding, worked on the MEM-FREE implementation, ran experiments with GPT-Neo, and contributed to paper writing.

## References

Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. 2016. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, pages 308–318.

Rohan Anil, Badih Ghazi, Vineet Gupta, Ravi Kumar, and Pasin Manurangsi. 2021. Large-scale differentially private BERT. *arXiv preprint arXiv:2108.01624*.

Borja Balle, Giovanni Cherubin, and Jamie Hayes. 2022. Reconstructing training data with informed adversaries. *arXiv preprint arXiv:2201.04845*.

Michael A Bender, Martin Farach-Colton, Mayank Goswami, Rob Johnson, Samuel McCauley, and Shikha Singh. 2018. Bloom filters, adaptivity, and the dictionary problem. In *2018 IEEE 59th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 182–193. IEEE.

Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. 2023. Pythia: A suite for analyzing large language models across training and scaling. *arXiv preprint arXiv:2304.01373*.

Steven Bird, Ewan Klein, and Edward Loper. 2009. *Natural language processing with Python: analyzing text with the natural language toolkit*. " O'Reilly Media, Inc.".

Sid Black, Leo Gao, Phil Wang, Connor Leahy, and Stella Biderman. 2021. Gpt-neo: Large scale autoregressive language modeling with mesh-tensorflow. *If you use this software, please cite it using these metadata*, 58.

Burton H Bloom. 1970. Space/time trade-offs in hash coding with allowable errors. *Communications of the ACM*, 13(7):422–426.

Hannah Brown, Katherine Lee, Fatemehsadat Mireshghallah, Reza Shokri, and Florian Tramèr. 2022. What does it mean for a language model to preserve privacy? Seoul, Korean. ACM Conference on Fairness, Accountability, and Transparency.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Jehoshua Bruck, Jie Gao, and Anxiao Jiang. 2006. Weighted Bloom filter. In *2006 IEEE International Symposium on Information Theory*, pages 2304–2308. IEEE.

Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. 2022. Quantifying memorization across neural language models. *arXiv preprint arXiv:2202.07646*.

Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. 2019. The secret sharer: Evaluating and testing unintended memorization in neural networks. In *28th USENIX Security Symposium (USENIX Security 19)*, pages 267–284.

Nicholas Carlini, Florian Tramer, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. 2021. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2633–2650.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*.

Ondřej Dušek, David M Howcroft, and Verena Rieser. 2019. Semantic noise matters for neural natural language generation. In *Proceedings of the 12th International Conference on Natural Language Generation*, pages 421–426.

Ondrej Dušek and Filip Jurcıcek. 2016. A context-aware natural language generation dataset for dialogue systems. In *Workshop on collecting and generating resources for chatbots and conversational agents-development and evaluation*, pages 6–9.

Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. 2015. Model inversion attacks that exploit confidence information and basic countermeasures. In *Proceedings of the 22nd ACM SIGSAC conference on computer and communications security*, pages 1322–1333.

Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. 2020. The Pile: An 800GB dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*.

Sebastian Gehrmann, Tosin Adewumi, Karmanya Aggarwal, Pawan Sasanka Ammanamanchi, Anuoluwapo Aremu, Antoine Bosselut, Khyathi Raghavi Chandu, Miruna Clinciu, Dipanjan Das, Kaustubh Dhole, et al. 2021. The gem benchmark: Natural language generation, its evaluation and metrics. In *Proceedings of the 1st Workshop on Natural Language Generation, Evaluation, and Metrics (GEM 2021)*, pages 96–120.

GitHub. 2022. About GitHub Copilot. https://docs.github.com/en/copilot/overview-of-github-copilot/about-github-copilot.

Niv Haim, Gal Vardi, Gilad Yehudai, Ohad Shamir, and Michal Irani. 2022. Reconstructing training data from trained neural networks. *arXiv preprint arXiv:2206.07758*.

Nikhil Kandpal, Eric Wallace, and Colin Raffel. 2022. Deduplicating training data mitigates privacy risks in language models. *arXiv preprint arXiv:2202.06539*.

Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. 2021. Deduplicating training

data makes language models better. *arXiv preprint arXiv:2107.06499*.

Linyong Nan, Dragomir Radev, Rui Zhang, Amrit Rau, Abhinand Sivaprasad, Chiachun Hsieh, Xiangru Tang, Aadit Vyas, Neha Verma, Pranav Krishna, et al. 2021. Dart: Open-domain structured data record to text generation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 432–447.

Ankur Parikh, Xuezhi Wang, Sebastian Gehrmann, Manaal Faruqui, Bhuwan Dhingra, Diyi Yang, and Dipanjan Das. 2020. Totto: A controlled table-to-text generation dataset. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1173–1186.

Martin Potthast, Benno Stein, Alberto Barrón-Cedeño, and Paolo Rosso. 2010. An evaluation framework for plagiarism detection. In *Coling 2010: Posters*, pages 997–1005.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J Liu, et al. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(140):1–67.

Swaroop Ramaswamy, Om Thakkar, Rajiv Mathews, Galen Andrew, H Brendan McMahan, and Françoise Beaufays. 2020. Training production language models without memorizing user data. *arXiv preprint arXiv:2009.10031*.

Chanchal K Roy, James R Cordy, and Rainer Koschke. 2009. Comparison and evaluation of code clone detection techniques and tools: A qualitative approach. *Science of computer programming*, 74(7):470–495.

Thomas Scialom, Paul-Alexis Dray, Sylvain Lamprier, Benjamin Piwowarski, and Jacopo Staiano. 2020. Mlsum: The multilingual summarization corpus. In *2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8051–8067. Association for Computational Linguistics.

Pierre Stock, Igor Shilov, Ilya Mironov, and Alexandre Sablayrolles. 2022. Defending against reconstruction attacks with Rényi differential privacy. *arXiv preprint arXiv:2202.07623*.

Sasu Tarkoma, Christian Esteve Rothenberg, and Eemil Lagerspetz. 2011. Theory and practice of bloom filters for distributed systems. *IEEE Communications Surveys & Tutorials*, 14(1):131–155.

Om Dipakbhai Thakkar, Swaroop Ramaswamy, Rajiv Mathews, and Francoise Beaufays. 2021. Understanding unintended memorization in language models under federated learning. In *Proceedings of the Third Workshop on Privacy in Natural Language Processing*, pages 1–10.

Kushal Tirumala, Aram H Markosyan, Luke Zettlemoyer, and Armen Aghajanyan. 2022. Memorization without overfitting: Analyzing the training dynamics of large language models. *arXiv preprint arXiv:2205.10770*.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Da Yu, Saurabh Naik, Arturs Backurs, Sivakanth Gopi, Huseyin A Inan, Gautam Kamath, Janardhan Kulkarni, Yin Tat Lee, Andre Manoel, Lukas Wutschitz, et al. 2021. Differentially private fine-tuning of language models. *arXiv preprint arXiv:2110.06500*.

Chiyuan Zhang, Daphne Ippolito, Katherine Lee, Matthew Jagielski, Florian Tramèr, and Nicholas Carlini. 2021. Counterfactual memorization in neural language models. *arXiv preprint arXiv:2112.12938*.

Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*.

Yuheng Zhang, Ruoxi Jia, Hengzhi Pei, Wenxiao Wang, Bo Li, and Dawn Song. 2020. The secret revealer: Generative model-inversion attacks against deep neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 253–261.

Xuandong Zhao, Lei Li, and Yu-Xiang Wang. 2022. Provably confidential language modelling. *arXiv preprint arXiv:2205.01863*.

## A  GitHub Copilot

At the time of this paper's writing, GitHub Copilot's memorization prevention mechanism is described in their FAQ at https://github.com/features/copilot. We copy the text here:

> **What can I do to reduce GitHub Copilot's suggestion of code that matches public code?**
>
> We built a filter to help detect and suppress the rare instances where a GitHub Copilot suggestion contains code that matches public code on GitHub. You have the choice to turn that filter on or off during setup. With the filter on, GitHub Copilot checks code suggestions with its surrounding code for matches or near matches (ignoring whitespace) against public code on GitHub of about 150 characters. If there is a match, the suggestion will not be shown to you. We plan on continuing to evolve this approach and welcome feedback and comment.

## B  Further Discussion of MEMFREE

### B.1  Formal Procedure

Algorithm 1 provides a formal procedure for MEMFREE decoding. In all our experiments, we used $\arg\max$ decoding as the sampling method for line 4.

---
**Algorithm 1** MEMFREE decoding algorithm.

---
1: **procedure** GREEDY MEMFREE DECODING(language model $f$, prefix $p$, gen length $n$, training set $D$)
2:     **repeat**
3:         logits $\leftarrow f(p) - \infty \cdot \{\mathbb{1}[(p||t) \in D] : t \in \text{vocab}\}$
4:         tok $\leftarrow$ sample from logits
5:         $p \leftarrow p||\text{tok}$
6:     **until** $n$ iterations
7: **end procedure**

---

### B.2  Choice of $n$-gram length

There are two tradeoffs to consider when choosing an $n$-gram length: the choice of $n$ changes the granularity of the memorization checking and the total number of substrings of the dataset that must be stored in the Bloom filter. with respect to the former, notice that short $n$-grams do not have sufficient novelty (loosely, entropy) to be considered memorizations, e.g., they are often commons words and phrases. However, too large also would not capture shorter sequences that have sufficient novelty. On the latter, notice that the universe of possible $n$-grams is exponential in $n$, but that the unique number of such sequences in a fixed dataset may decrease with large $n$. This total number of unique sequences impacts the required size of the Bloom filter to maintain a fixed false positive rate. With $N$ the number of unique $n$-grams and $fp$ a decimal probability of the false positive rate, the size of the filter in bits is:

$$m = \left\lceil \frac{-(N * \log(fp))}{\log(2)^2} \right\rceil.$$

Then, $k$ the number of Bloom hash functions can be calculated from the number of bits per element, i.e., $m/N$, as:

$$k = \lceil ((m/N) * log(2)) \rceil.$$

This determines the cost of inserting and looking up into the Bloom filter as $\mathcal{O}(k)$. But, because $k$ typically remains small (in our case, $k = 7$), this can be treated as a small constant-time operation. See Tarkoma et al. (2011) for the full calculations, which the ones listed here are taken from.

We err on the side of caution and select $n$=10 for our experiments. This does prevent the model from generating common words or phrases which consist of 10 or more tokens, such as "The quick brown fox

jumped over the lazy dog." or "supercalifragilisticexpialidocious". We find qualitatively that the impact of this is low, and that this also presents a balanced trade-off with the Bloom filter size.

## B.3 Choice of Minimum Frequency

Ideally, we want $n$ large enough so that we do not prevent common phrases and small enough so that we catch all (though practically, most) possible memorizations. Optimizing $n$ for this task is both non-trivial, as the objective is not clear, and computationally expensive. Instead, we choose $n = 10$ based on qualitative experience that this does not prevent many common phrases. Further, we do so to also limit the storage cost of the Bloom filter, because $n$ too large leads to a blow up in the number of elements, $N$.

It is important to note that using MEMFREE with a lower $n$ will result in worse performance on standard benchmarks than using it with a higher $n$. This is because a lower $n$ means more true answers prevents from being generated.



Figure 9: **(left) Most generations have few Bloom queries,** as observed by the small quartiles; however, there is a long tail of few generations with many Bloom hits (12.6% of generations had beyond 50 hits with a max of 1111). **(right) Some positions had significantly more hits**, e.g., the first and tenth tokens. **(both)** are histograms from 6000 generations of 50 tokens each using MEMFREE decoding on GPT-Neo 6B.

## B.4 Python Implementation

Figure 10 contains a Python implementation of MEMFREE using the HuggingFace Transformers[3] API.

## B.5 Impact of MEMFREE on Downstream Task Performance

In this section, we discuss the worst-case impact MEMFREE could have on performance on downstream tasks. We measure this by looking at the targets, the groundtruth text a model's outputs are compared against, for three abstractive summarization tasks, three question answers tasks, and the 12 tasks in the GEM natural language generation benchmark (Gehrmann et al., 2021). On all these tasks, a model would score perfectly on the validation set if it exactly outputted the groundtruth target sequence. By measuring how many of the 10-grams in each of these target sequences are present in the bloom filter used by MEMFREE, we can assess the worst-case impact MEMFREE would have on model performance at these tasks. The results of this analysis are shown in table 1

We see that for most of these tasks, the percentage of 10-grams which are present in the bloom filter is not too much above 1%, the false positive rate of our bloom filter. Tasks where the target sequences come from documents likely to be present in the Pile are the most affected by MEMFREE usage. For example, for the BillSum and Arxiv summarization tasks, over 86% of their validation set examples have a 10-gram in the bloom filter. Non-English tasks, which are labeled with an asterisk in Table 1 were also significantly affected. The drop in performance for non-English tasks is due to the fact that GPT-Neo's vocabulary is built off of English. This means that non-English phrases end up being broken into many more tokens on average than English ones, and a single common word in a non-English language might take up several tokens. This can be seen in the bloom hit examples for the MLSum-de task.

---

[3] https://github.com/huggingface/transformers

```
banned = None

model = ## huggingface model loader here
bloom = ## set-like bloom filter

num_tokens_in_filter = 10

def ban_bloom(input_ids, scores):
    """input_ids is the tokens of the prompt. scores is the logits outputted by the model given these input_ids."""
    input_ids = input_ids.cpu().detach().numpy()

    # Order the tokens by their likelihood.
    order = torch.argsort(-scores, 1)
    order = order.cpu().detach().numpy()

    batch_size = input_ids.shape[0]

    # Set the likelihood to 0 for all the most likely next tokens which would create an ngram in the bloom filter.
    for ex in range(batch_size):
        for i in order[ex]:
            sequence_to_check = (input_ids[ex].tolist() + [int(i)])
            if sequence_to_check[-num_tokens_in_filter:] in bloom:
                scores[ex,i] -= 1000
            else:
                break
    return scores

prior_processor = model._get_logits_processor
def fn(*args, **kwargs):
    prior = prior_processor(*args, **kwargs)
    prior.append(ban_bloom)
    return prior

model._get_logits_processor = fn

# Proceed with calling model.generate as normal.
```

Figure 10: Implementation of MemFree in HuggingFace

There are easy strategies to reduce the effect MEMFREE has on benchmark performance. First, one could deliberately choose to omit from the bloom filter datasets which one decides are acceptable to memorize from, such as Wikipedia and legal documents. Second, one could increase the $n$-gram size of the bloom filter. As shown in the qualitative examples in Table 1, $n$=10 is perhaps too stringent for fact-based task, where names of proper nouns can take up 10-tokens or more. Third, one could reduce the error rate of the bloom filter so as to emit fewer false positives.

## B.6   Performance of MEMFREE

In this section, we study two questions: (1) "does MEMFREE maintain model utility?" and (2) "does our optimized MEMFREE prevent memorization release".

Along question (1), recall that MEMFREE can admit false positives, which may degrade the utility of the language model. Fortunately, the false positive rate can be computed exactly, e.g., see Tarkoma et al. (2011), and a long literature has proposed optimizations to account for non-uniform distributions (Bruck et al., 2006) and to adaptively correct for false positives (Bender et al., 2018).

Here, we study how, under reasonable computational constraints and inference times, the observed rates impact model utility. As we will show, we observe that MEMFREE maintains the highest utility (no observable impact) while being the most efficient defense.

Along question (2), we study if our optimizations lead to a substantial increase in the false negative rate. To do this, we repeat the experiment from (Carlini et al., 2022), which prompted GPT-Neo models with examples from its training data. We compute how many examples are verbatim memorized when MEMFREE decoding is used. The 6B parameter GPT-Neo model memorizes more than 12,000 of these documents, but, after applying MEMFREE, it only outputs 4 verbatim memorizations. These 4 remaining verbatim memorizations are repeated fewer than 10 times in the training data, and so were not added to our Bloom filter. Nonetheless, this strategy reduced verbatim memorization by over $3000\times$.

## B.7   Bloom Filter Statistics

Figure 11 shows the distribution in number of tokens (out of 50 generated) that were changed by MEMFREE from the token that would have been generated using undefended greedy decoding.

41

| Task | % ex with len>10 | % ex with bloom hit | % 10grams with bloom hit | Example 10-grams with bloom hit |
|---|---|---|---|---|
| **Summarization Tasks** | | | | |
| TIFU | 92.0 | 16.9 | 1.3 | stall windows, get new mouse, keyboard and cup • my freezer and now my home is the bog of • went to a concert five hours away as the dd |
| Arxiv | 100.0 | 86.8 | 1.38 | of a bose gas below the critical temperature. • in this paper, we develop a structure - preserving • consider a model of diffusion where the individuals behavior is |
| Pubmed | 100.0 | 92.3 | 1.7 | normal alanine aminotransferase • the prevalence of osteoporosis in postmen • www.cs.tau.ac.il |
| BillSum | 100.0 | 88.6 | 3.0 | Employee Retirement Income Security Act of 1974 and the Internal • Congressional Budget and Impoundment Control Act of 1974 • Federal Meat Inspection Act, the Poultry Products Inspection |
| **Question-Answering Tasks** | | | | |
| SQuAD2.0 | 9.8 | 1.1 | 5.9 | E. Mann, Raymond S. Bradley and Malcolm • CTLs (cytotoxic T lymph • in 1975. It went public in 1979 and was |
| WebQuestions | 2.4 | 0.9 | 9.8 | Academia de Bellas Artes de San Fernando • Paris Saint-Germain F.C. • The Mating Habits of the Earthbound Human |
| CoQA | 4.0 | 0.5 | 10.6 | Kingdom of Serbs, Croats and Sloven • Sheikh Mohammed bin Rashid Al Maktou • grabbed the rest of the pickle and ran |
| **GEM Benchmark** | | | | |
| CommonGen | 81.9 | 5.7 | 1.4 | You ride the horse around the area near the fence • children walk with their dog on a leash down the] • she wears a helmet & sits on the motorcycle. |
| Chezch Restaurant* (Dušek et al., 2019) | 99.6 | 23.5 | 1.7 | jemnou restauraci BarBar, kter • jsou v různých • Bohužel, poblí |
| DART (Nan et al., 2021) | 97.1 | 20.1 | 1.7 | in New York City. He was a member of • a low-priced family restaurant located near Raja • a Member of the U.S. House of |
| E2E clean (Dušek and Jurcıcek, 2016) | 99.9 | 88.9 | 1.0 | near Rainbow Vegetarian Café in the city center. • Phoenix is a cheap French restaurant in riverside. • a French restaurant with a moderate price range, but |
| MLSum-de* (Scialom et al., 2020) | 100.0 | 58.7 | 2.58 | zum neuen Vorsitzenden • für verfassungswidrig. • längst überfäll |
| MLSum-es* (Scialom et al., 2020) | 100.0 | 42.3 | 2.2 | del pacto y no de la confrontación • selección española de f • investigación sobre la desaparici |
| Schema-Guided Dialog | 63.3 | 7.5 | 1.3 | The Lord of the Rings: The Return of the • tyard By Marriott Sacramento Cal Expo has a 3 star • with Southwest Airlines. The flight takes off at 7 |
| ToTTo (Parikh et al., 2020) | 98.0 | 20.9 | 3.2 | and was broadcast on Venevisión. • As of the census of 2000, there were 133 • on the U.S. Billboard 200 chart. |
| XSum | 99.4 | 18.5 | 1.6 | stressed will not increase your risk of dying, according • Two drug dealing brothers taken back to court for mocking • the Institute of Directors (IoD) has |
| WebNLG-en | 97.9 | 27.4 | 4.7 | written by J.R.R. Tolkien, • play in the Campeonato Brasileiro • is affiliated with Visvesvaraya Technological University |
| WebNLG-ru* | 100.0 | 99.6 | 42.9 | • • |
| WikiAuto + Turk/ASSET | 96.5 | 16.7 | 2.2 | pop-punk, surf rock, ska, • was discovered by a team of astronomers from the University • cover of Sgt. Pepper's Lonely Hearts Club Band |

Table 1: **Some benchmark tasks could be significantly affected by MEMFREE.** For several standard benchmark tasks commonly used to evaluate language models, we report the percentage of test set target sequences which consist of at least one 10-gram (meaning hitting the bloom filter is possible), the percentage of test set target sequences which contain at least one 10-gram present in the bloom filter, and the percentage of all the 10-grams in the test set targets which can be found in the bloom filter. We also show 3 example 10-grams (delineated by '•') which are present in both the test set and the bloom filter. (For QA tasks, we only consider the first answer for each question.) **The numbers here reflect the worst case scenario: the fraction of examples a language model that perfectly memorized the test set would be incapable of getting exactly correct when used with MEMFREE.**

Figure 11: **Most generations require few ($< 5$) changes to pass MEMFREE checks.** Data for histogram from 6000, 50-token generations using MEMFREE decoding on GPT-Neo 6b.

Figure 9 presented some of the query patterns of the MEMFREE decoder to investigate when and how it impacts decoding. First, we observe that MEMFREE is trivial to run in terms of compute: it takes only 49.8 milliseconds to run 10,000 queries on one CPU core. From Figure 9 (left), all generations required significantly fewer queries (mean $= 42.1$ queries / generation)—even running batches of many hundreds or thousands of queries would incur less than a few seconds additional overhead. Second, we find that the Bloom filter is often hit at the first and tenth tokens after the prompt. We see many hits at the first token because all our prompts are from the training data—so there are relatively fewer single token additions that generate a novel $n$-gram. Third, we find that most generations need only a few ($< 5$) alterations due to MEMFREE decoding.

## C  More Details on Measuring Approximate Memorization

### C.1  Similarity Metrics Implementations

As noted in Section 5, we identify instances of approximate memorization by measuring the similarity between a generated continuation and the groundtruth continuation for a prompt. We do so using BLEU-score and character-level edit distance.

We computed BLEU score using NLTK's BLEU computation (`nltk.translate.bleu_score`) with the default parameters (averaging equally BLEU-1, BLEU-2, BLEU-3, and BLEU-4)(Bird et al., 2009). Edit distance was computed with the 'editdistance' pip package. Normalized edit similarity between two strings $x$ and $y$ is defined as:

$$\text{EDITSIM}(x, y) = \frac{\text{EDITDISTANCE}(x, y)}{\max(|x|, |y|)}$$

### C.2  BLEU Score Threshold Selection

We chose to use a BLEU score of 0.75 or higher to indicate that a generation substantially memorized from the ground-truth continuation. We choose to threshold BLEU score rather than edit-distance since it is more interpretable to NLP researchers. Table 12 presents several examples of generations that were very close to the 0.75 BLEU score threshold. This threshold was deliberately set high so as to avoid false-negatives, allowing us to provide a reliable lower bound on memorization. In our qualitative analysis of the data, all examples that were above BLEU of 0.75 showed unambiguous memorization. We also observed memorization in some examples with lower BLEU scores.

43

| True Continuation | LM Generation | BLEU |
|---|---|---|
| about a deeper bond with the product. Nostalgia. It's delicate, but potent. Sweetheart.\n\nTeddy told me that in Greek, "nostalgia" literally means "the pain from an old wound". It's a twinge in your heart far more | about a deeper bond with a product. Nostalgia. It's delicate, but potent. Teddy told me that in Greek, nostalgia literally means the pain from an old wound. It's a twinge in your heart, far more powerful than memory alone. This | 0.79 |
| And I'll show you\nYou want the moves like jagger\nI got the moves like jagger\nI got the moooooooves like jagger\nI don't even try to control you\nLook into my eyes and I'll own you\nYou with the moves | And I'll show you\nAll the moves like Jagger\nI've got the moves like Jagger\nI've got the moves like Jagger\nI don't need to try to control you\nLook into my eyes and I'll own you\nWith the | 0.79 |
| Purpose\n\nThe laws of most jurisdictions throughout the world automatically confer\nexclusive Copyright and Related Rights (defined below) upon the creator\nand subsequent owner(s) (each and all, an "owner") of an original work of\nauthorship and/or a database | Purpose\n\n The laws of most jurisdictions throughout the world automatically confer\n exclusive Copyright and Related Rights (defined below) upon the creator\n and subsequent owner(s) of an original work of authorship (the "Work").\n Certain jurisdictions do not recognize a | 0.76 |

Figure 12: Examples of model generations with BLEU score close to 0.75, the threshold we used to declare that approximate memorization had occurred.

# D   Experiments with Large English Language Models

## D.1   Prompt Selection Process

Famous speeches were selected from the "Top 100 Speeches" list found at https://www.americanrhetoric.com/newtop100speeches.htm. Monologues were selected from the list of two-minute monologues found at http://www.monologuedb.com/tag/2-minute-monologues/. Novels were selected from the Time Magazine's Top 100 All-Time Novels list found at https://www.goodreads.com/list/show/2681.Time_Magazine_s_All_Time_100_Novels. The opening paragraphs of the first chapter (skipping over prefaces, introductions, and boilerplate) were used as each example. The 2011 and 2021 song lyrics were selected from the Billboard Year-End Hot 100 singles lists found at https://en.wikipedia.org/wiki/Billboard_Year-End_Hot_100_singles_of_2011 and https://en.wikipedia.org/wiki/Billboard_Year-End_Hot_100_singles_of_2012.

For each document, the first 100 *words* were used as a prompt, and the first 50 generated *words* were compared with the first 50 words of the true continuation. This approach has the ramification that not all prompts were the same length in *tokens*. However, this approach was necessary for fairness across style transfers because an all-uppercased string is going to be many subword tokens longer than the lowercased version of the same string.

Figure 13: **"Style-transfer" prompting divulges approximate memorization in two versions of GPT-3 and two sizes of PaLM.** Note that generations also follow the same style as the prompt. Generations were characterized as memorized if they had a BLEU score of at least 0.75 with the ground-truth continuation.

| Domain with $n$ total prompts | Model | # Prompts Memorized per Style-Transfer Type | | | |
|---|---|---|---|---|---|
| | | Original | Two Spaces | Lower | Upper |
| Open-Source Licenses ($n$=42) | GPT-3 DaVinci Original | 23 | 8 | 14 | 1 |
| | GPT-3 DaVinci v2 | 30 | 25 | 25 | 13 |
| Famous Speeches ($n$=100) | GPT-3 DaVinci Original | 20 | 1 | 14 | 0 |
| | GPT-3 DaVinci v2 | 12 | 6 | 11 | 3 |
| Famous Monologues ($n$=60) | GPT-3 DaVinci Original | 3 | 0 | 1 | 0 |
| | GPT-3 DaVinci v2 | 4 | 3 | 4 | 0 |
| Novel Openings ($n$=77) | GPT-3 DaVinci Original | 9 | 0 | 3 | 0 |
| | GPT-3 DaVinci v2 | 7 | 4 | 5 | 0 |
| Lyrics 2011 ($n$=11) | GPT-3 DaVinci Original | 7 | 2 | 6 | 2 |
| | GPT-3 DaVinci v2 | 14 | 11 | 14 | 4 |
| Lyrics 2021 ($n$=11) | GPT-3 DaVinci Original | 3 | 3 | 3 | 2 |
| | GPT-3 DaVinci v2 | 4 | 2 | 4 | 4 |

Table 2: **"Style-transfer" prompting surfaces approximate memorization in GPT-3.** We explore $n$ prompts for each domain. Note that generations also follow the same style as the prompt.

# E   Experiments with MEMFREE and Other Model Families

In addition to running experiment using the GPT-Neo family, we also ran them with the Pyhia model family (Biderman et al., 2023). Like GPT-Neo, Pythia was trained on the Pile. There are two versions of Pythia, one trained on the same version of the Pile as GPT-Neo, and another trained on a deduplicated version of the Pile.

Figure 14 shows the amount of memorization in each of these three model families, with and without MEMFREE. Figure 15 shows the same scatter plots as in Figure 6, but using the 6.9B-parameter Pythia. We see that Pythia exhibits more approximate memorization than GPT-Neo. Though MEMFREE is still effective at reducing approximate memorization, it is slightly less effective than it was on GPT-Neo.



Figure 14: Approximate memorization on the base and deduped versions of Pythia, compared with GPT-Neo.



(a) BLEU (word-level)

(b) Edit similarity (char-level)

Figure 15: MEMFREE is also effective at reducing approximate memorization for the deduped Pythia 6.9B model.

# F   Qualitative Examples that Bypass Copilot's Filter

In Figure 16, we show more examples that bypass Copilot's memorization filter.

**Standard Prompting**

```
/* low -> Starting index, high -> Ending index
    */
quickSort(arr[], low, high)
{
    if (low < high)
    {
        /* pi is partitioning index, arr[p] is
            now
            at right place */
        pi = partition(arr, low, high);


Copilot no longer generates continuations
```

**Comment Prompting**

```
# /* low -> Starting index, high -> Ending index
#    */
# quickSort(arr[], low, high)
# {
#     if (low < high)
#     {
#         /* pi is partitioning index, arr[p] is
#    now
#            at right place */
#         pi = partition(arr, low, high);
#         quickSort(arr, low, pi - 1);   # Before
#    pi
#         quickSort(arr, pi + 1, high); # After pi
#     }
# }
```

**Naming Convention Prompting**

```
/* _low -> Starting index, _high -> Ending
    index */
quickSort(arr[], _low, _high)
{
    if (_low < _high)
    {
        /* pi is partitioning index, arr[p] is
            now
            at right place */
        pi = partition(arr, _low, _high);
        quick_sort(arr, _low, pi - 1);  //
            Before pi
        quick_sort(arr, pi + 1, _high); //
            After pi
    }
}
```

**Language Prompting**

```
/* depart -> index de départ, fin -> index de fin
    */
quickSort(arr[], depart, fin)
{
    if (depart < fin)
    {
        /* pi est l'index de partitionnement, arr[
            p] est maintenant
            à la bonne place */
        pi = partition(arr, depart, fin);
        // Trier les éléments séparément avant et
            après la partition
        quick_sort(arr, depart, pi - 1);
        quick_sort(arr, pi + 1, fin);
    }
}
```

**Naming Convention Prompting**

```
/* LOW -> Starting index, HIGH -> Ending
    index */
quickSort(arr[], LOW , HIGH )
{
    if (LOW  <  HIGH )
    {
        /* pi is partitioning index, arr[p]
            is now
            at right place */
        pi = partition(arr, LOW , HIGH );
        quickSort(arr, LOW , pi - 1);  //
            Before pi
        quickSort(arr, pi + 1, HIGH );  //
            After pi
    }
}
```

**Comment Prompting**

```
/* low -> Starting index, high -> Ending index */
    //1
quickSort(arr[], low, high) //2
{
    if (low < high) //3
    {
        /* pi is partitioning index, arr[p] is now
            at right place */
        pi = partition(arr, low, high); //4
        quickSort(arr, low, pi - 1); //5
        quickSort(arr, pi + 1, high); //6
    }
}
```

Figure 16: Preventing GitHub's Copilot language model from emitting verbatim memorized content often results in it emitting memorized, but not verbatim, content. Trivially modifying the prompt bypasses these output filters. Prompts are highlighted in blue. Model evaluated with the option "block suggestions matching public code" enabled. For brevity we removed comments from model outputs.

# G   Qualitative Examples of Memorization Despite MEMFREE

Table 17 gives examples of approximate memorization which occured despite using MEMFREE decoding. Below , we present a longer list of full examples of MEMFREE decoding. In the following paragraphs, the gray highlighted texts are the prompts, and the texts that follow the prompts are generated by the model. For easier reading, we merged the individual tokens to form text spans, except at tokens where bloom-filter rejection happens. In this case, the symbol ▷ indicate connection between tokens. Moreover, red highlighted texts indicate one or more rejected tokens, which are always followed by one accepted token highlighted with green. The texts are lightly edited mostly to properly encode unicode symbols (some garbled symbols are replaced as ⊠) that LATEX does not handle gracefully.

**Split word "Activity"**  viewHolder.swipeLayout.setOnDoubleClick  Listener(new  SwipeLayout.DoubleClickListener() { Override public void onDoubleClick(SwipeLayout layout, boolean surface) { Toast.makeText(mContext, "DoubleClick : " + position, ▷ Toast ▷ .LENGTH ▷ _ ). ▷ show(); ▷ } // }); if return m ▷ Context.start ▷ Activity Service Activ ▷ ities(new Intent(mContext, ▷ Activity ▷ .class)); ▷ } m ▷ Context.fin ▷ ish ishing ▷ ();

**Typo "Wildlife" ⇒ "Wildife"**  should promptly announce a phase-in of non-lead ammunition for all hunting throughout the rest of the state to protect other wildlife poisoned by lead and to safeguard human health. Hunters can play a critical role in the recovery of the condor and also keep lead from poisoning the wild game they eat." The lead-free ammunition regulations are designed to reduce lead poisonings of the iconic and extremely endangered California condor. Condors, eagles, and other scavengers such as ravens ▷ , crows, and hawks ▷ are have ▷ been poisoned by lead ammunition. The California Department of Fish and ▷ Wildlife Game Wild ▷ ife (CDFW) is the lead agency responsible for implementing the regulations. The regulations are based on the Condor Protection Act,

**Singular to plural: "claim" ⇒ "claims"**  ) No 1924/2006. neurotransmission and muscle contraction including heart muscle 29 Magnesium Magnesium contributes to normal protein synthesis The claim may be used only for food which is at least a source of magnesium as referred to in the claim SOURCE OF [NAME OF VITAMIN/S] AND/OR [NAME OF MINERAL/S] as listed in the Annex to Regulation (EC) No 1924/2006. protein ▷ synthesis 30 ▷ M N ▷ iacin Niac ▷ in ins ▷ contribute to normal protein synthesis The ▷ claim claims ▷ may be used only for food which is at ▷ least lest ▷ a source of niacin as referred to ▷ in to ▷ in the claim SOURCE OF [

**Uppercase to lowercase**  cuisine of the city's long-established Hum Vegetarian Cafe & Restaurant to a central garden location. Settle into the elegant and verdant space and enjoy dishes including papaya and banana flower salads, mushrooms steamed in coconut, and the subtle combination of braised tofu with star anise and cinnamon. Nha Hang NgonVietnamese$$ (MAP GOOGLE MAP); %028-3827 7131; 160 ⊠ Pasteur; meals 60 ▷ ,000-100,000đ; h7am-10pm; ▷ ⊠ ) ▷ This is the place to come for ▷ a authentic ▷ Vietnamese food in the French ▷ Quarter quarter ▷ . The menu is a mix of Vietnamese and French dishes, and the food is

**Lower case to upper case**  broker = createBroker(deleteAllMessagesOnStartup); broker.start(); } public BrokerService createBroker(boolean deleteAllMessagesOnStartup) throws Exception { return createBroker(deleteAllMessagesOnStartup, TRANSPORT_URI); } public BrokerService createBroker ▷ ( With ▷ TransportURI(boolean deleteAllMessagesOnStart ▷ up Up ▷ ) throws Exception { ▷ ⊠ return ▷ create new broker ▷ Service.createBroker(deleteAllMessages ▷ On , ▷ TRANSPORT_URI); } ▷ ⊠ ⊠ ▷ } ▷ ⊠ ⊠ <|endoftext|>

**Change from "agree" to 'Ag-reeableness"**  a person is imaginative or independent, high openness can be perceived as unpredictability or lack of focus. Moreover, individuals with high openness are said to pursue self-actualization specifically by seeking out intense, euphoric experiences, such as skydiving, living abroad, gambling, et cetera. Conversely, those with low openness seek to gain fulfillment through perseverance, some disagreement remains about how to interpret and contextualize the openness factor. A tendency to be organized and dependable, show self- ▷ discipline, and be goal-oriented is also associated

with high openness. The openness factor is also associated with the Big Five personality traits of ▷ agree Ag ▷ reeableness, Conscientiousness ▷ , and ▷ Emotional Stability. See also Openness

**Passive voice to active voice (grammar error)**   still wouldn't shock me at all if the Chiefs pulled off a Giants- or Ravens-esque Super Bowl run to cap off this five-year window of the Smith/Reid era with a ring. While Pittsburgh has been this team's bugaboo, maybe they can avoid them in January thanks to a Jacksonville upset, and I still say this is the AFC team most likely to win a playoff game in New England. While so many were quick to write the Chiefs off, they just opened up ▷ a new chapter in their history. The Chiefs are the AFC's best team, and they're going to be a force to be ▷ reckoned reckon ▷ with for years to come. 1. New England Patriots ▷ ⊠ The ▷ Patriots are the AFC's best

**Change of protocol (email still get generated)**   ="https://groups.google.com/group/django-developers" target="_blank" rel="nofollow" onmousedown="this.href=&#39;https://groups.google.com/group/django-developers&#39;;return true;" onclick="this.href=&#39;https://groups.google.com/group/django-developers&#39;;return true;">https://groups ▷ . ... ▷ </a> <a href=" ▷ https http mail ▷ to ▷ : :// ▷ david@davidwalsh.name" target="_blank" rel="n ▷ of ore ... ▷ </a> <a href="mailto://david@davidw

**Synonyms**   ken interior. The seats were heavily cushioned black velvet. On the windows, the Darkling's symbol had been cut into the glass: two overlapping circles, the sun in eclipse. Across from me, the two Grisha were studying me with open curiosity. Their red kefta were of the finest wool, embroidered lavishly in black and lined in black fur. The fair-haired Heartrender was lanky and had a long, melancholy face. Ivan was taller, broader, ▷ and had a face like a bulldog's. "You are ▷ the a ▷ Gr very ▷ pretty girl," Ivan said. "Thank you," I ▷ said replied answered ▷ . ▷ ⊠ " ▷ I'm not a girl." "You are a girl," he said.

**Synonyms**   severing any such bond. In re L.M., 923 A.2d 505, 511 (Pa. Super. 2007) (citing 23 Pa.C.S. § 2511) (some citations omitted). Section 2511(a) provides in pertinent part: (a) General rule.-The rights of a parent in regard to a child ▷ ▷ may are ▷ not terminated by a proceeding brought under ▷ this ▷ part chapter section sub subsection [ article paragraph ▷ or ▷ paragraph section ▷ 2512 or 2513(a) or (b), or any ▷ ⊠ _____ ▷ _____ ▷ ⊠ â ▷ *Retired Senior Judge

**Synonyms**   " "Do Androids Dream of Electric Sheep?" (the original of "Blade Runner"), and his masterpiece, "Ubik." Dick's fans are not modest in their claims. Nor are they especially precise: Borges, Calvino, Kafka, Robertson Davies are cited, in the blurbs and introductions, as his peers. A note of inconsistency inflects these claims-Calvino and Robertson Davies? ▷ -but the point is clear: Dick is the most important writer of the last century. The book is divided into three ▷ sections parts main categories ▷ : "Themes," "Themes and Themes," and "Themes

**Rejecting multiple candidates**   s den." Scott is aware of the impact his race and size has on the way people – particularly authority figures like law enforcement officers – perceive him. He is big. He is dark-skinned. "They look at us like we don't know how to control ourselves and we just get angry quick," he said. "It's not even like that. They criminalize us for no reason." Scott⊠ ▷ ⊠s mother, who is white, said she has been stopped by police for no reason. "I' ▷ ve m ll d ▷ be say get like never just have ask rather been tell pull ▷ over and they'd be like, ⊠ ▷ ⊠ ⊠ L ▷ What are you doing?'" she said

**Other examples**   Suzy is great! She helped me buy my condo at a great price (foreclosure) and then was super patient with my husband and me 4 years later when we were on a search for a house. She helped us get our... Suzie H., Jacksonville Goes above and beyond Suzy has helped me close on my third property in 3 years. First she found me my dream pool home at the Beach then she helped me find two investment townhomes in ▷ the same area ▷ . and ▷ now she is helping me find my dream home. She is always available to answer any questions I have and goes above and beyond to help me find the perfect home. I ▷ would highly ▷ recommend her! Suzy H., Jacksonville

from this new programme. I have also been reminded of the role of tax measures in supporting urban development. With us in the gallery today is Mr Vuyisa Qabaka, a Cape Town entrepreneur and co-founder of an organisation called the Good Neighbourhoods Foundation. His advice is that "Government should encourage township investment. For instance, it could promote urban development and regeneration

through accelerated depreciation allowances for new building constructions or refurbishment of existing buildings." ▷ I am sure that ▷ the many ▷ of you ▷ will have in ▷ this Chamber will agree with him. I am also sure that many of you will agree with the Minister of Finance, who has said that the tax system should be used to support the ▷ development growth economy ▷ and to create

m off on some details.) Unelma keltaisesta kuninkaasta. Fastaval is not your average convention – it specializes in incredibly tight auteur-designed roleplaying scenarios. A bunch of people run each scenario for players, not just the creator. There's awards for best scenarios in different categories. The Society for Nordic Roleplaying published a collection of these scenarios translated into Finnish a few years ago, called Unelma keltais ▷ esta kuninkaasta. It's a great book, ▷ and but with ▷ a lot of great scenarios. ▷ I The ▷ book is available in English, but it's not cheap. I've been looking for a copy for a while

disappoint Jimmy. Then, I slept like a baby. SoFortWorthIt Oscars Swag GIVEAWAY!!! The Oscars are exhausting, y'all. I'll definitely be cheering for all the stars this year, especially since I know the kind of caviar-Champagne-and-swag-filled night they're experiencing. And you know what? I want you to experience what it's like to get arm-loads of ▷ free stuff. So, I'm ▷ giving doing going partnering ▷ with the folks at the FortWorthIt Oscars Swag Giveaway to give away a $100 Visa gift card to one lucky winner. To enter, all you have ▷ to do ▷ is

decision." "It will go down to destruction... or else, it will survive." "This is their moment of trial." "They've got to show themselves worthy of everything we gods have given them." "But evil is dark and strong." "And it may be that the scales of fate... are not yet in full balance." "What can I do to equalize both sides of the struggle, Athena?" "If you don't want to increase the powers of all men... then why don ▷ 't you just give me the power to destroy them?" "I can't do that." ▷ " "[ ▷ Thunderclap]" "I'm sorry." "I'm ▷ sorry not so afraid ▷ I it you that the ▷ gods have decreed... that the balance of power must be maintained." "I'm

give him a minute between removing the first tray and replacing it with the second - and you can come up with all sorts of theories to explain your findings. You can even throw a person in an MRI machine, study the flickering images on your computer screen. But the brain is the ultimate black box. Eventually, to grasp the first cut, you'll have to make another. The car pulls into the parking lot of the nursing home, noses into an empty space. Annese and Cork ▷ y get out, and Annese goes to the trunk to get the wheelchair. Corky is still standing, leaning on the car. "I'm going ▷ to in ▷ ," ▷ he she An ▷ nese says. "I'll come ▷ with in ▷ with ▷ you ya

| True Continuation | MEMFREE | Undefended |
|---|---|---|
| attendee list reflects the minimum number of seats reserved for your game. Preregistered attendees for each of your events are listed below. If one attendee reserved multiple seats for your event, that attendee's name will only be listed once. Therefore, the attendee list reflects the minimum number of seats reserved for your game. Preregistered attendees for each of your events are listed below. If one attendee reserved multiple seats for your event, that attendee's | attendee list reflects the minimum number of seats reserved for your game. Preregistered attendees for each of your events are listed below. If one attendee reserved multiple seats for your event, that attendee's name will only be listed once. Therefore, the attendee list reflects the minimum number of seats reserved for your game. Preregistered attendees for each of your events are listed below. If one attendee reserved multiple seats for your event, that attendee's **BLEU: 1** | attendee list reflects the minimum number of seats reserved for your game. Preregistered attendees for each of your events are listed below. If one attendee reserved multiple seats for your event, that attendee's name will only be listed once. Therefore, the attendee list reflects the minimum number of seats reserved for your game. Preregistered attendees for each of your events are listed below. If one attendee reserved multiple seats for your event, that attendee's **BLEU: 1** |
| 293. 294. 295. 296. 297. 298. 299. 300. 301. 302. 303. 304. 305. 306. 307. 308. 309. 310. 311. 312. | 293. 294. 295. 296. 297. 298. 299. 300. 301. 302. 3034. 3043. 305. 3064. 3076. 308. 3097. 31009. 31108. 31210. **BLEU: 0.95** | 293. 294. 295. 296. 297. 298. 299. 300. 301. 302. 303. 304. 305. 306. 307. 308. 309. 310. 311. 312. **BLEU: 1** |
| 0x058f8f8aU, 0x3f9292adU, 0x219d9dbcU, 0x70383848U, 0xf1f5f504U, 0x63bcbcdfU, 0x77b6b6c1U, 0xafdada75U, 0x42212163U, 0x20101030U, 0 | 0x058f8f8aU, 0x3f9292adU, 0x219d9dbcU, 0x70383848U, 0xf1f5f504U, 0x63bcbcdfU, 0x77b6b6c1UL, 0xafdada75U,0x42212163U, 0x20101030U, 0 **BLEU: 0.93** | 0x058f8f8aU, 0x3f9292adU, 0x219d9dbcU, 0x70383848U, 0xf1f5f504U, 0x63bcbcdfU, 0x77b6b6c1U, 0xafdada75U, 0x42212163U, 0x20101030U, 0 **BLEU: 1** |
| 7, calc(sin((pi/180)*a7))) define(cea0, calc(cos((pi/180)*ea0))) define(cea1, calc(cos((pi/180)*ea1))) define(cea2, calc(cos((pi/180)*ea2))) define(cea3, calc(cos((pi/180)*ea3))) define(cea4, calc(cos((pi/180 | 7, calc(sin((pi/180)*a7))) define(cea0, calc(cos((pi/180)*ea0))) define(cea1, calc(cos((pi/180)*ea1))) define(cea2, calc(cos((pi/180)*ea2))) define(cea3, calc(cos((pi/180)*ea3))) define(cea4, calc(cos((pi/180 **BLEU: 0.95** | 7, calc(sin((pi/180)*a7))) define(cea0, calc(cos((pi/180)*ea0))) define(cea1, calc(cos((pi/180)*ea1))) define(cea2, calc(cos((pi/180)*ea2))) define(cea3, calc(cos((pi/180)*ea3))) define(cea4, calc(cos((pi/180 **BLEU: 1** |

Figure 17: Random sample of MEMFREE generations where the BLEU score with the true continuation > 0.9. Most of these examples are repetitive and/or lists of numbers. In the MEMFREE column, we use highlights to show the difference from the true continuation: red means deleted text, and green means added text.

# H  Author Ordering Algorithm

```
import hashlib
import numpy as np

def hash(x):
  h=hashlib.new("md5")
  h.update(bytes(x,"ascii"))
  return int(h.hexdigest(),16)

names = ("Nicholas Daphne " +
  "Katherine Matthew " +
  "Florian Chiyuan Milad " +
  "Christopher").split()

for i in range(0,10000):
  s = str(i)
  l = [hash(x+s) for x in names]
  o = np.argsort(l)
  if names[o[0]] != "Daphne":
    continue
  if names[o[-1]] != "Nicholas":
    continue
  print([names[x] for x in o])
  exit(0)
```

Figure 18: Author ordering algorithm

# Fine-Tuning GPT-3 for Synthetic Danish News Generation

**Mina Almasi** and **Anton Drasbæk Schiønning**

Aarhus University, Department of Linguistics, Cognitive Science and Semiotics
mina.almasi@post.au.dk, drasbaek@post.au.dk

## Abstract

While GPT-3 has garnered significant attention for its capabilities in natural language generation, research on its use outside of English is still relatively limited. We focus on how GPT-3 can be fine-tuned for generating synthetic news articles in a low-resource language, namely Danish. The model's performance is evaluated on the dimensions of human and machine detection in two separate experiments. When presented with either a real or GPT-3 generated news article, human participants achieve a 58.1% classification accuracy. Contrarily, a fine-tuned BERT classifier obtains a 92.7% accuracy on the same task. This discrepancy likely pertains to the fine-tuned GPT-3 model oversampling high-likelihood tokens in its text generation. Although this is undetectable to the human eye, it leaves a statistical discrepancy for machine classifiers to detect. We address how decisions in the experimental design favoured the machine classifiers over the human evaluators, and whether the produced synthetic articles are applicable in a real-world context.

## 1 Introduction

In recent years, rapid development in natural language processing, particularly in the area of pretrained language models, has led to significant advancements in various language tasks. State-of-the-art models, such as GPT-3 (Brown et al., 2020) and BERT (Devlin et al., 2019), have excelled in tasks such as classification of documents (Kong et al., 2022), text completion (Balkus and Yan, 2022), language translation (Yan et al., 2022) and text summarization (Wazery et al., 2022). These advances have even led some to suggest that we are currently experiencing a whole paradigm shift in NLP with the introduction of pretrained language models (Min et al., 2021).

However, most contemporary studies using GPT-3 focus on its performance in English. This is to be expected as the model was almost exclusively trained on English with less than 8% of training data being non-English (OpenAI, 2020). Still, the few investigations on GPT-3 for non-English text generation show promising results (Kraft et al., 2022; Müller and Laurent, 2022). This even holds for low-resource languages such as Catalan (Armengol-Estapé et al., 2021).

Yet, the more prevalent approach in NLP for low-resource languages has been using smaller language-specific models or multilingual models such as mBERT (Doddapaneni et al., 2021). This is despite multilingual models seemingly lacking in natural language generation tasks, especially for the Nordic languages and other low-resource languages (Rönnqvist et al., 2019; Wu and Dredze, 2020). In terms of language-specific models, this development has also occurred in Danish NLP with several Danish models appearing based on the likes of BERT and ELECTRA (e.g., Tamimi-Sarnikowski, 2021 and Møllerhøj, 2021). Nevertheless, such models are miniscule in size compared to the state-of-the-art language models. For instance, the Danish BERT model by Møllerhøj (2021) is trained on 9.7 billion characters. Comparatively, GPT-3's total training data corresponds to 1.1 trillion characters (OpenAI, 2020).

In this paper, we seek to understand how well GPT-3 can perform for a low-resource language such as Danish when optimized for that language through fine-tuning. To our knowledge, this is the first structured assessment of GPT-3's capabilities in a Danish NLP task. Concretely, we investigate whether GPT-3 can be fine-tuned to produce synthetic news articles that are indistinguishable to real news articles written by journalists. Generating news articles with GPT-3 is a common task with previous work showing remarkable results in English (Brown et al., 2020; Uchendu et al., 2021).

Inspired by a similar study from Ippolito et al. (2020), we make a two-fold evaluation of the model's performance:

(A) **Human Detection**: Can untrained human participants distinguish between real and synthetic articles in an experimental setting?

(B) **Machine Detection**: Can machine classifiers be trained to distinguish between real and synthetic articles?

As human and machine detection methods presumably apply distinct techniques to spot synthetically generated text (Ippolito et al., 2020), a two dimensional evaluation provides a more nuanced insight into how GPT-3 performs on the task.

Our findings suggest that a fine-tuned GPT-3 can generate convincing Danish synthetic news, deceiving human readers while being identifiable by a BERT classifier. This demonstrates GPT-3's capacity to perform succesfully in the context of low-resource languages, but with the drawback of heightened machine-detectability due to an overuse of high-probability tokens.

## 2 Related Work

### 2.1 Natural Language Generation with Pretrained Language Models

Natural language generation (NLG) is a subfield of NLP concerned with the process of producing intelligible language. However, even within this subfield, there are a diverse range of related sub-tasks. Examples of such tasks, which have natural language as the input and output, are summarization, question answering and translation (Celikyilmaz et al., 2021).

Similar to other fields in NLP, text generation has evolved rapidly with the paradigm of pretrained language models. These models have been critical for advancing NLG as they understand natural language, express it fluently and are capable of being fine-tuned for a specific domain (Li et al., 2021). Importantly, pretrained language models can generate natural language that is novel rather than just outputting text memorized from the training data. This was demonstrated in McCoy et al. (2021) who found GPT-2 and Transformer-XL to produce novel words and unique syntactic structures not found in the training data.

The demonstrated successes of GPT-3 in NLG cannot only be attributed to the sheer amount of data it has seen, but also to the underlying decoder-transformer architecture. GPT-based models are built using only decoder blocks which possess a masked self-attention layer that prevents the language model from considering future context (Wang et al., 2022). This architecture is more easily applicable to NLG tasks than the alternative encoder-only structures found in BERT-based models (Lewis et al., 2020).

### 2.2 The Fine-Tuning Approach

The groundbreaking paper introducing GPT-3 titled, "Language Models are Few-Shot Learners" highlighted its ability to achieve strong performance on various NLP tasks after only seeing a few examples (Brown et al., 2020). This few-shot learning approach can be contrasted with fine-tuning in which the model is updated through re-training with task-specific data. Although GPT-3 arguably excels at text generation from few-shot learning, OpenAI recommends fine-tuning the model for custom applications citing advantages such as higher quality results.[1]

Related research has also already demonstrated the flexibility of GPT-3 by fine-tuning it for a wide variety of tasks. Perhaps the most ambitious fine-tune of GPT-3 is OpenAI's Codex which was trained on 159 GB of Python files from 54 million GitHub repositories. As a result of this fine-tune, Codex outperformed base GPT-3 on a benchmark on several different coding tasks (Chen et al., 2021). However, fine-tunes of GPT-3 extend beyond just massive applications. A study by Zong and Krishnamachari (2022) on extracting equations from math word problems found an 80% accuracy for a fine-tuned GPT-3 model compared to only 40% accuracy with 3-shot learning. Contrary to the enormous Codex fine-tune, their fine-tune had just seen 1000 examples. Similar small-scale fine-tunes of GPT-3 improved its abilities for assessing students' short answer questions (Moore et al., 2022) and writing less biased job advertisements (Borchers et al., 2022).

The takeaway from these studies is that GPT-3 can improve performance through fine-tuning for specific downstream tasks despite its generalized task excellence from just few-shot learning.

---

[1]https://platform.openai.com/docs/guides/fine-tuning

### 2.3 Evaluating Synthetically Generated Text

#### 2.3.1 Human Detection

Evaluating whether artificial intelligence can deceive humans dates back to the Turing Test (Turing, 1950). While the current state of AI is incomparable to the 1950s, the underlying idea of judging machines on their human-like performance is still relevant. Much of research using this approach evaluated language models by asking participants to classify whether text excerpts were human-written or synthetically generated (Bogaert et al., 2022; Brown et al., 2020; Uchendu et al., 2021).

Although these classifications provide valuable insight into a language model's capabilities, they leave many questions as to why and how these models excel. For this reason, other studies ask participants to rate various qualities of the text without knowing whether the text is synthetic or real. The exact qualities that are rated differ across studies. For instance, some studies judge the overall text quality (Zhang et al., 2020) or fluency (Adelani et al., 2020) on a Likert scale. Dou et al. (2022)'s SCARECROW framework offers a more systematic approach to analyzing synthetic text, accessible to laypeople with basic training. It groups common error types within categories, like language errors for grammar and incoherence, and factual errors for incorrect or nonsensical information.

#### 2.3.2 Machine Detection

Although SCARECROW provides a standardized human evaluation of language models, human detection may not be ideal for detecting GPT-3 news articles as low accuracies would suggest. For instance, Clark et al. (2021) found that human evaluators only unmasked GPT-3 news stories with 56% accuracy despite them being trained for the task. Yet, this does not imply that synthetic text cannot be detected at all. In fact, past research on synthetic text detection has found machines to be superior to humans (Ippolito et al., 2020; Meyer et al., 2022; Uchendu et al., 2021). For example, Ippolito et al. (2020) utilized both a bag-of-words logistic regression and a fine-tuned BERT, reporting much greater performance than human evaluators. While the BERT model was optimal, the bag-of-words model did not lag far behind. As formulated by the study, the high performing machine detectors are likely due to the sampling method of language models being skewed towards high-likelihood words. Therefore, synthetic text is more

easily distinguishable from human language which has greater variability in word choice (Holtzman et al., 2020). This linguistic difference is also noted in other research (Gehrmann et al., 2019; Tay et al., 2020).

Nevertheless, models relying solely on word probabilities are still inferior to more complex language models such as BERT. This may indicate that there are other factors which differentiate real and synthetic articles that language models pick up on with fine-tuning. Just like Ippolito et al. (2020), Uchendu et al. (2021) found that the fine-tuned BERT was the best performing detector across text generated by 19 language models including GPT-3.

## 3 Data

The real news stories were all sourced from the Danish news site tv2.dk. In October 2022, TV2's news platform boasted over 3 million unique users (Danske Medier Research, 2022), which is more than half of Denmark's population. Hence, it makes an excellent representation of typical news content consumed by Danes. These articles were obtained via two channels: directly scraping from TV2 and employing the DaNewsRoom Danish news database (Varab and Schluter, 2020).

In the selection process, only article bodies with a minimum length of 100 words were considered, and longer articles were shortened to a maximum of 150 words. Although the exact threshold is somewhat arbitrary, it was kept in this range for two reasons. Firstly, accumulating costs for generating articles with the fine-tuned GPT-3 necessitated that we kept the articles short. Also, using longer articles would entail that each participant would evaluate fewer articles as their time was limited.

In total, 1866 real Danish news articles from TV2 were sourced and used for three purposes: Fine-tuning GPT-3 (1209 real articles), providing training/validation data for machine classifiers (609 real articles), and serving as test data in the experiments (48 real articles). Additionally, 657 synthetic articles were generated by the fine-tuned GPT-3 for training the classifiers (609 synthetic articles) and test data in the experiments (48 synthetic articles).

## 4 Methods

### 4.1 Fine-Tuning GPT-3

GPT-3, specifically *text-davinci-002*, was fine-tuned with 1209 pre-processed real news articles

Figure 1: Training loss for fine-tuning GPT-3. The dashed lines indicate an epoch ending (1209 examples).

using OpenAI's API in Python. All articles were formatted to JSONL in accordance with the API documentation.[2] The headlines and subheadings were combined to be the prompts for GPT-3, and the corresponding article bodies were completions. All hyperparameters set for the fine-tune are detailed in Appendix A.1. As the training loss for the fine-tune plateaued during the fourth epoch (Figure 1), we ended model training after this epoch.

### 4.2 Generating Synthetic News Articles

The fine-tuned GPT-3 was then applied to generate synthetic news articles. As in the training phase, the prompts (headline and subheading) came from real news articles.[3] When generating the text completions, we modified several of the default hyperparameters based on previous research for similar cases and OpenAI's general recommendations.[4]

Firstly, GPT-3's temperature sampling method was adjusted by setting the temperature parameter to zero. In temperature sampling, a high temperature means that low probability tokens are more frequently sampled. By setting temperature to zero, the model becomes deterministic, always sampling the most probable token when generating text. We made this adjustment because a high temperature may lead to factual errors as the model "takes more risks". Also, Ippolito et al. (2020) found that a zero temperature in conjunction with a frequency penalty was the most successful for generating English news articles. This parameter penalizes a new token based on how frequently it appears in the generated text so far. It can be used along with a

presence penalty (penalizes solely based on presence rather than frequency) to decrease the likelihood of sampling repetitive token sequences. As repetitiveness is also a particular concern for GPT-3's text generation (Dehouche, 2021), we add small presence and frequency penalties of 0.2. The full specification including ranges for the hyperparameters are in Appendix A.2.

The text completions formed the synthetic news articles, utilized as training/validation data for machine detection and test data for both experiments. Sanity checks were made to verify that these articles were similar to the real news articles in length and structure, but we made no modifications to them whatsoever.

## 5 Experiment A: Human Detection

Experiment A is designed as a binary classification task where human participants must distinguish between real articles written by humans and GPT-3's synthetic ones.

### 5.1 Participants

120 participants (66.6% female, age: M = 30.0, SD = 13.7) voluntarily took part in the online study. The study was run on the online platform SoSci Survey (Leiner, 2022) for one week in October 2022. To ensure a wide participant reach, the study was optimized for both computer and smartphone use. Complying with the prerequisites for the study, all participants were adult Danish native speakers.

### 5.2 Experimental Procedure

In each experimental trial, participants saw a page with one news article and four questions to be answered (see Appendix A.5). Participants were firstly asked to evaluate whether they believed the article body to be written by a human or an artificial intelligence. Subsequently, participants had to rate their confidence on a 5-point Likert scale from *completely unsure* (1) to *completely sure* (5). Finally, participants were asked to label whether the article had any distracting language or factual errors. These error types were inspired by the SCARE-CROW framework but simplified as the full framework would be too complex for untrained evaluators. To ensure participants understood what the error types implied, examples were written beneath each question. The articles were formatted to be closer in appearance with a real news article. This was done by differentiating in the size and color of

---

[2]https://platform.openai.com/docs/guides/fine-tuning/prepare-training-data

[3]To avoid double-dipping, these headlines and subheadings came from articles that were not part of the 1866 real articles mentioned in the data section.

[4]platform.openai.com/docs/api-reference

the headline, subheading and the article body (Appendix A.5). Importantly, it was clearly stated that only the body should be evaluated, not the headline and subheading as those always originated from real news stories.

In total, each participant evaluated 16 articles (8 real and 8 synthetic) in a randomly shuffled order. To cover the wide topical variance within news articles, 96 articles were used across all participants. That is, each participant only assessed a sixth of the total article pool, which corresponds to every article being evaluated by 20 unique participants.

### 5.3 Results

#### 5.3.1 Human Accuracy

With 20 assessments of 96 articles, the human results are based on 1920 total classifications. The overall classification accuracy was just 58.1%. This means that participants only performed eight percentage points over chance level which is a comparable result to similar studies conducted in English (see 2.3). Interestingly, when presented with a synthetic news article, participants correctly labeled it as machine-written 53.6% of the time. Contrarily, a true positive rate of 62.6% indicates that participants were better at identifying real news articles as human-written. In addition, it should also be underlined that none of the 96 articles were exclusively classified correctly or incorrectly. The articles that were the easiest to identify were classified correctly 95% of the time, whereas there were only 15% correct classifications for the hardest ones.

Moreover, none of the 120 participants answered correctly on all 16 articles that they saw, with all of them misclassifying at least one synthetic news article as real news. This implies that the synthetic news articles have fooled all 120 participants to some extent.

Furthermore, all participants were screened on their news consumption level and prior knowledge of GPT-3. To see whether domain expertise caused enhanced performance, a mixed effects logistic regression model was run with media consumption level and GPT-3 knowledge as fixed effects. The news article ID is used as a random effect to account for variance that is specific to the articles.[5]

The full model output is displayed in Appendix A.6. The baseline/intercept in the model corresponds to a participant who never reads news and

---

[5]accuracy ~ news consumption + gpt-3 knowledge + (1 | article ID)



Figure 2: Confidence rating distribution of all trials. The fill indicates whether the corresponding classifications were correct or not.

never had heard of GPT-3 prior to the experiment. The output reveals that a higher level of news consumption does not lead to significantly higher accuracies. However, compared to the baseline, we see significant improvements for participants that have heard of GPT-3 before ($\beta = 0.327, odds\,ratio = 0.581, SE_\beta = 0.131, p = 0.013$) and those participants that have additionally read GPT-3 texts ($\beta = 0.478, odds\,ratio = 0.617, SE_\beta = 0.146, p = 0.001$). This suggests that having GPT-3 knowledge may give an advantage in demarcating real from synthetic news, although participants who had worked with GPT-3 (highest level of GPT-3 knowledge) did not outperform the baseline.

#### 5.3.2 Confidence and Error Identification

Participants were also asked to rate their confidence in the classification as well as marking error types for each trial. Figure 2 reveals that participants typically abstain from the most extreme confidence ratings of *Completely unsure (1)* and *Completely sure (5)*. As expected, participants' accuracy is around the chance level for low confidences. However, even when claiming to be *Completely sure*, the fraction of correct answers only increases to 69%. For confidences of *Fairly sure (4)*, this drops to only 60% correct answers.

We also see some interesting patterns in error type responses by the participants. Figure 3 illustrates which errors were marked for real and synthetic articles respectively. Overall, the pat-

terns are strikingly similar. The figure reveals that participants most often did not find errors in the articles. When errors then were marked, there was a propensity to find language errors over factual errors for both real and synthetic articles. Despite the similarities, participants were more inclined to identify both factual and language errors for synthetic articles than for real ones. However, this did not necessarily lead to correct classifications. For instance, when participants marked *Both* errors, articles were almost exclusively labeled as synthetic (89.7% of cases) although Figure 3 reveals that this was often incorrect.

In sum, participants struggled with demarcating real news articles from synthetic ones in Experiment A. The overall accuracy was only 58% with classifications of synthetic news articles approaching chance level. Also, all 120 participants were fooled by at least one synthetic article and even the most confident classifications frequently led to wrong responses. Finally, patterns in error types marked by participants are similar for real and synthetic articles which shows the participants' inability to demarcate the articles by style and content.

## 6 Experiment B: Machine Detection

Experiment B explores whether it is possible to construct machine classifiers that are capable of distinguishing between real and synthetic articles. This is approached with logistic regression using bag-of-words (BOW) and TF-IDF as baseline models. The more advanced language model, NB-BERT-LARGE, is then fine-tuned, tested and evaluated against the baselines and human participants.

### 6.1 Building Classifiers

Two baseline classifiers are constructed using logistic regression with BOW and TF-IDF numerical representations of the vocabulary within the entire corpus (see Appendix A.3 for their hyperparameters). The BOW classifier is the most simple baseline, solely representing word frequencies within each document. TF-IDF provides a more detailed representation by also accounting for a word's rarity in relation to the entire set of documents.

Expanding beyond purely vocabulary-based classification, we fine-tune the BERT model, NB-BERT-LARGE (Kummervold et al., 2021), for the binary classification task. This BERT model was pretrained on the Norwegian Colossal Corpus which is a diverse collection of textual data



Figure 3: The marked error types by participants. On the left, the responses to real articles are found, and the synthetic responses are on the right. The fill indicates how participants classified the corresponding articles.

(Kummervold et al., 2022). Although Norwegian is the primary language of the corpus, the collection contains several languages. Notably, Danish is the biggest language after Norwegian with 13.6% of the corpus being in Danish. We fine-tuned the model with the Trainer API using Hugging Face's transformers package (Wolf et al., 2020) in Python. The fine-tuning dataset comprised 1218 labeled articles split into a training and validation set (training: 75%, validation: 25%). Half of these were the real news articles from TV2 and the other half synthetic news articles. The test data comprised the same 96 articles that humans evaluated in Experiment A.

The hyperparameters for the fine-tuning of BERT are detailed in Appendix A.4. Resulting from an early stopping callback,[6] the model was fine-tuned for two epochs, obtaining a validation accuracy of 95.7%.

### 6.2 Results

#### 6.2.1 Classification Accuracies

Table 1 shows the results of both the machine and human detection on the test data of 96 articles. The fine-tuned BERT model outclasses humans at the task with a 92.7% accuracy on the test set as well as the highest F1-score. Also, even the baseline BOW and TF-IDF models performed substantially better than the human average accuracy with accuracies around 80%, indicating that vocabulary discrepancies can demarcate the real and synthetic articles to an extent.

---

[6] based on the validation accuracy

59

| Classifier | Accuracy | F1 | Precision | Recall | TP | TN | FP | FN |
|---|---|---|---|---|---|---|---|---|
| Human | 0.581 | 0.599 | 0.575 | 0.626 | 62.6% | 53.6% | 46.4% | 37.4% |
| BOW | 0.802 | 0.796 | 0.822 | 0.771 | 77.1% | 83.3% | 16.7% | 22.9% |
| TF-IDF | 0.802 | 0.800 | 0.809 | 0.792 | 79.2% | 81.3% | 18.8% | 20.8% |
| BERT (fine-tuned) | 0.927 | 0.927 | 0.932 | 0.927 | 87.5% | 97.9% | 2.1% | 12.5% |

Table 1: Evaluation metrics for all classifiers on the test data of 96 articles.

An interesting similarity between all machine classifiers is their tendency to classify articles as synthetic. This is most noticeable with fine-tuned BERT which has 12.5% false negatives as opposed to just 2.1% false positives. Remarkably, BERT's true negative classifications of 97.9% means that the model has only classified a single synthetic article wrong. This propensity to classify articles as synthetic contrasts human participants, who had a bias towards classifying most articles as real.

### 6.2.2 Classifier Agreement

We turn to examine classifier agreement quantitatively by evaluating their inter-rater reliability using Cohen's Kappa. Unsurprisingly, this metric reveals that TF-IDF and BOW have an almost perfect agreement, $\kappa = 0.91, z = 3.37, p < 0.05$. Moreover, both TF-IDF ($\kappa = 0.62, z = 6.14$) and BOW ($\kappa = 0.62, z = 6.11$) have a substantial agreement with BERT that is greater than would be expected by chance (both $p < 0.05$).

Table 2 gives a qualitative insight into the agreements with examples of how four test articles were classified. Article A was the most commonly misclassified article for humans (17 out of 20 misclassifications). However, interestingly, all three classifiers correctly identified it as synthetic news. Additionally, article B is one of four instances where BERT correctly identified a synthetic news article while both BOW and TF-IDF failed. Oppositely, article C provides an example of BERT's over-inclination to classify as synthetic. It is one of three articles where BERT misclassified a real news article while BOW and TF-IDF did not. Finally, article D is the only synthetic article that BERT misclassified. However, as Table 2 shows, BOW and TF-IDF also struggled with this article.

The overall takeaway remains that these machine detections performed vastly better than human participants. This improvement was clear even for the two baseline models based on BOW and TF-IDF. Still, the more sophisticated fine-tuned BERT classifier performed the best by far, with an impressive 92.7% overall accuracy and just a single

misclassification of the 48 synthetic articles.

## 7 Limitations

A few limitations must be addressed in relation to these results. Firstly, several design decisions presumably favoured the machine detectors over the human evaluators. Whereas 78.3% of human participants had never seen GPT-3 produced texts before, all machine classifiers received extensive training on over 1000 labelled articles prior to the final testing. Also, the zero temperature token sampling for generating synthetic articles created an overrepresentation of high-likelihood tokens. This may be identified by the machine detectors, whereas such patterns are are probably too subtle to notice for humans (Ippolito et al., 2020). Also, Dou et al. (2022) show that higher temperatures are associated with GPT-3 making off-prompt errors. Such errors would not be captured by the machine classifiers, whereas humans would more likely identify these more semantic shortcomings.

Moreover, it must be addressed that human classifications are possibly influenced from being conducted in an experimental setting. Contrary to the machine classifiers, the human participants saw the headline and subheading for all articles. Despite being repeatedly told not to evaluate them, it cannot be dismissed that these extra elements still could have influenced their decision-making process. For instance, a familiar headline could have evoked an intuition for the article being real before reading the article body. On the other hand, one could argue that this was beneficial for humans as they could improve assessments by comparing contents in the headline and subheading to the article body.

Still, these methodological decisions systematically favored the machine classifiers over the human evaluators. However, asserting that the machine superiority would evaporate based on these considerations is a reach considering how vast the performance gap was.

Another limitation relates to the generalizability of the synthetic news articles. Due to experimen-

| Article A | | | | |
|---|---|---|---|---|
| **Correct** | **Human** | **BOW** | **TF-IDF** | **BERT** |
| *Synthetic* | **Real** | **Synthetic** | **Synthetic** | **Synthetic** |

Greenland's government has decided not to apply for permission for further oil drilling in the coming year. This is announced by the Greenlandic Ministry of Nature, Environment and Agriculture in a press release. "We have decided not to apply for oil drilling in 2023, because we want to spend time developing a new strategy for the Greenlandic economy, which will form the basis for a new oil and gas strategy," it says. The government also emphasizes that it will maintain its "vision of a fossil-free Greenland". The decision comes after a meeting on Tuesday between the government's four parties. It is mainly the consideration for the climate that has led the government to drop further oil drilling.

| Article B | | | | |
|---|---|---|---|---|
| **Correct** | **Human** | **BOW** | **TF-IDF** | **BERT** |
| *Synthetic* | **Synthetic** | **Real** | **Real** | **Synthetic** |

Two photographers and a culture minister are now criticized by the Press Council for having participated in a photo series where they posed with weapons. The Press Council writes this in a press release. In the case against Culture Minister Ane Halsboe-Jørgensen (S), the council has assessed that she has violated good press ethics by participating in the photo series 'The Gun Series'. "By participating in a photo series with weapons and ammunition, the Culture Minister has expressed that it is acceptable to carry weapons, whether it is in connection with artistic photography or not," the decision states. The decision against photographer Rasmus Flindt Pedersen and Jim Lyngvild is more stringent. Both have violated good press ethics by participating in the photo series, says the Press Council.

| Article C | | | | |
|---|---|---|---|---|
| **Correct** | **Human** | **BOW** | **TF-IDF** | **BERT** |
| *Real* | **Real** | **Real** | **Real** | **Synthetic** |

Consideration for endangered animal species such as hazel dormice, birch mice, and bats in Denmark is now temporarily hindering a massive natural gas project that aims to supply Poland with natural gas from Norway. The Environmental and Food Appeals Board has annulled the project's environmental permit, thereby halting the construction work of the Baltic Pipe pipeline across Denmark. "We are very disappointed with the decision," says Marian Kaagh, the deputy director of the company Energinet, which is responsible for the construction work in Denmark. In a press release, she says that Energinet has been working on a number of initiatives to ensure good living conditions for the animals in the areas where the pipeline is being built. This was a requirement when the Environmental Protection Agency granted the environmental permit for the Baltic Pipe project in 2019. However, according to the Environmental and Food Appeals Board, the conditions should have been thoroughly investigated before the permit was issued and construction work could begin.

| Article D | | | | |
|---|---|---|---|---|
| **Correct** | **Human** | **BOW** | **TF-IDF** | **BERT** |
| *Synthetic* | **Real** | **Real** | **Real** | **Real** |

The upcoming super hospitals are meant to help improve the healthcare system in Denmark. However, they will not be completed on time. On average, the 16 hospital constructions are almost two years delayed, according to a statement from the Quality Fund for hospital constructions, which TV 2 has obtained access to. It is an expression of "a number of challenges," as the fund's director, Morten Hjortenberg, puts it. "We had hoped for better results halfway through the construction period. It raises concerns and reflections on some of the decisions and priorities that were made during the projects," he says. The fund's task is to provide funding for hospital constructions and ensure high quality – that is, what is often called "quality funds." The total budget for the constructions is over 30 billion Danish kroner – of which the state accounts for 23 billion and the regions' self-financing contribution is 9 billion.

Table 2: Four article bodies from the test data translated to English. Predictions are marked in green if they were correct and red for incorrect. The human prediction is based on the majority classification for the 20 participants for an article (classified as real if split 50/50). See Appendix A.7 for the original articles in Danish.

tal constraints, articles were shortened greatly, and may therefore not be comparable to what we consider news in a real-world context. In addition, even if it could write longer articles, our fine-tuned GPT-3 model's capabilities are practically useless in a journalistic context despite producing human-like outputs. This is because inferring a factually correct article body from just a headline requires additional, current context about the world which is inaccessible in this setup. Instead, the only thinkable purposes for this "headline-to-article news generator" have malicious undertones such as automating fake news production.

## 8 Conclusion

As advancements in natural language processing continue to progress rapidly, it is crucial to remember the importance of including and improving upon NLP in low-resource languages. This paper acknowledges this need by conducting a structured assessment of GPT-3's abilities for Danish natural language generation when fine-tuned for the task.

Our study shows that GPT-3 can be fine-tuned to produce Danish synthetic news articles that are virtually indistinguishable to real news articles for humans. However, this does not imply that the articles are actually indistinguishable as the human eye is not all-seeing. By constructing a fine-tuned BERT model for the same discrimination task, we

find that machine detection of the synthetic news articles was possible to a great extent. Hence, there must have been underlying flaws in GPT-3's article generations, likely relating to an oversampling of high-likelihood words.

The introduction of ChatGPT and GPT-4 will likely impact the findings presented in this paper, lowering detection accuracies further for both humans and machines. Although, as those models are closed-sourced, it would be troublesome to assess whether the testing articles are already part of the training data which poses a methodological challenge. Regardless, as our findings for Danish conform with similar studies in English, we encourage future work on low-resource languages to develop machine detectors which possibly stand the test when human evaluators are deceived.

*Supplementary Materials Availability Statement:* All source code used in the project is available from GitHub at https://github.com/drasbaek/finetuning-gpt3-danish-news. A dataset with the synthetic articles as well as classifications made by machine detectors is also available on the GitHub. The dataset containing human responses from Experiment A cannot be made available due to GDPR regulations. The real news articles from TV2 are also not made publicly available due to copyright limitations. In the interest of reproducibility, dummy data is made available on the GitHub which mimics the actual data to the greatest possible extent under the circumstances. Contact the authors for more information on the project.

## Ethical Considerations

In this paper, we have created a GPT-3 fine-tune that is capable of producing synthetic news. As it may be possible to use it for malicious purposes, the fine-tuned model will not be available to anyone besides the authors. Per January 4, 2024, the authors will also lose access to the model as OpenAI announced all davinci models, including fine-tunes, will depreciate. [7] Nonetheless, we acknowledge that this paper demonstrates the ease of producing such a model, but also how it may be detected.

Finally, we recognize that the synthetic news produced for this paper could potentially contain societal biases from GPT-3's training data or from the real news articles used for fine-tuning.

---

[7] https://openai.com/blog/gpt-4-api-general-availability

## References

David Ifeoluwa Adelani, Haotian Mai, Fuming Fang, Huy H. Nguyen, Junichi Yamagishi, and Isao Echizen. 2020. Generating Sentiment-Preserving Fake Online Reviews Using Neural Language Models and Their Human- and Machine-Based Detection. In *Advanced Information Networking and Applications*, Advances in Intelligent Systems and Computing, pages 1341–1354, Cham. Springer International Publishing.

Jordi Armengol-Estapé, Ona de Gibert Bonet, and Maite Melero. 2021. On the Multilingual Capabilities of Very Large-Scale English Language Models. ArXiv:2108.13349 [cs].

Salvador Balkus and Donghui Yan. 2022. Improving Short Text Classification With Augmented Data Using GPT-3. ArXiv:2205.10981 [cs].

Jérémie Bogaert, Marie-Catherine de Marneffe, Antonin Descampe, and Francois-Xavier Standaert. 2022. Automatic and Manual Detection of Generated News: Case Study, Limitations and Challenges. In *Proceedings of the 1st International Workshop on Multimedia AI against Disinformation*, MAD '22, pages 18–26, New York, NY, USA. Association for Computing Machinery.

Conrad Borchers, Dalia Gala, Benjamin Gilburt, Eduard Oravkin, Wilfried Bounsi, Yuki M Asano, and Hannah Kirk. 2022. Looking for a Handsome Carpenter! Debiasing GPT-3 Job Advertisements. In *Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)*, pages 212–224, Seattle, Washington. Association for Computational Linguistics.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020.

Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Asli Celikyilmaz, Elizabeth Clark, and Jianfeng Gao. 2021. Evaluation of Text Generation: A Survey. ArXiv:2006.14799 [cs].

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. Evaluating Large Language Models Trained on Code. ArXiv:2107.03374 [cs].

Elizabeth Clark, Tal August, Sofia Serrano, Nikita Haduong, Suchin Gururangan, and Noah A. Smith. 2021. All That's 'Human' Is Not Gold: Evaluating Human Evaluation of Generated Text. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 7282–7296, Online. Association for Computational Linguistics.

Danske Medier Research. 2022. Toplisten.

N Dehouche. 2021. Plagiarism in the age of massive Generative Pre-trained Transformers (GPT-3). *Ethics in Science and Environmental Politics*, 21:17–23.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Sumanth Doddapaneni, Gowtham Ramesh, Mitesh M. Khapra, Anoop Kunchukuttan, and Pratyush Kumar. 2021. A Primer on Pretrained Multilingual Language Models. ArXiv:2107.00676 [cs].

Yao Dou, Maxwell Forbes, Rik Koncel-Kedziorski, Noah A. Smith, and Yejin Choi. 2022. Is GPT-3

Text Indistinguishable from Human Text? Scarecrow: A Framework for Scrutinizing Machine Text. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7250–7274, Dublin, Ireland. Association for Computational Linguistics.

Sebastian Gehrmann, Hendrik Strobelt, and Alexander Rush. 2019. GLTR: Statistical Detection and Visualization of Generated Text. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 111–116, Florence, Italy. Association for Computational Linguistics.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The Curious Case of Neural Text Degeneration. In *International Conference on Learning Representations*.

Daphne Ippolito, Daniel Duckworth, Chris Callison-Burch, and Douglas Eck. 2020. Automatic Detection of Generated Text is Easiest when Humans are Fooled. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1808–1822, Online. Association for Computational Linguistics.

Jun Kong, Jin Wang, and Xuejie Zhang. 2022. Hierarchical BERT with an adaptive fine-tuning strategy for document classification. *Knowledge-Based Systems*, 238:107872.

Angelie Kraft, Hans-Peter Zorn, Pascal Fecht, Judith Simon, Chris Biemann, and Ricardo Usbeck. 2022. *Measuring Gender Bias in German Language Generation*. Gesellschaft für Informatik, Bonn. Accepted: 2022-09-28T17:10:03Z ISSN: 1617-5468.

Per E Kummervold, Javier De la Rosa, Freddy Wetjen, and Svein Arne Brygfjeld. 2021. Operationalizing a National Digital Library: The Case for a Norwegian Transformer Model. In *Proceedings of the 23rd Nordic Conference on Computational Linguistics (NoDaLiDa)*, pages 20–29, Reykjavik, Iceland (Online). Linköping University Electronic Press, Sweden.

Per E Kummervold, Javier De la Rosa, Freddy Wetjen, and Svein Arne Brygfjeld. 2022. Norwegian Colossal Corpus Description.

D.J. Leiner. 2022. SoSci Survey.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Junyi Li, Tianyi Tang, Wayne Xin Zhao, and Ji-Rong Wen. 2021. Pretrained Language Model for Text Generation: A Survey. volume 5, pages 4492–4499. ISSN: 1045-0823.

R. Thomas McCoy, Paul Smolensky, Tal Linzen, Jianfeng Gao, and Asli Celikyilmaz. 2021. How much do language models copy from their training data? Evaluating linguistic novelty in text generation using RAVEN. ArXiv:2111.09509 [cs].

Selina Meyer, David Elsweiler, Bernd Ludwig, Marcos Fernandez-Pichel, and David E. Losada. 2022. Do We Still Need Human Assessors? Prompt-Based GPT-3 User Simulation in Conversational AI. In Proceedings of the 4th Conference on Conversational User Interfaces, CUI '22, pages 1–6, New York, NY, USA. Association for Computing Machinery.

Bonan Min, Hayley Ross, Elior Sulem, Amir Pouran Ben Veyseh, Thien Huu Nguyen, Oscar Sainz, Eneko Agirre, Ilana Heinz, and Dan Roth. 2021. Recent Advances in Natural Language Processing via Large Pre-Trained Language Models: A Survey. ArXiv:2111.01243 [cs].

Steven Moore, Huy A. Nguyen, Norman Bier, Tanvi Domadia, and John Stamper. 2022. Assessing the Quality of Student-Generated Short Answer Questions Using GPT-3. In Educating for a New Future: Making Sense of Technology-Enhanced Learning Adoption, Lecture Notes in Computer Science, pages 243–257, Cham. Springer International Publishing.

Jens Dahl Møllerhøj. 2021. certainlyio/nordic_bert.

Martin Müller and Florian Laurent. 2022. Cedille: A large autoregressive French language model. ArXiv:2202.03371 [cs].

OpenAI. 2020. openai/gpt-3: Languages by Character Count.

Samuel Rönnqvist, Jenna Kanerva, Tapio Salakoski, and Filip Ginter. 2019. Is Multilingual BERT Fluent in Language Generation? In Proceedings of the First NLPL Workshop on Deep Learning for Natural Language Processing, pages 29–36, Turku, Finland. Linköping University Electronic Press.

Phillip Tamimi-Sarnikowski. 2021. sarnikowski/danish_transformers.

Yi Tay, Dara Bahri, Che Zheng, Clifford Brunk, Donald Metzler, and Andrew Tomkins. 2020. Reverse Engineering Configurations of Neural Text Generation Models. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 275–279, Online. Association for Computational Linguistics.

A. M. Turing. 1950. Computing Machinery and Intelligence. Mind, 59(236):433–460. Publisher: [Oxford University Press, Mind Association].

Adaku Uchendu, Zeyu Ma, Thai Le, Rui Zhang, and Dongwon Lee. 2021. TURINGBENCH: A Benchmark Environment for Turing Test in the Age of Neural Text Generation. In Findings of the Association for Computational Linguistics: EMNLP 2021, pages 2001–2016, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Daniel Varab and Natalie Schluter. 2020. DaNewsroom: A Large-scale Danish Summarisation Dataset. In Proceedings of the Twelfth Language Resources and Evaluation Conference, pages 6731–6739, Marseille, France. European Language Resources Association.

Haifeng Wang, Jiwei Li, Hua Wu, Eduard Hovy, and Yu Sun. 2022. Pre-Trained Language Models and Their Applications. Engineering.

Y. M. Wazery, Marwa E. Saleh, Abdullah Alharbi, and Abdelmgeid A. Ali. 2022. Abstractive Arabic Text Summarization Based on Deep Learning. Computational Intelligence and Neuroscience, 2022:e1566890. Publisher: Hindawi.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-Art Natural Language Processing. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 38–45, Online. Association for Computational Linguistics.

Shijie Wu and Mark Dredze. 2020. Are All Languages Created Equal in Multilingual BERT? In Proceedings of the 5th Workshop on Representation Learning for NLP, pages 120–130, Online. Association for Computational Linguistics.

Rong Yan, Jiang Li, Xiangdong Su, Xiaoming Wang, and Guanglai Gao. 2022. Boosting the Transformer with the BERT Supervision in Low-Resource Machine Translation. Applied Sciences, 12(14):7195. Number: 14 Publisher: Multidisciplinary Digital Publishing Institute.

Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter Liu. 2020. PEGASUS: Pre-training with Extracted Gap-sentences for Abstractive Summarization. In Proceedings of the 37th International Conference on Machine Learning, pages 11328–11339. PMLR. ISSN: 2640-3498.

Mingyu Zong and Bhaskar Krishnamachari. 2022. Solving Math Word Problems Concerning Systems of Equations with GPT-3. Proceedings of the Thirteenth AAAI Symposium on Educational Advances in Artificial Intelligence, page 8.

# A    Appendix

## A.1    Fine-tuning Parameters for GPT-3

| Parameters | Value |
|---|---|
| Batch Size | 2 |
| Learning Rate Multiplier | 0.2 |
| Prompt Loss Weight | 0.01 |
| Epochs | 4 |

## A.2    Text Generation Parameters for GPT-3

| Parameters | Value | Value Range |
|---|---|---|
| Temperature | 0 | 0 to 1 |
| Frequency Penalty | 0.2 | -2 to 2 |
| Presence Penalty | 0.2 | -2 to 2 |
| Max Tokens | 400 | 0 to 2048 |

## A.3    Logistic Regression Parameters for BOW and TFIDF

| Parameters | Value |
|---|---|
| Solver | lfbgs |
| C | 10 |
| Random State | 2 |
| Max Iterations | 250 |

## A.4    Fine-Tuning Parameters for BERT

| Parameters | Value |
|---|---|
| Learning Rate | 2e-5 |
| Weight Decay | 0.01 |
| Batch Size | 24 |
| Epochs run | 2 |
| Max Epochs | 5 |

The fine-tuned BERT was defined with an early stopping callback which stopped model training if the validation accuracy did not improve for 3 epochs. The final model used for inference was thus only run for 2 epochs.

## A.5  Experimental Procedure

**University in massive data leak: - The criminals have all the information now, says expert**

**A student discovered that with just a few clicks, he could see others' social security numbers.**

A large number of sensitive personal details about Danish students are currently available to anyone who wants to search the internet. This is because a database containing information on about 100,000 students at the University of Copenhagen has been leaked. It happened on Thursday evening when a student at the University of Copenhagen discovered that with just a few clicks, he could see other people's social security numbers. The leak was hidden behind a blurred address on the internet, and it required a so-called reverse lookup service to find it. TV 2, using this service, found the address, and it can be seen that it contains a large number of files with information about the approximately 100,000 students.

**Do you think that the article body is written by a human or artificial intelligence ?**

☐ **Human**

☐ **Artificial Intelligence**

**How sure are you of your answer?**

| Completely unsure | Slightly sure | Somewhat sure | Fairly sure | Completely sure |
|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 |
| ☐ | ☐ | ☐ | ☐ | ☐ |

**Are there any distracting language errors?**
E.g.,  spelling mistakes,  wrong punctuation, incoherent or repetitive language

☐ **Yes**

☐ **No**

**Are there any distracting factual errors?**
E.g.,  contradicting information or factual mistakes about individuals or events

☐ **Yes**

☐ **No**

Illustration of a trial from experiment A. All text was written in Danish in the actual experiment. The article body in the example is synthetically generated. The captions "HEADLINE", "SUBHEADING" and "ARTICLE BODY" did not appear in the actual experiment.

## A.6 Logistic Regression Model Output for Predicting Accuracy

| Fixed Effect | Estimate | Standard Error | Z-value | P-value |
|---|---|---|---|---|
| Intercept | 0.33668 | 0.39335 | 0.856 | 0.39204 |
| News_Consumption_2 | -0.50311 | 0.43260 | -1.163 | 0.24484 |
| News_Consumption _3 | -0.03473 | 0.39697 | -0.087 | 0.93028 |
| News_Consumption_4 | -0.27516 | 0.40664 | -0.677 | 0.49862 |
| News_Consumption_5 | -0.10105 | 0.39719 | -0.254 | 0.79817 |
| GPT_Knowledge_2 | 0.32738 | 0.13130 | 2.493 | 0.01266 |
| GPT_Knowledge_3 | 0.47842 | 0.14626 | 3.271 | 0.00107 |
| GPT_Knowledge_4 | 0.37824 | 0.22513 | 1.680 | 0.09293 |

| Fixed Effect Level | Participant Response (translated) |
|---|---|
| News_Consumption_1 | Never read the news |
| News_Consumption_2 | Very rarely read the news |
| News_Consumption_3 | Read news every week but not daily |
| News_Consumption_4 | Read news once every day |
| News_Consumption_5 | Read news multiple times a day |
| GPT_Knowledge_1 | Never heard of GPT-3 |
| GPT_Knowledge_2 | Heard of GPT-3, but never read anything it wrote or worked with it |
| GPT_Knowledge_3 | Heard of GPT-3 and read texts it wrote, but never worked with it |
| GPT_Knowledge_4 | Heard of GPT-3, read texts it wrote and worked with it |

| | Participant Response (original) |
|---|---|
| | Læser aldrig nyheder |
| | Læser meget sjældent nyheder |
| | Læser nyheder hver uge men ikke dagligt |
| | Læser nyheder en gang om dagen |
| | Læser nyheder flere gange om dagen |
| | Aldrig hørt om GPT-3 |
| | Hørt om GPT-3, men aldrig læst noget den har skrevet eller arbejdet med den |
| | Hørt om GPT-3 og læst tekster den har skrevet, men aldrig arbejdet med den |
| | Hørt om GPT-3, læst tekster den har skrevet og arbejdet med den |

## A.7   Classifier Agreement (Table 2) Danish Original Text

| Article A | | | | |
|---|---|---|---|---|
| **Correct** | **Human** | **BOW** | **TF-IDF** | **BERT** |
| *Synthetic* | **Real** | **Synthetic** | **Synthetic** | **Synthetic** |
| Grønlands regering har besluttet sig for ikke at søge om tilladelse til yderligere olieboringer i det kommende år. Det oplyser det grønlandske ministerium for natur, miljø og landbrug i en pressemeddelelse. - Vi har besluttet os for ikke at søge om olieboringer i 2023, fordi vi vil bruge tid på at udvikle en ny strategi for den grønlandske økonomi, som skal danne grundlag for en ny olie- og gasstrategi, lyder det. Regeringen understreger samtidig, at den vil fastholde sin "vision om et fossilfrit Grønland". Beslutningen kommer efter et møde tirsdag mellem regeringens fire partier. Det er især hensynet til klimaet, der har fået regeringen til at droppe yderligere olieboringer. | | | | |

| Article B | | | | |
|---|---|---|---|---|
| **Correct** | **Human** | **BOW** | **TF-IDF** | **BERT** |
| *Synthetic* | **Synthetic** | **Real** | **Real** | **Synthetic** |
| To fotografer og en kulturminister får nu kritik af Pressenævnet for at have deltaget i en billedserie, hvor de poserede med våben.  Det skriver Pressenævnet i en pressemeddelelse.  I sagen mod kulturminister Ane Halsboe-Jørgensen (S) har nævnet vurderet, at hun har brudt god presseskik ved at deltage i billedserien 'The Gun Series'. - Kulturministeren har ved deltagelse i billedserie med våben og ammunition givet udtryk for, at det er acceptabelt at bære våben, uanset om det er i forbindelse med kunstnerisk fotografering eller ej, lyder det i afgørelsen. Afgørelsen mod fotografen Rasmus Flindt Pedersen og Jim Lyngvild er mere knibsk. Begge har brudt god presseskik ved deltagelse i billedserien, mener Pressenævnet. | | | | |

| Article C | | | | |
|---|---|---|---|---|
| **Correct** | **Human** | **BOW** | **TF-IDF** | **BERT** |
| *Real* | **Real** | **Real** | **Real** | **Synthetic** |
| Hensynet til truede dyrearter som hasselmus, birkemus og flagermus i Danmark stikker nu en midlertidig kæp i hjulet på et enormt naturgasprojekt, der skal forsyne Polen med naturgas fra Norge. Det er Miljø- og Fødevareklagenævnet, der har annulleret projektets miljøtilladelse, og dermed har sat en stopper for anlægsarbejdet af rørledningen Baltic Pipe på tværs af Danmark.  - Vi er meget kede af afgørelsen, siger Marian Kaagh, der er vicedirektør i selskabet Energinet, der står for anlægsarbejdet I Danmark. I en pressemeddelelse siger hun, at Energinet har arbejdet med en række tiltag for at sikre gode levevilkår for dyrene de steder, hvor rørledningen bliver anlagt. Det var et krav, da Miljøstyrelsen i 2019 gav miljøtilladelse til Baltic Pipe-projektet. Men ifølge Miljø- og Fødevareklagenævnet burde forholdene være grundigt undersøgt, allerede inden tilladelsen blev udstedt, og anlægsarbejdet kunne begynde. | | | | |

| Article D | | | | |
|---|---|---|---|---|
| **Correct** | **Human** | **BOW** | **TF-IDF** | **BERT** |
| *Synthetic* | **Real** | **Real** | **Real** | **Real** |
| De kommende supersygehuse skal være med til at løfte sundhedsvæsenet i Danmark. Men de bliver ikke klar til tiden. Gennemsnitligt er de 16 sygehusbyggerier knap to år forsinkede, viser en opgørelse fra Kvalitetsfonden for sygehusbyggerierne, som TV 2 har fået aktindsigt i. Det er et udtryk for, at der er "en del udfordringer", som fondens direktør, Morten Hjortenberg, siger det. - Vi havde håbet på bedre resultater her halvvejs inde i byggeperioden. Det giver anledning til bekymring og eftertanke om nogle af de beslutninger og prioriteringer, der blev truffet under projekterne, siger han. Fondens opgave er at stille penge til rådighed for sygehusbyggerierne og sikre en høj kvalitet – altså det man ofte kalder "kvalitetsfonde". Byggeriernes samlede budget er på over 30 milliarder kroner – heraf står staten for 23 milliarder og regionernes selvfinansierende bidrag på 9 milliarder. | | | | |

# GAN-LM: Generative Adversarial Network using Language Models for Downstream Applications

**Dae Yon Hwang    Yaroslav Nechaev    Cyprien De Lichy    Renxian Zhang**
Amazon Alexa AI
{dyhwang, nechaey, cllichy, renxiz}@amazon.com

## Abstract

In this work, we investigate Data Augmentation methods to improve the performance of state-of-the-art models for four different downstream tasks. Specifically, we propose Generative Adversarial Network using Language Models (GAN-LM) approach that combines a deep generative model with a pre-trained language model to produce diverse augmentations. We compare the GAN-LM to various conventional methods in non-contextual- and contextual-levels on four public datasets: *ZESHEL* for zero-shot entity linking, *TREC* for question classification, *STS-B* for sentence pairs semantic textual similarity (STS), and *mSTS* for multilingual sentence pairs STS. Additionally, we subsample these datasets to study the impact of such augmentations in low-resource settings where limited amounts of training data is available. Compared to the state-of-the-art methods in downstream tasks, we mostly achieve the best performance using GAN-LM approach. Finally, we investigate the way of combining the GAN-LM with other augmentation methods to complement our proposed approach. The developed code for reproducibility is included in the supplementary material.[1]

## 1 Introduction

Nowadays, the availability of large unsupervised corpora and computational resources has led to development of large language models (LMs) that are now employed across a wide variety of natural language processing (NLP) tasks including but not limited to entity linking (EL), text summarization, question classification (QC) and semantic textual similarity (STS). While such models can sometimes work well for tasks where little (few-shot) or no (zero-shot) supervised data is available, the performance loss in such low-resource settings can be substantial compared to their high-resource counterparts. This gap is even larger for

low-resource languages. Thus, scientists in both industry and academia still have to rely on a multitude of methods Hedderich et al. (2021), such as Data Augmentation (DA), to yield sufficient levels of performance on their low-resource tasks.

DA allows to artificially increase the size of a dataset by generating additional synthetic examples from the existing ones. A large amount of diverse training data is important to ensure the generalization of a model but it is not always possible to collect due to cost and time constraints or lack of target language data and task experts. To mitigate this issue, DA can be used to improve performance.

In this work, we test a wide variety of DA approaches, both the ones found in the literature as well as our own approach, on four different tasks: Zero-shot EL with *ZESHEL* dataset, QC with *TREC* database, STS with *STS-B* dataset, and multilingual STS with *mSTS* database. Two different levels of augmentations are considered: (1) Non-contextual, or word-level, and (2) Contextual, where full sentence is considered for DA. To further highlight the impact of different DA approaches, we produced low-resource versions of the above-mentioned tasks by subsampling training sets and removing rich textual contexts where applicable. We propose a novel Generative Adversarial Network using Language Models (GAN-LM) which employs GAN with Wasserstein distance to improve the stability of training and uses the pre-trained LM for generating synthetic textual data to extend its usability. We also introduce tunable thresholds and a decoding method to control the diversity and lexical similarity of synthetic data to mitigate the mode collapse problem in GAN. Compared to other DAs, GAN-LM employs an adversarial training with the offered data in each task to learn the characteristic of it which generates suitable synthetic data to aid in downstream tasks (covered in Section 5.6). Even if we used pre-trained LM in GAN-LM, we do not use its

---

[1] https://github.com/amazon-science/data-augmentation-for-entity-resolution

69

generation ability (e.g. paraphrase, text generation) for downstream tasks. To complement our approach, we mix GAN-LM with other DAs (e.g. Back-translation, GPT) to enhance further in low-resource languages and limited entity linking task. The source code used to train the GAN model and produce augmentations listed in this paper is publicly released and attached with a paper.

## 2 Related Work

Originally, DAs for NLP relied on synonyms to increase diversity and dataset size. Synonyms could be found in various resources like WordNet Miller et al. (1990) and PPDB Ganitkevitch and Callison-Burch (2014). In Wang and Yang (2015), they considered word embedding with K-Nearest-Neighbor (KNN) and cosine similarity to search and substitute similar words. Other pre-trained word embeddings such as Word2Vec Mikolov et al. (2013), GloVe Pennington et al. (2014) and fastText Bojanowski et al. (2017) have been leveraged for that purpose. Furthermore, the authors in Wei and Zou (2019) generated synthetic texts by changing the words through synonym replacement or random insertions, substitutions and deletions where Shou et al. (2022) include the abstract meaning representation graph along with it for STS task. In Pruthi et al. (2019), the authors simulated spelling mistakes by random insertions, substitutions, character swaps and deletions to enhance the robustness of the model for sentiment analysis. Also, punctuation as DA was considered in Karimi et al. (2021) for QC task. Later, back-translation with Neural Machine Translation (NMT) was employed to generate variations of target words Sennrich et al. (2016).

More emerging techniques for DA are using deep neural networks which mostly use auto-regressive language model to predict words from a given context, e.g. GPT-2 Radford et al. (2019), XLNet Yang et al. (2019) and BART Lewis et al. (2020) which have been used for DA in diverse applications such as question-answering, text classification and machine translation. Using LMs, KNN-based DA with knowledge distillation Kamalloo et al. (2021) is considered for QC task. There are also works related to the adversarial learning to understand their effects on language models. Alzantot et al. (2018) proposed a black-box population-based optimization to generate the imperceptible adversarial examples to fool the models. In Zhang et al. (2019a), they considered Metropolis-Hastings

attack to generate the adversarial examples which were tested in terms of attack and training.

However, there are relatively few works using GANs for text generation even if it is one of the most notable approaches in other domains Antoniou et al. (2017). In Kusner and Hernández-Lobato (2016), the authors used a GAN model with Gumbel-Softmax to have a differentiable sampling distribution approximating a categorical one. In Subramanian et al. (2017), diverse GANs with recurrent and convolutional architectures were evaluated for text augmentation at word and character-levels. Yu et al. (2017) proposed a sequence GAN with reinforcement learning to address the problem of assessing a partially generated sequence. Another work in Nie et al. (2018) developed a GAN model consisting of relational memory-based generator, the Gumbel-Softmax relaxation, and multiple embedded representations in the discriminator. In Golovneva and Peris (2020), authors explored a data generation for the bootstrapping of a new language and the handling of low-resource features using a sequential GAN. Croce et al. (2020) used the fine-tuning of BERT with unlabeled data in a generative adversarial setting to reduce the time consuming of annotating the data but did not extend to the DA application. Similarly, Thakur et al. (2021) use the cross-encoder to label the new inputs for training a bi-encoder model. Marek et al. (2021) focus on out-of-domain data generation with a sequential GAN to build the robust dialog system. Compared to these works, GAN-LM combines a large pre-trained model and GAN with tunable thresholds to suitably control the diversity and similarity of generated data and it was tested on various downstream tasks to assess its generalizability. Also, we can use any pre-trained LM on top of the GAN part which extends its applicability to various tasks. To highlight the effectiveness of DA, the low-resource settings are investigated in Shi et al. (2021) and Hedderich et al. (2021) where we mainly investigate the different size of training set and suggest a way to define the optimal size of augmented data.

## 3 Data Augmentations

### 3.1 Non-Contextual-Level Augmentation

In this work, we utilize four augmentation approaches as non-contextual-level. ***Lexical:*** We use WordNet Miller et al. (1990) to replace each word in the original text with a synonym. ***Spelling:*** We generate alternate texts from common misspellings

of the original words Coulombe (2018). **Character:** Here, we randomly change characters in the original tokens with four different ways: Insertions, substitutions, swaps and deletions Pruthi et al. (2019). For lexical, spelling and character-based methods, we use the implementation in *nlpaug*[2] with 10% replacement. **Token-LM:** To understand the effectiveness of GAN part in GAN-LM, we consider pre-trained LMs solely. To generate the synthetic data: (1) Use LM to get token embeddings for input text and (2) perform nearest neighbor search for each token to find alternate tokens that meets the similarity thresholds. We search the synthetic tokens which satisfied these thresholds to balance the analogy and diversity, compared to the original token. The similarity thresholds are defined empirically (e.g. Table 5). We did not insert the noise on the input embedding as GAN-LM since the generated data is far from the original one.

## 3.2 Contextual-Level Augmentation

To extend our work, we experiment three methods as contextual-level augmentation. **Text Generation:** This is a typical auto-regressive generation which uses the original text as the initial context and extends it Yang et al. (2020). For this, we employ GPT-2 Radford et al. (2019) and OPT Zhang et al. (2022) for English-based datasets, and mGPT Tan et al. (2021) for multilingual dataset. **Paraphrase:** This augmentation transforms a sentence with similar semantic meaning but a different syntactic form where we consider the fine-tuned T5 model Raffel et al. (2020) on Google PAWS Zhang et al. (2019b) for English-based tasks and Prism model Thompson and Post (2020a,b) for multilingual-based task. **Back-translation:** It is a process of retranslating content from the target language back to its source language to generate a sentence variant. For this augmentation, we employ multiple pre-trained neural translation models trained on OPUS data Helsinki-NLP (2023) with *nlpaug*.

## 3.3 Generative Adversarial Network

GAN is basically coming from the adversarial learning which aims to trick the model by providing deceptive input. GAN targets to correctly classify both unmodified and adversarial examples to receive the rewards. It consists of two neural

networks, generator and discriminator, where each of them tries to outplay the other. The goal of generator is to artificially manufacture outputs that could be hard to distinguish from real data. The discriminator is similar to the usual classification model that aims to differentiate between real and synthetic data from generator. Using GAN, we target to achieve eminent performances with only offered train set in each downstream task.

Specifically, we considered a WGAN-GP Gulrajani et al. (2017) which uses the Wasserstein distance as loss to capitalize on the probability distributions from fake and real data rather than labeled samples. Compared to the vanilla GAN, it is robust to vanishing gradient and mode collapse through smoother gradient updates from its loss functions.

## 3.4 GAN-LM

To extend the usability of GAN in NLP domain, we propose GAN-LM which combines GAN with pre-trained LM regardless of non-contextualized and contextualized models. In this work, we focus on the latter one which promises the better result. Loss function of GAN-LM is covered in Equation (1) and its structure is shown in Appendix.

$$
\begin{aligned}
R &= LM_{encoder}(\text{Input Text}) \\
\epsilon &\sim \text{Uniform}(0,1),\ \eta \sim N(0,1) \\
F &= G(R+\eta), \quad \hat{F} = \epsilon \cdot R + (1-\epsilon) \cdot F \\
L_d &= D(F) - D(R) + \lambda \cdot (||\nabla_{\hat{F}} D(\hat{F})||_2 - 1)^2 \\
L_g &= -D(F)
\end{aligned}
$$
(1)

where $LM_{encoder}$ is the encoder of LM to generate the embedding of input text for augmentation. $\epsilon$ and $\eta$ are random numbers from the uniform and Gaussian distributions respectively. Also, $R$ is the real embedding and $F$ is the fake embedding generated from the generator, $G(\cdot)$. $\hat{F}$ is weighted embedding from real and fake embeddings. $D(\cdot)$ means the discriminator output for embeddings. $L_d$ and $L_g$ refer to the loss functions of the discriminator and generator respectively. $D(F) - D(R)$ describes the 1-Wasserstein distance and $\lambda \cdot (||\nabla_{\hat{F}} D(\hat{F})||_2 - 1)^2$ is called gradient penalty used for mitigating the vanishing gradient where we use $\lambda = 10$ based on the suggestion in Gulrajani et al. (2017).

Figure 1 illustrates the flowchart of the overall algorithm in GAN-LM. First, we generated embeddings for each input text to serve as a training set for the GAN. For embeddings, we use the pre-

Figure 1: GAN-LM with pre-trained LM. Before decoding from LM, we additionally consider the low and high thresholds for similarity matching between the averaged synthetic embedding from GAN-LM and the candidate embedding from defined dictionary in pre-trained LM to control the diversity and lexical similarity of synthetic text.

trained BART-base Lewis et al. (2020) and mBART-large-50 Tang et al. (2020) as encoder and decoder according to the empirical results (e.g. Table 5). However, GAN-LM can be applied with other approaches that can encode arbitrary text into embeddings and comprise of a well-defined dictionary to map generated embedding back to tokens. Using these transformers, we can decode the synthetic embedding into a text which can be different from the original input. We fix the sentence lengths to 27, 25, 36, 24 tokens for *ZESHEL, TREC, STS-B, mSTS* datasets respectively, which cover 99% of the data for each dataset. To express the text into an input data for GAN training, we stacked each token as the dimension with padding the remainder of the input with zero values (more covered in Appendix).

Figure 1 shows the steps of the algorithm in GAN-LM. In training part (red area), we encode the input text into embeddings using $LM_{encoder}$, then we add Gaussian noise on top and input resulting embeddings to the generator. Next, the generator produces synthetic embeddings which should resemble real ones and feeds those to the discriminator which tries to distinguish between real and synthetic embeddings. In synthetic data generation pipeline (yellow area), we feed the target text, for which we want to generate and alternate form, to the encoder and add Gaussian noise

to that embedding. The generator will produce the synthetic embedding for that target text and then we average the original and synthetic embeddings to maintain the structure of original text. To decode, we perform nearest neighbor search for each token using those generated synthetic embeddings. Finally, we introduce upper and lower thresholds on similarity to select tokens that are diverse yet still possess similar semantics compared to the original and don't accidentally change the meaning of the input text. The augmented tokens are selected randomly from tokens that meet those thresholds with the bias towards tokens for high similarity score. Since BART's vocabulary consists of subword units, GAN-LM is able to come up with new valid-looking words that were never in the original training set (see Table 4). To clarify, we employ the pre-trained LM for tokenization and detokenization in GAN-LM but we do not use its text generation methods (e.g. paraphrase) for downstream tasks.

## 4 Datasets and Employed Models

We experiment with four different downstream tasks where (1-3) are English-based databases and (4) is multilingual-based dataset: (1) *ZESHEL* – a zero-shot learning dataset for EL, (2) *TREC* – a text retrieval dataset for QC, (3) *STS-B* - an integrated version of STS tasks for measuring the semantic

similarity between two sentences, and (4) *mSTS* - multilingual version of STS task.

*ZESHEL* introduced by Logeswaran et al. (2019) is based on Wikia where there are non-overlapping domains in train/validation/test sets to simulate zero-shot learning. For this task, we employ BLINK's Wu et al. (2019) bi-encoder model from scratch. *TREC* shown in Li and Roth (2002) is collected from Hovy et al. (2001) where questions were manually created with 50 fine class labels. For this application, we use fine-tuned BERT-Tiny Turc et al. (2019) with training data in *TREC*. *STS-B* covered in Cer et al. (2017) includes news headlines, image captions and user forum posts. In each sentence pair, semantic similarity labels are provided by a number between 0 and 5. For this task, we use SentenceTransformers Reimers and Gurevych (2019) from scratch using the mean pooling layer with XLM-RoBERTa Conneau et al. (2020). *mSTS* introduced in Cer et al. (2017); Reimers and Gurevych (2020) has sentence pairs in different languages with semantic similarity scores between 0 and 5. For train set, we used the offered monolingual pairs of AR-AR, ES-ES and the translated sentences of ES-ES into EN, DE, TR, FR, IT, NL using Google Translator since we do not have monolingual pairs for them. The provided EN-EN dataset was eliminated from train set since most cross-lingual datasets were made from translating one sentence of EN-EN Reimers and Gurevych (2020). All the cross-lingual pairs are considered as test set. For this application, we employ the mean pooling of outputs for the pre-trained multilingual BERT (mBERT) Devlin et al. (2019) with fine-tuning from train set.

# 5 Results and Discussion

## 5.1 Experimental Setting

For all downstream tasks, in addition to the original size, we construct a low-resource version (i.e. limited train set) to highlight augmentation impact. In addition, *ZESHEL* contains rich textual context for both entity mentions and catalog entities, which provide additional information for EL. To isolate the impact of DA approaches, we test model performance with and without those contexts.

For augmentation, in *ZESHEL* we consider both the entity and mention to generate synthetic data, in *TREC* we synthesize an alternate question sentence, in *STS-B and mSTS* we generate an alternate sentence from one of the pair. To build the GAN-LM, we use pre-trained BART and mBART with 0.3-0.7 and 0.5-0.9 similarity thresholds respectively which give a good diversity of generated data while remaining close to the original semantics. The thresholds are decided from empirical results such as Table 5. Also, the size of augmentation is determined from the empirical results in validation set (e.g. Table 6) where we cover the optimal size in this work. Compared to other tasks, we additionally fine-tune GAN-LM with a target language (e.g. AR-AR) in *mSTS* after training with multilingual sentences to boost the quality of synthetic data.

In all tasks, we use the same target metrics as found in the literature. For *ZESHEL* we use recall@$k$, for *TREC* F1 score, for *STS-B and mSTS* the spearman's rank correlation (SRC) between the cosine similarity of sentence pairs embeddings and ground-truth labels. In all experiments, we retrain target model 3 times with different seeds and report average results with 95% confidence interval (CI). Finally, the computational cost for GAN-LM is covered in Appendix where it takes a longer training time compared to non-contextual-level methods, and comparable time to contextual-level approaches. However, GAN-LM promises the better performances in most cases and utilizes LMs without fine-tuning for generation purpose.

## 5.2 Results in Entity Linking

Table 1 shows the results for *ZESHEL*. In this application, we target to find the generalized augmentations for zero-shot learning task. In low-resource cases, the amount of train set with augmentation is 5K generated from 1K baseline while in full training data, the size of the training set after augmentation is 69K from 49K baseline.

Overall, improvements after augmentation in normal case are lower than for the low-resource scenario which confirms the importance of augmentation in the limited data setting. Including synthetic data can have an effect of inferring the unseen data which might have a different distribution from train set. Few training set samples in low-resource mean insufficient variation of data to help the models, especially high capacity ones, generalize well. Thus, data augmentation often improves more in low-resource scenarios, compared to normal case.

When we consider scenarios without context, we can see that there are large improvements in performance using augmentations, especially

Table 1: Recall values in *ZESHEL* with 95% CI. Baseline describes the performance of model without augmentation and change denotes the performances against baseline in absolute term. In each scenario, bold means the best results and underline denotes the proposed methods.

| Scenarios | Type | R@1 | R@8 | R@32 | R@64 | CI | Change |
|---|---|---|---|---|---|---|---|
| Normal without context | **_GAN-LM-GPT_** | **28.91%** | **54.83%** | **64.77%** | **69.38%** | 1.71% | **7.94%** |
| | _GAN-LM_ | 24.2% | 48.96% | 60.85% | 66.16% | **1.65%** | 3.51% |
| | GPT | 28.32% | 54.14% | 63.31% | 67.46% | 1.89% | 6.77% |
| | OPT | 27.54% | 53.28% | 62.81% | 67.15% | 1.89% | 6.16% |
| | Paraphrase | 22.1% | 46.89% | 59.1% | 64.73% | 2.03% | 1.67% |
| | Back-Translation | 20.7% | 44.77% | 57.13% | 62.99% | 2.06% | -0.14% |
| | Token-LM | 21.33% | 45.52% | 57.55% | 63.29% | 1.83% | 0.39% |
| | Char | 22.11% | 46.36% | 58.5% | 64.07% | 4.38% | 1.22% |
| | Spel | 21.52% | 45.76% | 58.22% | 63.88% | 2.25% | 0.81% |
| | Lexical | 20.67% | 44.8% | 57.23% | 62.91% | 2.01% | -0.13% |
| Low-resource without context | **_GAN-LM-GPT_** | **25.25%** | **50.94%** | **59.9%** | **63.8%** | 2.3% | **15.11%** |
| | _GAN-LM_ | 18.67% | 42.43% | 55.21% | 61.03% | **1.97%** | 9.47% |
| | GPT | 22.52% | 47.52% | 58.23% | 62.62% | 2.37% | 12.86% |
| | OPT | 19.76% | 45.07% | 57.06% | 61.82% | 2.33% | 11.07% |
| | Paraphrase | 17.83% | 41.16% | 53.79% | 60% | 2.41% | 8.33% |
| | Back-Translation | 16.14% | 37.71% | 50.63% | 56.82% | 2.84% | 5.46% |
| | Token-LM | 15.86% | 36.9% | 49.98% | 56.2% | 2.9% | 4.87% |
| | Char | 16.52% | 37.91% | 51.34% | 57.53% | 2.67% | 5.96% |
| | Spel | 16.11% | 37.44% | 50.63% | 56.87% | 3.88% | 5.4% |
| | Lexical | 15.56% | 36.67% | 49.9% | 56.01% | 2.24% | 4.67% |
| | Baseline - Low | 12.4% | 31.24% | 44.65% | 51.16% | 3.09% | - |
| | Baseline - Normal | 20.57% | 44.89% | 57.56% | 63.13% | 1.92% | - |

| Scenarios | Type | R@1 | R@8 | R@32 | R@64 | CI | Change |
|---|---|---|---|---|---|---|---|
| Normal with context | **_GAN-LM_** | **39.13%** | **66.45%** | **76.3%** | **79.98%** | **0.65%** | **1.23%** |
| | GPT | 37.36% | 65.31% | 74.78% | 78.65% | 1.54% | -0.21% |
| | OPT | 37.63% | 65.37% | 74.88% | 78.77% | 0.93% | -0.08% |
| | Paraphrase | 37.88% | 65.35% | 74.94% | 78.7% | 0.76% | -0.02% |
| | Back-Translation | 37.73% | 65.26% | 74.95% | 78.73% | 1.25% | -0.07% |
| | Token-LM | 37.53% | 64.58% | 74.49% | 78.41% | 1.27% | -0.49% |
| | Char | 37.53% | 64.68% | 74.6% | 78.56% | 1.37% | -0.4% |
| | Spel | 37.27% | 64.42% | 74.42% | 78.38% | 1.19% | -0.62% |
| | Lexical | 37.49% | 64.86% | 74.89% | 78.66% | 1.66% | -0.27% |
| Low-resource with context | **_GAN-LM_** | **23.93%** | **49.79%** | **61.5%** | **66.75%** | **1.29%** | **3.71%** |
| | GPT | 21.57% | 47.75% | 59.75% | 64.69% | 2.05% | 1.66% |
| | OPT | 22.84% | 47.99% | 60.47% | 65.38% | 1.68% | 2.39% |
| | Paraphrase | 20.13% | 45.59% | 58.36% | 63.62% | 1.75% | 0.14% |
| | Back-Translation | 17.6% | 42.25% | 54.86% | 60.84% | 1.98% | -2.9% |
| | Token-LM | 13.76% | 35.95% | 48.64% | 54.97% | 1.62% | -8.45% |
| | Char | 14.92% | 38.11% | 51.17% | 57.35% | 2.85% | -6.4% |
| | Spel | 19.46% | 44.46% | 56.85% | 62.54% | 4.71% | -0.96% |
| | Lexical | 17.59% | 41.68% | 54.03% | 60.18% | 2.62% | -3.41% |
| | Baseline - Low | 20.92% | 45.19% | 57.63% | 63.39% | 1.59% | - |
| | Baseline - Normal | 37.93% | 65% | 75.08% | 78.95% | 1.19% | - |

with contextual-level, and GAN-LM mostly outperforms others, except for GPT and OPT. In this case, EL model has been trained on only entity in train set to infer the entity with its contexts in test set. Thus, it can be beneficial to use the augmented data with additional descriptions to imitate the context of it which can be done by GPT and OPT. We further investigate the augmentation from a combination between GAN-LM and GPT, called GAN-LM-GPT. In this approach, GAN-LM generates alternate forms from the original inputs at the token-level and GPT adds new textual content after that. We observe improvements after combinations of both methods, especially in the low-resource case. Therefore, we can also consider GAN-LM-GPT augmentation when train data is limited without additional contexts in entity linking (EL) task since it helps to include the diverse variations in test set which cannot be covered by the considered train set. For scenarios with context, most augmentations, especially with non-contextual-level, decrease the performance since synthetic data from these approaches could be less related to the available contexts which could be harmful to EL. However, GAN-LM has tunable thresholds to control the diversity and similarity of synthetic data which finally promises the improvements. In conclusion, we observe that GAN-LM and its complement, GAN-LM-GPT, are the best choices for EL task whether in low-resource or normal setting. In *ZESHEL*, domains in test set are not overlapped with the ones in train set, which confirms that GAN-LM is fairly compared with other augmentations.

## 5.3 Results in Question Classification

Now, we test the influence of DAs for question classification (QC) task covered in Table 2 left side. In this task, we need label-invariant augmentations to improve the performance. The size of training data for augmentations is 1K from 109 baseline in low-resource and 8K from 2K baseline in half-train set case. Interestingly, the improvements after augmentations in both scenarios have a similar pattern: Contextual-level augmentations outperforms the non-contextual ones, except for spelling and lexical (only for low-resource) while GAN-LM is always the best performing approach. In addition, the improvements in half-train set scenario are higher than the ones in the low-resource. From our investigation, the result without augmentation in half-train set is considerably worse than the one in normal case (i.e. 8.84% F1 difference), meaning the effect of augmentation can be huge in half-train set to improve further. Also, adding synthetic data on model can be noticed as inferring the possible variations in test set which needs some degree of real traffic from original data to suitably utilize the augmented data. Still, GAN-LM works the best in

Table 2: F1 and SRC values in *TREC* and *STS-B* with 95% CI. Here, we did not cover the normal case for augmentation since we already achieve the better or similar performance with half-train set, compared to full training set without augmentation (i.e. Baseline - Normal). In normal scenario, GAN-LM gives 34.28% F1 score in *TREC* and 79.84% SRC in *STS-B*. Denotations are identical as Table 1.

**Question Classification in *TREC***

| Scenarios | Type | F1 | CI | Change |
|---|---|---|---|---|
| Half-train set | **GAN-LM** | **32.14%** | 2.23% | **16.01%** |
| | GPT | 29.16% | 2.66% | 13.03% |
| | OPT | 28.75% | 2.7% | 12.62% |
| | Paraphrase | 28.39% | 3% | 12.26% |
| | Back-Translation | 28.03% | 2.36% | 11.9% |
| | Token-LM | 27.16% | **1.67%** | 11.03% |
| | Char | 25.5% | 7.02% | 9.37% |
| | Spel | 29.05% | 2.16% | 12.92% |
| | Lexical | 26.93% | 5.02% | 10.8% |
| Low-resource | **GAN-LM** | **10.15%** | 1.95% | **9.27%** |
| | GPT | 8.48% | 3.61% | 7.6% |
| | OPT | 8.17% | 1.9% | 7.29% |
| | Paraphrase | 5.93% | 2.42% | 5.05% |
| | Back-Translation | 7.27% | 1.59% | 6.39% |
| | Token-LM | 5.26% | 3.72% | 4.38% |
| | Char | 4.19% | **1.42%** | 3.31% |
| | Spel | 7.68% | 4.03% | 6.8% |
| | Lexical | 6.09% | 3.3% | 5.21% |
| | Baseline - Low | 0.88% | 1.54% | - |
| | Baseline - Half | 16.13% | 1.16% | - |
| | Baseline - Normal | 24.97% | 2.27% | - |

**Semantic Textual Similarity in *STS-B***

| Scenarios | Type | SRC | CI | Change |
|---|---|---|---|---|
| Half-train set | **GAN-LM** | **78.02%** | 0.96% | **4.44%** |
| | GPT | 76.94% | 0.83% | 3.36% |
| | OPT | 76.97% | 1.65% | 3.39% |
| | Paraphrase | 77.07% | 2.01% | 3.49% |
| | Back-Translation | 77.1% | 2.4% | 3.52% |
| | Token-LM | 76.11% | **0.57%** | 2.53% |
| | Char | 75.43% | 0.86% | 1.85% |
| | Spel | 76.61% | 2.13% | 3.03% |
| | Lexical | 76.74% | 1.39% | 3.16% |
| Low-resource | **GAN-LM** | **61.66%** | 1.46% | **23.44%** |
| | GPT | 58.11% | 6.38% | 19.89% |
| | OPT | 59.17% | 3.95% | 20.95% |
| | Paraphrase | 57.9% | 3.1% | 19.68% |
| | Back-Translation | 58.02% | 6.72% | 19.8% |
| | Token-LM | 56.66% | 2.59% | 18.44% |
| | Char | 53.32% | 1.6% | 15.1% |
| | Spel | 54.52% | 5.07% | 16.3% |
| | Lexical | 57.77% | 5.17% | 19.55% |
| | Baseline - Low | 38.22% | 10.61% | - |
| | Baseline - Half | 73.58% | 4.08% | - |
| | Baseline - Normal | 78.49% | 0.28% | - |

this environment and can be a top pick which has 7.17% F1 improvement with half-train set against full train set without augmentation.

## 5.4 Results in Semantic Textual Similarity

Table 2 right side covers the results on the *STS-B* dataset. In this application, we need various and semantically closed augmented data to improve the result. The size of training data is 1K from 115 baseline in low-resource and 8K from 2K baseline in half-train set scenario. In low-resource, we can achieve great improvements, especially with contextual-level and GAN-LM approaches. In half-train set, the improvement is smaller than the one in low-resource setting but we can see consistent improvements by including synthetic data. Again, contextual-level augmentations outperforms non-contextual-level and GAN-LM yields the best performance for semantic textual similarity (STS) task which gives a closed performance as the result from full train set without augmentation.

## 5.5 Results in Multilingual Semantic Textual Similarity

Lastly, we extend our work to multilingual task. Table 3 shows the results on the *mSTS* dataset. In this task, we target diverse and semantically consistent augmented samples in multilingual to

enhance the performance. The amount of train set is 800 from 200 baseline in low-resource, and 4K from 2K baseline in normal scenario. In low-resource, we can confirm that all augmentations improve the overall performance, especially with GAN-LM. Compared to low-resource, the improvement in normal is lower but still, GAN-LM mostly gives the best results, except for EN-AR. This is because GAN-LM is mostly trained on Indo-European languages (i.e. EN, DE, NL, FR, ES, IT) which enhances the generation ability for these languages. Interestingly, GAN-LM works well in EN-TR since the performance without augmentation in this pair is very low and it has a large gap to be improved by augmentation, especially with GAN-LM which saves original structure with similarity thresholds and does token-level tweaking with affordable diversity learned from train set. We can find that back-translation works the best in EN-AR because it directly uses the well-defined neural translation models for augmentation which finally decreases the unsuitable assigned languages (e.g. code-switching) suffered by other augmentations. To complement our approach, we combine GAN-LM with back-translation, called GAN-LM-Back, to enhance the performance. In this method, we generate the synthetic data for AR-AR and EN-EN using back-translation and other monolingual pairs

Table 3: SRC values in *mSTS* with 95% CI. Here, we focus on the contextual-level augmentations which promise the superior performances in STS task. Denotations are identical as Table 1.

| Scenarios | Type | EN-AR | ES-EN | EN-DE | EN-TR | FR-EN | IT-EN | NL-EN | CI | Change |
|---|---|---|---|---|---|---|---|---|---|---|
| Normal | **GAN-LM -Back** | 46.18% | **55.92%** | **59.23%** | **43.72%** | 60.93% | **57.32%** | 53.9% | 2.64% | **2.38%** |
| | GAN-LM | 44.44% | 53.6% | 59.2% | 42.62% | **61.48%** | 55.31% | **53.96%** | 2.62% | 1.43% |
| | mGPT | 45.24% | 50.86% | 59.2% | 42.52% | 60.51% | 53.07% | 53.86% | 2.71% | 0.67% |
| | Paraphrase | 45.21% | 48.69% | 58.06% | 40.9% | 60.67% | 54.12% | 53.32% | 2.92% | 0.06% |
| | Back-Translation | **46.36%** | 50.62% | 57.26% | 41.82% | 58.64% | 53.48% | 52.98% | 2.72% | 0.08% |
| Low-resource | **GAN-LM** | **31.75%** | **37.05%** | **44.71%** | **24.21%** | **43.12%** | **39.96%** | **43.96%** | 3.06% | **5.43%** |
| | mGPT | 30.29% | 34.33% | 38.11% | 19.64% | 34.9% | 33.37% | 39.19% | 4.83% | 0.44% |
| | Paraphrase | 28.67% | 35.93% | 37.76% | 22.04% | 35.4% | 32.63% | 35.24% | 3.59% | 0.13% |
| | Back-Translation | 31.01% | 34.44% | 36.67% | 21.94% | 36.28% | 31.7% | 37.15% | 4.49% | 0.35% |
| | Baseline - Low | 29.95% | 33.13% | 36.04% | 18.23% | 37.26% | 34.68% | 37.46% | 3.85% | - |
| | Baseline - Normal | 45.08% | 50.52% | 56.9% | 40.94% | 60.89% | 53.16% | 53.08% | 2.47% | - |

using GAN-LM to fine-tune the mean pooling of mBERT. Using GAN-LM-Back, we achieve the overall enhancements. Thus, we can understand that GAN-LM and its extension, GAN-LM-Back, are the best approaches for mutlilingual STS task.

## 5.6 Analysis of Synthetic Data

In this section, we analyze the synthetic data from each augmentation method. Table 4 shows examples of synthetic data in *TREC* dataset. Lexical-based finds the synonym of the word, spelling and character-based tweak the words, and token-LM-based changes the lowercase word and auxiliary verb. Both back-translation and paraphrase restate a text with different orders and words while both OPT and GPT adds a new context after original statement. Interestingly, GAN-LM focuses on changing question word which is the main factor to increase the performance. It also finds a semantically similar word. Lastly, we can see that GAN-LM-GPT is the combination between GAN-LM and GPT. From our findings, DA improves model robustness to unseen noisy inputs in downstream tasks. Augmentations containing grammatical mistakes, speech recognition errors, semantically similar terms help the model generalize better. With GAN-LM, we preserve both the semantics and the structure of the input text, while providing diverse augmentations. More examples of augmented data are covered in Appendix.

## 5.7 Ablation Study

In Table 5, we example the ablation study of GAN-LM where 0.3-0.7 range as the similarity thresholds and BART as pre-trained LM are the best choices in *STS-B*. The similarity thresholds control the analogy and diversity of synthetic data where 0.3-0.7 range was the top choice in *STS-B* to balance these

Table 4: Examples of generated augmentations. Bold texts in each cell mean the changed parts.

| Type | Example |
|---|---|
| Original | Why do heavier objects travel downhill faster ? |
| Lexical | Why do heavier object travel downhill **quicker**? |
| Spelling | **Whay** do heavier objects travel downhill faster? |
| Character | Why do heavier **osbjects tralvel downhzill** faster? |
| Token-LM | **WHY does** heavier objects travel downhill faster ? |
| Back-Translation | Why **are the** heavier objects **moving down faster?** |
| Paraphrase | Why do heavier objects **go faster downhill?** |
| OPT | Why do heavier objects travel downhill faster ? **Because they're heavier** |
| GPT | Why do heavier objects travel downhill faster ? **Or slow down to 2 km h** |
| GAN-LM | **HOW** do heavier objects travel **down** faster ? |
| GAN-LM-GPT | **HOW** do heavier objects travel **down** faster ? **Or slow down to 2 km h** |

two terms for achieving the best performance. Similar patterns are observed in other downstream tasks, except for *mSTS* where mBART and 0.5-0.9 range are selected. Additional architectural ablation study is shown in Appendix.

In addition, we investigate the effect of size of augmented data in Table 6. We consider the validation set (or cross-validation for dataset without validation set) to determine the optimal size of augmented data and find that there is a specific point when the validation performance becomes stable. Our findings indicate that performances (SRC - Test in Table 6) are stabilized after this certain point, implying that the generated synthetic data offers sufficient diversity to improve the model's generalization capabilities. The size of augmentation in other tasks are determined by same approach.

## 6 Conclusions

In this work, we investigate the effect of different DAs to improve the performance on various tasks. We study both techniques found in the literature as well as the proposed GAN-LM in different scenarios: We subsample training sets to study model per-

Table 5: GAN-LM study in *STS-B* with half-train set.

| Type | SRC |
|---|---|
| GAN-LM with BART (0.3 - 0.7) | **78.02%** |
| GAN-LM with BART (0.1 - 0.5) | 75.57% |
| GAN-LM with BART (0.5 - 0.9) | 77.49% |
| GAN-LM with BERT (0.3 - 0.7) | 71.33% |
| GAN-LM with XLNet (0.3 - 0.7) | 74.21% |

Table 6: Investigation on the size of train set with GAN-LM. Validation and Test describe each set in *STS-B*.

| Low-resource in *STS-B* | | | |
|---|---|---|---|
| Type | Size of Train set | SRC - Validation | SRC - Test |
| | 690 | 65.56% | 56.81% |
| | 920 | 68.93% | 60.16% |
| GAN-LM | 1150 (same as Table 2) | **71%** | **61.66%** |
| | 1380 | 70.89% | 61.61% |

formance under low-resource conditions and use half or full training set to understand under different conditions. In most experiments, GAN-LM clearly gives the better results than non-contextual and contextual-level augmentations. In addition to apply GAN-LM solely, we combine it with GPT and back-translation to supplement the performance.

## 7 Limitations

There are three predictable limitations in the developed GAN-LM. First, the convergence of training process in GAN-LM should be investigated carefully. Different datasets have different distribution of data and characteristics which can affect the GAN-LM's convergence and we need a few iterations of training to confirm the suitable epochs for each task. Second, there can be a machine bias since each model is trained on machine generated synthetic data. Therefore, searching the suitable pre-trained model is important to be considered at the beginning. Last, while we did a thorough evaluation of GAN-LM on four downstream tasks, it is still a general-purpose approach and its effectiveness on specific tasks or domains may vary. Thus, further research is required to fully understand its capabilities and limitations in different contexts.

***Supplementary Materials Availability Statement:***
Source code is included in supplementary materials. Notes on reproducibility (e.g. computational budget and used hyperparameters) are included in Appendix. Additional ablation study and augmented examples are covered in Appendix. Links for considered datasets and models are shown in Appendix.

## References

Moustafa Alzantot, Yash Sharma, Ahmed Elgohary, Bo-Jhang Ho, Mani Srivastava, and Kai-Wei Chang. 2018. Generating natural language adversarial examples.

Antreas Antoniou, Amos Storkey, and Harrison Edwards. 2017. Data augmentation generative adversarial networks.

Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146.

Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. 2017. SemEval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 1–14, Vancouver, Canada. Association for Computational Linguistics.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Claude Coulombe. 2018. Text data augmentation made simple by leveraging nlp cloud apis. *arXiv preprint arXiv:1812.04718*.

Danilo Croce, Giuseppe Castellucci, and Roberto Basili. 2020. GAN-BERT: Generative adversarial learning for robust text classification with a bunch of labeled examples. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2114–2119, Online. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Juri Ganitkevitch and Chris Callison-Burch. 2014. The multilingual paraphrase database. In *The 9th edition*

*of the Language Resources and Evaluation Conference*, Reykjavik, Iceland. European Language Resources Association.

O. Yu. Golovneva and Charith S. Peris. 2020. Generative adversarial networks for annotated data augmentation in data sparse nlu. In *ICON*.

Ishaan Gulrajani, Faruk Ahmed, Martín Arjovsky, Vincent Dumoulin, and Aaron C. Courville. 2017. Improved training of wasserstein gans. In *NIPS*.

Michael A. Hedderich, Lukas Lange, Heike Adel, Jannik Strötgen, and Dietrich Klakow. 2021. A survey on recent approaches for natural language processing in low-resource scenarios. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2545–2568, Online. Association for Computational Linguistics.

Helsinki-NLP. 2023. Github - helsinki-nlp/opus-mt: Open neural machine translation models and web services.

Eduard Hovy, Laurie Gerber, Ulf Hermjakob, Chin-Yew Lin, and Deepak Ravichandran. 2001. Toward semantics-based answer pinpointing. In *Proceedings of the First International Conference on Human Language Technology Research*.

Ehsan Kamalloo, Mehdi Rezagholizadeh, Peyman Passban, and Ali Ghodsi. 2021. Not far away, not so close: Sample efficient nearest neighbour data augmentation via MiniMax. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3522–3533, Online. Association for Computational Linguistics.

Akbar Karimi, Leonardo Rossi, and Andrea Prati. 2021. AEDA: An easier data augmentation technique for text classification. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2748–2754, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Matt J Kusner and José Miguel Hernández-Lobato. 2016. Gans for sequences of discrete elements with the gumbel-softmax distribution. *arXiv preprint arXiv:1611.04051*.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *ArXiv*, abs/1910.13461.

Xin Li and Dan Roth. 2002. Learning question classifiers. In *COLING 2002: The 19th International Conference on Computational Linguistics*.

Lajanugen Logeswaran, Ming-Wei Chang, Kenton Lee, Kristina Toutanova, Jacob Devlin, and Honglak Lee. 2019. Zero-shot entity linking by reading entity descriptions. *arXiv preprint arXiv:1906.07348*.

Petro Marek, Vishal Ishwar Naik, Vincent Auvray, and Anuj Goyal. 2021. Oodgan: Generative adversarial network for out-of-domain data generation. *ArXiv*, abs/2104.02484.

Tomas Mikolov, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. In *ICLR*.

George A. Miller, Richard Beckwith, Christiane D. Fellbaum, Derek Gross, and Katherine J. Miller. 1990. Introduction to wordnet: An on-line lexical database. *International Journal of Lexicography*, 3:235–244.

Weili Nie, Nina Narodytska, and Ankit Patel. 2018. Relgan: Relational generative adversarial networks for text generation. In *International conference on learning representations*.

Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. GloVe: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.

Danish Pruthi, Bhuwan Dhingra, and Zachary Chase Lipton. 2019. Combating adversarial misspellings with robust word recognition. In *ACL*.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67.

Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks.

Nils Reimers and Iryna Gurevych. 2020. Making monolingual sentence embeddings multilingual using knowledge distillation.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. *ArXiv*, abs/1511.06709.

Haoyue Shi, Karen Livescu, and Kevin Gimpel. 2021. Substructure substitution: Structured data augmentation for NLP. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3494–3508, Online. Association for Computational Linguistics.

Ziyi Shou, Yuxin Jiang, and Fangzhen Lin. 2022. AMR-DA: Data augmentation by Abstract Meaning Representation. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 3082–3098, Dublin, Ireland. Association for Computational Linguistics.

Sandeep Subramanian, Sai Rajeswar, Francis Dutil, Christopher Pal, and Aaron Courville. 2017. Adversarial generation of natural language. In *Proceedings of the 2nd Workshop on Representation Learning for NLP*, pages 241–251.

Zhixing Tan, Xiangwen Zhang, Shuo Wang, and Yang Liu. 2021. Msp: Multi-stage prompting for making pre-trained language models better translators.

Yuqing Tang, Chau Tran, Xian Li, Peng-Jen Chen, Naman Goyal, Vishrav Chaudhary, Jiatao Gu, and Angela Fan. 2020. Multilingual translation with extensible multilingual pretraining and finetuning.

Nandan Thakur, Nils Reimers, Johannes Daxenberger, and Iryna Gurevych. 2021. Augmented SBERT: Data augmentation method for improving bi-encoders for pairwise sentence scoring tasks. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 296–310, Online. Association for Computational Linguistics.

Brian Thompson and Matt Post. 2020a. Automatic machine translation evaluation in many languages via zero-shot paraphrasing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Online. Association for Computational Linguistics.

Brian Thompson and Matt Post. 2020b. Paraphrase generation as zero-shot multilingual translation: Disentangling semantic similarity from lexical and syntactic diversity. In *Proceedings of the Fifth Conference on Machine Translation (Volume 1: Research Papers)*, Online. Association for Computational Linguistics.

Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Well-read students learn better: On the importance of pre-training compact models.

William Yang Wang and Diyi Yang. 2015. That's so annoying!!!: A lexical and frame-semantic embedding based data augmentation approach to automatic categorization of annoying behaviors using #petpeeve tweets. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 2557–2563, Lisbon, Portugal. Association for Computational Linguistics.

Jason Wei and Kai Zou. 2019. Eda: Easy data augmentation techniques for boosting performance on text classification tasks.

Ledell Wu, Fabio Petroni, Martin Josifoski, Sebastian Riedel, and Luke Zettlemoyer. 2019. Scalable zero-shot entity linking with dense entity retrieval. *arXiv preprint arXiv:1911.03814*.

Yiben Yang, Chaitanya Malaviya, Jared Fernandez, Swabha Swayamdipta, Ronan Le Bras, Ji-Ping Wang, Chandra Bhagavatula, Yejin Choi, and Doug Downey. 2020. Generative data augmentation for commonsense reasoning. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1008–1025, Online. Association for Computational Linguistics.

Zhilin Yang, Zihang Dai, Yiming Yang, Jaime G. Carbonell, Ruslan Salakhutdinov, and Quoc V. Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. In *NeurIPS*.

Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. 2017. Seqgan: Sequence generative adversarial nets with policy gradient. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31.

Huangzhao Zhang, Hao Zhou, Ning Miao, and Lei Li. 2019a. Generating fluent adversarial examples for natural languages. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5564–5569, Florence, Italy. Association for Computational Linguistics.

Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022. Opt: Open pretrained transformer language models.

Yuan Zhang, Jason Baldridge, and Luheng He. 2019b. PAWS: Paraphrase Adversaries from Word Scrambling. In *Proc. of NAACL*.

# Summaries as Captions: Generating Figure Captions for Scientific Documents with Automated Text Summarization

**Chieh-Yang Huang,[1]\* Ting-Yao Hsu,[1]\* Ryan Rossi,[2] Ani Nenkova,[2] Sungchul Kim,[2]**
**Gromit Yeuk-Yin Chan,[2] Eunyee Koh,[2] Clyde Lee Giles,[1] Ting-Hao 'Kenneth' Huang[1]**
[1]Pennsylvania State University, University Park, PA, USA.
{chiehyang,txh357,clg20,txh710}@psu.edu
[2]Adobe Research, San Francisco, CA, USA.
{ryrossi,nenkova,sukim,ychan,eunyee}@adobe.com

## Abstract

Good figure captions help paper readers understand complex scientific figures. Unfortunately, even published papers often have poorly written captions. Automatic caption generation could aid paper writers by providing good starting captions that can be refined for better quality. Prior work often treated figure caption generation as a vision-to-language task. In this paper, we show that it can be more effectively tackled as a **text summarization** task in scientific documents. We fine-tuned PEGASUS, a pre-trained abstractive summarization model, to specifically summarize figure-referencing paragraphs (*e.g.*, "Figure 3 shows...") into figure captions. Experiments on large-scale arXiv figures show that our method outperforms prior vision methods in both automatic and human evaluations. We further conducted an in-depth investigation focused on two key challenges: *(i)* the common presence of low-quality author-written captions and *(ii)* the lack of clear standards for good captions. Our code and data are available at: https://github.com/Crowd-AI-Lab/Generating-Figure-Captions-as-a-Text-Summarization-Task.

## 1 Introduction

In scientific documents, effective figure captions help readers understand complex figures like bar charts, line charts, or pie charts. These captions describe the images and often include necessary context from the document's full text (Durbin, 2004). Unfortunately, even published papers often have poorly-written captions. As per our analysis (Section 8.2), around 53.88% of line charts' captions in arXiv cs.CL papers are found to be unhelpful for NLP readers. Automatic caption generation could aid paper writers by providing good starting captions that can be refined for better quality.

Previous research typically approached figure caption generation as a *vision-to-language* task, *i.e.*,

creating captions based on the image. For instance, Hsu et al. (2021) used an end-to-end approach with CNN+RNN structures, which extracted feature representation from the image and converted it into caption text. Qian et al. (2021) took a slightly different approach: first understanding what is in the image, pulling out key information, and then using a preset template to create the caption. However, although achieving some success in synthetic data (Kahou et al., 2017; Kafle et al., 2018; Chen et al., 2020a; Zhu et al., 2021), these approaches often struggled to caption real-world figures. For example, Hsu et al. (2021)'s end-to-end approach, trained and tested using arXiv figures, achieved a BLEU-4 score of only 2.91.

In this paper, we argue that figure captioning in scientific documents can be more effectively tackled as a **text-summarization task**: The caption can be generated by **summarizing the paragraphs mentioning the figure** (as shown in Figure 1.) Scientific figures typically come with extensive text in the scientific document that can aid caption generation. Our analysis (Section 5) shows that, in arXiv, over 75% of words in figure captions can be aligned with the words in the paragraphs referencing those figures, which motivates our approach. The automatic evaluation shows that summarizing paragraphs referencing the figures results in better captions than prior vision-based methods. In a human evaluation by external domain experts, our best-performing model's captions were preferred over the original captions 46.67% of the time.

We further conducted an in-depth investigation focused on two key challenges: *(i)* the common presence of low-quality author-written captions and *(ii)* the lack of clear standards for good captions. Surprisingly, 53.88% of the author-written captions in our sample was deemed unhelpful. This has implications for the design of future captioning systems, underscoring the influence of data quality on captioning performance.

*\*Equal contribution.*

80

Figure 1: Figure captioning can be addressed as a text-summarization task. The figure's caption can be generated by summarizing the paragraphs mentioning the figure. The caption is generated by the model Pegasus$_{P+O+B}$. The example shown in this figure is extracted from the paper (Doulaty et al., 2015).

## 2 Related Work

Prior figure captioning works can be broadly categorized into two approaches: caption generation *(i)* based on the image of the figure or *(ii)* based on the data chart underlying the figure.

Earlier image-based approaches focused on automated image understanding, which involved parsing images to extract the figure's key attributes and converting parsed data into captions, *e.g.*, using predefined templates (Kahou et al., 2017; Kafle et al., 2018; Methani et al., 2020; Qian et al., 2021; Siegel et al., 2016). Recently, with the advance of deep learning, more works are adopting an end-to-end paradigm, generating captions straight from the neural representations of images (Mahinpei et al., 2022; Pelka et al., 2021; Hsu et al., 2021; Chen et al., 2019; Kantharaj et al., 2022; Chen et al., 2020a). Our work contrasts with prior studies by focusing on text to generate captions instead of visuals. To the best of our knowledge, no existing figure-caption datasets explicitly contain the figures' accompanying documents (Pelka et al., 2021; Hsu et al., 2021; Chen et al., 2019), as this task has generally been approached as a vision task. Most recently, a knowledge-augmented image captioning method that uses both image and text data was introduced (Yang et al., 2023), suggesting the potential of using text from documents.

Some approaches generate captions using the underlying tabular data of a figure rather than the figure's image. Earlier approaches often employed rule-based techniques (Corio and Lapalme, 1999; Demir et al., 2008; Fasciano and Lapalme, 1996; Mittal et al., 1998), while newer ones favor learning-based methods (Barzilay and Lapata, 2005; Wiseman et al., 2017; Moraes et al., 2014; Zhu et al., 2021; Kantharaj et al., 2022; Obeid and Hoque, 2020; Reiter et al., 2005; Parikh et al.,

2020; Chen et al., 2020b; Gong et al., 2019; Su et al., 2021; Chen et al., 2020c). Despite these approaches' ability to utilize tabular and meta data, they necessitate access to the figure's raw data. Contrarily, our work uses the rich textual information in scientific documents to generate captions.

## 3 Problem Statement and Terminology

A document $D$ contains $n$ figures, $F_1$ to $F_n$, where $F_i$ has a caption $C_i$ that was written by the document author. In document $D$, $j$ sentences, $M_{i,1}$ to $M_{i,j}$, explicitly mention $F_i$ (*e.g.*, "As shown in $F_i$..."). The objective of this work is to automatically generate a high-quality caption, $C_i'$, for figure $F_i$ using only its mentions ($M_{i,1}$ to $M_{i,j}$) and the surrounding text of the mentions in document $D$.

In the rest of the paper, we use these terms:

- A **"Mention"** refers to a sentence in a document that explicitly mentions the target figure, *e.g.*, "As shown in Figure 6..." If there are multiple Mentions, the first Mention is referred to.

- A **"Paragraph"** refers to a section of text containing a Mention. In this work, the boundaries of a Paragraph are determined by the <p> tag produced by PDF parsing.

- Sentences near a Mention may contain relevant information, so we extracted $n$ preceding sentences and $m$ following sentences to form the **"Window[n, m]"** text snippet. For instance, "Window[1, 2]" refers to a snippet of four sentences, including one preceding sentence, the Mention sentence, and two following sentences.

- An **"OCR"** refers to the textual information (*e.g.*, legends, labels, etc.) extracted from the image, by optical character recognition (OCR) software.

| Source | Random | | Mention | | Paragraph | | OCR | Window[0, 1] | | Window[1, 0] | | Window[1, 1] | | Window[2, 2] | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | S | P | - | +OCR | - | +OCR | - | - | +OCR | - | +OCR | - | +OCR | - | +OCR |
| **Caption** | 35.23 | 44.43 | 53.43 | 60.16 | 75.19 | **76.68** | 34.75 | 60.85 | 65.43 | 59.09 | 64.19 | 65.20 | 68.73 | 69.09 | 71.77 |
| **Source** | 32.52 | 19.52 | **39.51** | 18.78 | 12.53 | 9.39 | 20.79 | 30.49 | 17.19 | 32.40 | 17.33 | 25.10 | 15.55 | 19.84 | 13.45 |

Table 1: Macro coverage rates (percentage) between captions and relevant texts (S: Sentence and P: Paragraph). Caption coverage gives the percentage of words in the caption that can be found in the source texts and vice versa (punctuation and stop words are excluded.) The results show that 76.68% of the words in captions could be found in Paragraph+OCR, motivating us to generate captions by text summarization.

## 4 Dataset

Before diving into our experiments and analyses, we first describe the dataset upon which our study is grounded. Our results are based on a scientific figure caption dataset, SCICAP, and several pre-processing steps to fit it into our workflow.

SCICAP is a dataset that contains over 416,000 line charts and captions extracted from more than 290,000 arXiv papers (Hsu et al., 2021). It was one of the first large-scale figure-captioning datasets based on real-world scientific figures. However, it does not contain the paragraphs that mention the figure. To address this, we downloaded all the PDF files of the original arXiv papers used in SCI-CAP and extracted all the Mentions and Paragraphs as outlined in Section 6.1. Detailed information on preprocessing, including the dataset resplit and OCR extraction, are described in Appendix B.

## 5 Motivating Analysis

To understand the correlation between mentions and captions, we performed a series of analyses using the data described in Section 4. Specifically, we investigated the extent to which the words in the figure captions are represented in the corresponding figure-mentioning paragraphs. We used awesome-align (Dou and Neubig, 2021) to obtain the alignment between the source texts (mentions, paragraphs, and OCRs) and captions. Awesome-align compared the similarity of the words' contextual embeddings and assigned an alignment between words if the similarity passed a threshold. We used SciBERT (Beltagy et al., 2019) to obtain contextual embeddings and softmax threshold = 0.99 to reduce false alignments.

After obtaining the alignments, we computed what percentage of information in the caption could be found in the source texts. The results shown in Table 1 indicate that **76.68% of the caption's information could be found in Paragraph and OCR,** motivating us to generate figure captions by

summarizing Paragraph. We also observed that a randomly selected sentence and paragraph from the same paper can cover 35.23% and 44.43% of the caption, respectively, showing that there was some generic information-sharing across the paper. We also conducted a study using the exact overlapping (*i.e.*, BLEU score) in Appendix A.

## 6 Generating Figure Captions as a Text Summarization Task

Figure 1 overviews the proposed pipeline. This section describes each step of the pipeline.

### 6.1 Extracting Mentions and Paragraphs

The system first extracts Mentions and their associated Paragraphs (as defined in Section 3.) In this paper, we used Grobid (kermitt2, 2022), a publicly-available tool for converting PDF files into structured XML documents, to extract plain text from the paragraphs (including the <p> tags) in each paper. This plain text was then segmented into sentences using BlingFire (microsoft, 2022). We developed regular expressions to identify sentences mentioning specific figures. For instance, sentences such as "As shown in Figure 6, ..." were first identified and then linked to Figure 6. To assess the performance of these regular expressions, we conducted a manual evaluation of 300 samples from our experimental dataset. The results showed a high level of precision (99.58%) and recall (94.44%).

### 6.2 Generating Captions Using Text Summarization Models

As shown in Figure 1, our system then automatically summarizes all the extracted Mentions (or Paragraphs) into a figure caption. In this work, we used PEGASUS, an abstractive summarization model (Zhang et al., 2020), and fine-tuned it on our dataset. Five Pegasus models, $Pegasus_M$, $Pegasus_P$, $Pegasus_O$, $Pegasus_{M+O}$,

and Pegasus$_{P+O}$, were trained utilizing five distinct input combinations, including *(i)* **Mention,** *(ii)* **Paragraph,** *(iii)* **OCR output of the target figure image,** *(iv)* **Mention+OCR, and** *(v)* **Paragraph+OCR.** Pegasus$_{P+O}$ encompasses the most of relevant information in the document and thus is expected to yield the optimal summary.

Additionally, we built Pegasus$_{P+O+B}$, a specialized version of the model designed to be trained on a subset of higher-quality captions, *(vi)* **Paragraph+OCR-Better**. Given the absence of reliable automated ways to assess the quality of captions, we followed a guideline from previous studies indicating that longer captions enhance reader comprehension (Hartley, 2003; Gelman et al., 2002). We trained the model using captions with 30 or more tokens. The average caption length was 26.8 tokens, so we set 30 tokens as the threshold. The training was performed using Paragraph+OCR inputs.

We identified two major challenges in generating captions for scientific figures in real-world scenarios. We discuss these challenges in the following subsections, with an in-depth analysis in Section 8.

### 6.2.1 Challenge 1: Addressing Unreliable Quality of Real-World Data

Low-quality captions often occur in scholarly articles. Our analysis (see Section 8.1) showed that 50% of line charts' author-written captions in arXiv cs.CL papers were deemed unhelpful by domain experts. The impact of this unreliable data quality is that developers could train and test captioning models with unhelpful captions. The lack of automatic methods for evaluating caption quality makes it hard to identify suitable training examples and eliminate poor ones. To address this issue, we included Pegasus$_{P+O+B}$ that was trained on longer captions, which is suggested by literature to be more helpful to readers (Hartley, 2003; Gelman et al., 2002). To account for low-quality test data, we conducted both human and automatic evaluations. The data quality of figure captions was analyzed and is presented in Section 8.2.

### 6.2.2 Challenge 2: Defining a Clear Standard for "Good" Figure Captions

The deeper issue is the lack of a set of well-defined and actionable criteria for determining the usefulness of a figure caption. Although there are guidelines for writing effective scientific figure captions (Rougier et al., 2014; Biegel and Kamat,

2019), their translation into algorithmic models can be challenging. From a modeling standpoint, the lack of a clear goal presents a challenge, as it is uncertain what to optimize for once fluency has been achieved. In this paper, we focus on demonstrating the feasibility of generating captions via text summarization. Although we did not incorporate specialized goals in the model, we examine the criteria for a "good" caption in Section 8.2.

## 7 Experimental Results

**A Simple Baseline: Using Extracted Mentions as Captions.** Motivated by our information overlap study (Section 5), we created the **Reuse** baselines. These baselines simply repurpose portions of the input text as the prediction.

**Vision-to-Language Baselines.** The vision-to-language generation treated this task as an image-captioning task that took the scientific figure image as input and generated a text to describe it. We compared two vision-to-language models as baselines. First, we built a sequence-to-sequence model by combining BEiT (Bao et al., 2022) and GPT-2 (Radford et al., 2019). We also selected the TrOCR (Li et al., 2021) model, a transformer-based sequence-to-sequence model pre-trained for OCR tasks. Compared to image encoders like ViT (Dosovitskiy et al., 2021) and BEiT (Bao et al., 2022), which were trained on photos, OCR models trained on printed and handwritten documents align more closely with the scientific paper domain. All figures from SCICAP (106,391 training samples) were used for training since no mentions were required.

**Experimental Setup.** A total of 14 methods were included for comparison: six reuse baselines with six input variations (M, P, W[0, 1], W[0, 2], W[1, 1], and W[2, 2]); five text summarization models with five inputs (M, M+O, P, P+O, and O); one text summarization model using P+O with controlled data quality; and two vision-to-language models (BEiT+GPT-2 and TrOCR). Note that we use subscripts of M, P, W, O, B to denote the input features: Mention, Paragraph, Window, OCR, and Better data quality, respectively. The model training details and decoding configuration are provided in Appendix C.

### 7.1 Automatic Evaluation Results

**Conventional Automatic Evaluation.** We used F1 of ROUGE-1, ROUGE-2, ROUGE-L (Lin,

| Model | Feature | Length | Rouge-1 (F1) | | Rouge-2 (F1) | | Rouge-L (F1) | | MoverScore | | BERTScore | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Score | Norm | Score | Norm | Score | Norm | Score | Norm | Score | Norm |
| **Reuse** | **M** | 33.2 | .291 | 1.346 | .139 | 1.790 | .239 | 1.401 | .535 | 1.023 | .628 | 1.064 |
| | **P** | 238.3 | .171 | 1.042 | .089 | 1.006 | .134 | 1.030 | .503 | 1.004 | .567 | 1.008 |
| | **W[0, 1]** | 50.3 | .281 | 1.216 | .132 | 1.509 | .224 | 1.273 | .529 | 1.016 | .620 | 1.048 |
| | **W[0, 2]** | 68.0 | .259 | 1.129 | .123 | 1.341 | .205 | 1.186 | .524 | 1.013 | .611 | 1.034 |
| | **W[1, 1]** | 67.8 | .266 | 1.156 | .124 | 1.346 | .204 | 1.183 | .524 | 1.012 | .613 | 1.037 |
| | **W[2, 2]** | 98.7 | .235 | 1.082 | .112 | 1.179 | .180 | 1.105 | .517 | 1.007 | .600 | 1.020 |
| **Pegasus** | **M** | 12.2 | .321 | 1.898 | .153 | 2.907 | .283 | 1.971 | .553 | 1.065 | .654 | 1.158 |
| | **M+O** | 12.8 | .331 | 1.909 | .161 | 2.945 | .292 | 1.993 | .556 | 1.071 | .661 | 1.166 |
| | **P** | 14.0 | .374 | 2.067 | .205 | 3.507 | .334 | 2.201 | .570 | 1.095 | .682 | 1.196 |
| | **P+O** | 14.0 | **.381** | **2.106** | **.212** | **3.635** | **.340** | **2.242** | **.571** | **1.097** | **.685** | **1.202** |
| | **P+O+B** | 38.3 | .321 | 1.452 | .154 | 1.916 | .265 | 1.537 | .546 | 1.044 | .639 | 1.082 |
| | **O** | 12.1 | .133 | 0.789 | .026 | 0.495 | .119 | 0.828 | .518 | 0.998 | .561 | 0.993 |
| **TrOCR** | **Figure** | 10.0 | .220 | 1.464 | .073 | 1.653 | .195 | 1.502 | .534 | 1.033 | .610 | 1.096 |
| **BEiT+GPT2** | | 15.8 | .164 | 0.864 | .042 | 0.666 | .144 | 0.917 | .529 | 1.013 | .592 | 1.031 |

Table 2: Task Performance with the **best** and <u>second-best</u> results highlighted. Pegasus$_{P+O}$, the text-summarization model with all available information (Paragraph+OCR), performed the best in all four metrics. Pegasus$_{P+O+B}$, the model trained with better captions, however, got lower scores.



Figure 2: The relationship between average text length and ROUGE-2 score (left: original ROUGE-2; right: normalized ROUGE-2). The random baseline in the left figure shows that text lengths and scores are not independent. For example, when the predicted text is shorter than 50 tokens, predicting longer texts generally results in a higher ROUGE-2 score. The normalized scores indicate the proposed system's performance gain over the random baseline of the same length. Pegasus$_{P+O+B}$ and Reuse$_M$ get closer to TrOCR after normalization, suggesting the need for normalization for accurate interpretation of results.

2004; Nallapati et al., 2016), MoverScore (Zhao et al., 2019), and BERTScore for automatic evaluation. When computing ROUGE scores using rouge-score (google research, 2022), we turned all text into lower case and stem words. As both MoverScore and BERTScore are based on the semantic similarity, we obtained contextual embeddings from SciBERT (Beltagy et al., 2019).

**Automatic Evaluation with Normalization Over Caption Length.** ROUGE F1 tends to favor longer texts within a certain length, leading to a skewed comparison where models generating longer texts receive higher scores (Sun et al., 2019). We followed Sun et al. (2019)'s approach of normalizing the scores with the corresponding random baseline that generates texts of the same length.

$$Score_{normalized} = \frac{Score}{Random(length)} \quad (1)$$

where $length$ is the average length of the texts generated by the target system. We estimated $Random(length)$ by applying linear interpolation on several (length, random score) pairs. The (length, random score) pairs were obtained by randomly selecting a certain number of sentences (1, 2, ..., 10 sentences) from the input paragraph as the prediction. To get random scores of texts shorter than a single sentence (around 30 tokens), we truncated sentences to the desired length (4, 6, ..., 30 tokens). For each length setting, we ran 10 different random seeds and reported the average. The Random line in Figure 2-Left shows the behavior of ROUGE-2 favoring longer texts within 50 tokens.[1] The normalized scores, as shown in Figure 2-Right, clearly indicate the superiority of our

---

[1] Similar trends for ROUGE-1, ROUGE-L, MoverScore, and BERTScore are included in Appendix E.

proposed model over the random baseline.

Table 2 shows the normalized automatic evaluation results. Overall, **Pegasus$_{P+O}$, the text-summarization model with all available information (Paragraph+OCR), achieved the best performance in all three metrics.** Pegasus$_{P+O+B}$, the model using the same information but trained on a better subset of captions (Paragraph+OCR-Better), did not perform well. We hypothesized this was due to half of the test data comprising poor captions (refer to Section 8.2). This was validated by examining performance shifts in different quality beams (Section 8.1) and conducting a human evaluation (Section 7.2). Meanwhile, Reuse$_M$, the Reuse baseline with Mention, outperforms other Reuse baselines. Its performance declined as context sizes grew and shifted.

## 7.2 Human Evaluation Results

**Pilot MTurk Study to Select Top Models.** Before the main human evaluation, we ran a pilot study on Amazon Mechanical Turk (MTurk) to identify any apparently underperforming baselines for exclusion in the final study, simplifying the main human evaluations. In this study, we asked MTurk workers to carefully read a figure and select the *worst* figure caption among *(i)* TrOCR, *(ii)* Pegasus$_{P+O}$, *(iii)* Pegasus$_{P+O+B}$, and *(iv)* ground-truth caption. Ninety figures without errors were randomly sampled from our annotated set (*i.e.*, figures from cs.CL arXiv papers in Section 8.2) for the study. For each of the figures, we recruited 20 MTurk workers to judge.[2] We report the number of majority votes (when tied, we counted all captions with the highest votes as the worst) and the average number of votes in Table 3. Results indicated that TrOCR's caption won the majority vote 41 out of 90 times, with its average vote count significantly exceeding others. Hence, we excluded TrOCR from our formal human evaluation.

**Main Human Evaluation with Domain Experts.** Three Ph.D. students with NLP backgrounds (who are not coauthors) were recruited as human judges, as it is hard for those without basic domain understanding to evaluate captions. This study has been approved by the IRB office of the authors' institute. The same 90 figures used in the pilot

[2] Four MTurk qualifications were used: Locale (US Only), HIT Approval Rate ($\geq$98%), Number of Approved HITs ($\geq$3000), and the Adult Content Qualification. The payment for each task was set to 0.09 (hourly wage = $10 dollars).

| n = 90 | #Maj. Votes↓ | Avg. Votes↓ | T-Test over Avg. Votes | | |
|---|---|---|---|---|---|
| | | | Peg$_{P+O}$ | Peg$_{P+O+B}$ | Caption |
| **TrOCR** | 41 | 5.99 | <.001*** | .006** | .001** |
| **Peg$_{P+O}$** | 20 | 4.54 | - | .253 | .973 |
| **Peg$_{P+O+B}$** | 24 | 4.93 | - | - | .318 |
| **Caption** | 19 | 4.53 | - | - | - |

Table 3: The result of the pilot Mturk study. When tied, all captions with the highest votes were counted as the worst for #Majority votes. TrOCR is significantly worse than other approaches when rated by crowd workers.

| n = 90 | Avg. Ranking↓ | T-Test on Avg. Ranking | |
|---|---|---|---|
| | | Peg$_{P+O+B}$ | Caption |
| **Peg$_{P+O}$** | 2.152 | .016* | .015* |
| **Peg$_{P+O+B}$** | 1.930 | - | .923 |
| **Caption** | **1.919** | - | - |

Table 4: Average ranking of the human evaluation. Pegasus$_{P+O+B}$ was rated significantly better than Pegasus$_{P+O}$ and was at the same level as the ground-truth caption.

MTurk study were used again. We asked the human judges to compare each figure's *(i)* Pegasus$_{P+O}$, *(ii)* Pegasus$_{P+O+B}$, and *(iii)* ground-truth caption. The judges were asked to rank the captions based on how strongly they agreed with this statement: "When I read the paper, this caption can help me understand the message that the figure tries to convey." Figure 5 (see Appendix D) shows the interface the human judges used.

Table 4 shows the results of average ranking (from 1 to 3). Overall, **the ground-truth caption and Pegasus$_{P+O+B}$ were ranked similarly** (1.919 vs. 1.930 with p-value = 0.923). **Humans also favored Pegasus$_{P+O+B}$ over Pegasus$_{P+O}$ significantly** (1.919 vs. 2.152 with p-value = 0.016). This supports our heuristic for automatically determining caption quality based on length and aligns with previous findings that longer captions improve reader comprehension (Hartley, 2003; Gelman et al., 2002). However, we found that the task of caption ranking poses a challenge, as evidenced by the lower correlations between raters, with Kendall's tau values of 0.133, 0.148, and 0.274, and Spearman's rho values of 0.128, 0.156, and 0.317. This highlights the complexity of the task and suggests that scaling human evaluation across domains might be difficult. Different preferences over captions, such as length, could lead to lower agreement among raters.

|  | Agree | Disagree | Total | Agree Percentage |
|---|---|---|---|---|
| **Helpfulness** | 184 | 215 | 399 | 46.12% |
| **Image-Text** | 338 | 61 | 399 | 84.71% |
| **Visual-Desc** | 64 | 335 | 399 | 16.04% |
| **Takeaway** | 74 | 325 | 399 | 18.55% |

Table 5: Results of the manual annotation. More than 50% of the captions were annotated as unhelpful. (Out of the initial 438 figure captions, we excluded those with extraction or classification errors, *e.g.*, incomplete images, leaving us with only 399 captions.)

## 8    In-Depth Analysis

We conducted an in-depth investigation focused on two key challenges: *(i)* the common presence of low-quality author-written captions and *(ii)* the lack of clear standards for good captions.

**Quality Annotation Procedure.**    We manually annotated 438 captions in the Computation and Language domain (cs.CL) from the test set. Figure 6 (see Appendix D) shows the interface we used, in which the title, abstract, and PDF file of the paper were shown alongside the target figure's image, caption, and questions. For each caption, we asked the annotators (coauthors) to rate four aspects using a five-point Likert scale:

- **Image-Text.** The caption included named entities or important words/numbers in the figure (*e.g.*, title, legends, labels, etc.).
- **Visual-Description.** The caption included some visual characteristics of the figure (*e.g.*, color, shape, trend, etc.).
- **Takeaway.** The caption explicitly stated the high-level takeaway message or the conclusion that the figure attempted to convey.
- **Helpfulness.** "The caption helped me understand the message that the figure attempted to convey".

The annotated data was consolidated by grouping "Strongly Agree" and "Agree" as "[Agree]" and grouping "Neutral", "Disagree", and "Strongly Disagree" as "[Disagree]". The results of this consolidation are presented in Table 5.

### 8.1    Challenge 1: Addressing Unreliable Quality of Real-World Data

Table 5 shows that **over 50% of the author-written captions in arXiv cs.CL papers were deemed unhelpful**. High unhelpful caption proportion may skew evaluation results, particularly for automatic evaluations that compare generated text to human-written captions. To address this, we

|  | #Sample | $\text{Peg}_{P+O}$ | $\text{Peg}_{P+O+B}$ | Caption |
|---|---|---|---|---|
| **Helpful** | 55 | 2.176 | 1.970 | **1.855** |
| **Unhelpful** | 35 | 2.114 | **1.867** | 2.019 |

Table 6: Human ranking results (lower is better) on helpful and unhelpful beams. Pegasus$_{P+O+B}$ received better rankings in the unhelpful beam.

evaluated models on different quality beams using the 399 annotated figure captions shown in Table 5. The captions were divided into the "helpful beam" (184 captions rated [Agree]) and the "unhelpful beam" (215 captions rated [Disagree]).

**Automatic Evaluation Over Beams of Different Quality.**    To validate the effect of low-quality captions, we re-performed the automatic evaluation for the helpful and unhelpful beam sets. Figure 3 shows the Normalized ROUGE-2 and MoverScore scores for each model in the helpful and unhelpful beam sets.[3] Most models performed better in the unhelpful beam, except Pegasus$_{P+O+B}$, which had better scores in the helpful beam. Pegasus$_{P+O+B}$ was trained on captions with more than 30 tokens. This result suggests that improving training data quality, such as by using only longer captions, can positively impact the model's behavior and result in a better generation of helpful captions.

**Human Evaluation Over Beams of Different Quality.**    We also re-evaluated human scores for both the helpful and unhelpful beams. The human evaluation in Section 7.2 only covered 90 figures, with 55 in the helpful beam and 35 in the unhelpful beam. Table 6 shows the results. On average, **Pegasus$_{P+O+B}$ (1.867) was ranked better than author-written captions (2.019) in the unhelpful beam, in which machine-generated captions were preferred by human judges 22 out of 35 times.** The results suggest that, with careful training data quality control, when author-written captions are not very helpful, machines could potentially generate better captions.

### 8.2    Challenge 2: What Constitutes a Good Figure Caption?

We calculated Pearson correlations (Rodgers and Nicewander, 1988) among the four aspects using raw five-point Likert ratings. The results are shown in Table 7. **The highest correlation was found between Takeaway and Helpfulness, suggesting that a high-quality caption accurately captures**

---

[3]In addition, ROUGE-1, ROUGE-L, and BERTScore scores can be found in Figure 12 in Appendix E.

Figure 3: Normalized ROUGE-2 and MoverScore for beams of different quality. Most of the generation models performed better in the unhelpful beam, suggesting that output may be more similar to bad captions. Only the model trained with **better** captions learned to generate good captions by showing a better score in the helpful beam.

|  | Image-Text | Visual-Desc | Takeaway | Length |
|---|---|---|---|---|
| **Helpfulness** | 0.206 | **0.523** | **0.686** | <u>0.383</u> |
| **Image-Text** | - | 0.177 | 0.186 | 0.248 |
| **Visual-Desc** | - | - | **0.625** | **0.535** |
| **Takeaway** | - | - | - | **0.514** |

Table 7: Pearson correlations between different aspects. We used the row scores (five-point Likert scale) to compute the correlation. **Strong correlation** (≥0.5) and <u>medium correlation</u> (0.3 to 0.5) are highlighted. Helpfulness is highly correlated with Visual-Description and Takeaway and is moderately correlated with Length.

**the main message of the figure.** There were also strong correlations between Helpfulness, Visual-Description, and Takeaway, indicating that a good caption effectively conveys visual information and summarizes the main message. However, Table 5 shows that only 16.04% and 18.55% of the captions described the visual characteristics and the takeaway message, respectively.

A moderate correlation between Helpfulness and Length supports previous research findings that longer captions are generally more helpful for readers (Hartley, 2003; Gelman et al., 2002).

### 8.3 Caption Length Distribution

Throughout this work's development, the length of captions emerged as a consistent issue. Despite existing literature indicating the benefits of longer captions for readers (Hartley, 2003; Gelman et al., 2002), space limitations often leave authors with no option but to craft shorter captions. To shed some light on this aspect and offer insight for future research, we analyzed the lengths of both author-created and machine-generated captions. We used Kernel Density Estimate (KDE) plots to investigate the distribution of caption lengths across different models and domains. As shown in Figure 4a, the majority of models demonstrate a common peak at 10 tokens, while Pegasus$_{P+O+B}$ presents a significant deviation with a peak near 30 tokens. Figure 4b presents the distribution of helpfulness

| Information | Image-Text | Visual-Desc | Takeaway |
|---|---|---|---|
| **#Tokens** | 0.181 | 0.428 | 0.357 |
| **Percentage** | 0.099 | 0.279 | 0.210 |

Table 8: Correlations between the amount of missing information from Paragraph and the quality aspects. The missing information is related to visual descriptions and takeaway messages.

scores, derived from quality annotation data (see Section 8.2). Captions rated with a maximum helpfulness score of 5 show a peak at 35 tokens. We can also see a clear shift in caption length with higher scores. In Figure 4c, we dug into the top 10 category taxonomy from arXiv. This figure suggests that a higher portion of the captions in cs, math, stat, and eess are shorter (10 tokens); while the rest of the categories (cond-mat, quant-ph, q-bio, *etc*) have higher probabilities for longer captions. However, within the cs domain (Figure 4d), the top 10 subcategories do not show significant differences regarding caption length distribution.

## 9 Discussion

**Is Text Really All You Need?** Our results demonstrate that summarizing figure-mentioning paragraphs is sufficient to generate captions, as shown by the similar scores of Pegasus$_P$ and Pegasus$_{P+O}$ in Table 2. Adding OCR had limited impact. Furthermore, in a recent study of scientific figure captioning conducted by Yang *et al.* (Yang et al., 2023), the best-performing model only considered figure-mentioning paragraphs and OCR tokens– note that their OCR tokens were visual features– without taking the figure's imagery into account. These results raise an interesting question: Do we need visual information at all? What for? The token alignment study (Section 5) showed that 75.19% of the caption information could be found in the Paragraphs, meaning 24.81% of the information was missing. Understanding this missing information could help improve the models' performance.

(a) All the examined generative models.



(b) Helpfulness scores in 5-point Likert scale.



(c) Top 10 frequent arXiv categories.



(d) Top 10 frequent arXiv CS subcategories and cs.HC.

Figure 4: The KDE plot illustrates diverse caption length distributions among models, the relationship with helpfulness, and variations across arXiv categories. The sample numbers are included in the legend.

Thus, we calculated the correlation between the amount of missing information and three aspect ratings (image-text, visual-description, and takeaway) in the quality annotation data (Section 8.2). The missing information was quantified as the number or percentage of tokens without aligning to any tokens in figure-mentioning paragraphs. Table 8 demonstrates a positive correlation between the extent of missing information and visual descriptions and takeaway messages. This suggests that incorporating visual descriptions (*e.g.*, "dashed line," "red line") is key to enhancing performance by filling in the gaps in information not covered by the article's text. Furthermore, the strong correlation between Helpfulness and Visual-Description in Table 5 also indicates that including image information is necessary for writing good captions. It should be noted that OCR is only capable of capturing image texts (*e.g.*, labels, legends) and not visual element information (*e.g.*, "dashed line"). A promising future direction is developing a multimodal model that can effectively incorporate both image and text.

**What is the *Best* Length for Captions?** Our research indicates that filtering shorter captions can facilitate the generation of more helpful captions. However, the resulting captions tend to be longer than usual, as shown in the Pegasus$_{P+O+B}$ shift to the right in Figure 4a. This raises a question: Is it fair to compare short and long captions on usefulness, given that longer captions inherently contain more information? While our automatic evaluation addressed this by implementing length normalization, our human evaluations and quality annotations did not specifically instruct the annotators to

consider caption lengths. Nevertheless, we argue that even if we asked annotators to consider caption lengths while identifying helpful captions, the "ideal" caption length would differ among annotators due to multiple factors. For example, as shown in Figure 4c, the length distributions of captions vary across domains. The low inter-agreement from our human evaluation (see Section 7.2) also suggests that personal preferences could influence ideal caption length (Lundgard and Satyanarayan, 2021). Moreover, the ideal length could also be dictated by the context: writers might favor shorter captions due to page constraints, while readers might prefer longer but informative ones (Stokes and Hearst, 2022; Sun et al., 2019). To tackle this issue, a potential future direction could be enabling models to generate captions of diverse lengths to suit different users and contexts.

## 10 Conclusion and Future Work

This work presented a new perspective on automatic figure captioning, demonstrating that a language-based approach, *i.e.*, summarizing figure-referring paragraphs, can outperform conventional vision-based methods. Our analysis further showed many unhelpful captions in arXiv papers, highlighting data quality's impact on captioning performance. This work lays the groundwork for further research, including exploring new data selection, revision, and augmentation strategies to mitigate the effects of low-quality data, developing new evaluation methods, and creating more robust models that better handle noisy data. We also aim to expand the technology's scope to cover a wider variety of figures and article types.

## Acknowledgements

## Limitations

Although our proposed methods have been shown to be effective, we are aware of several limitations. First, our approach requires mentions in order to produce captions, but it is not always easy to automatically identify the mentions for a given figure in real-world data. There were 18.81% of figures in the original SciCap that did not have any identified mentions, which we excluded from this work. Many factors contributed to the gap, including errors caused by upstream components such as image extraction or image type classification (*e.g.*, table), unexpected figure index formats (*e.g.*, "Figure VIII", "Figure C·1","Fig.Fig. 4(b)"), PDF parsing errors, or the figure never being mentioned in the paper. Second, our method uses texts instead of images as the primary information source, so, naturally, it inherits all the constraints of text. Our method can not capture any visual element in the figure that the text never mentioned; it struggles when the text is poorly written. Finally, this paper focused on non-compound line charts in arXiv papers; the human evaluation only focused on NLP papers. More research is needed to examine the generalizability.

## Ethics Statement

We consider the proposed technology to impose little risk to readers, as it only summarizes what has already been presented in the paper. However, when the generated caption contains some inaccurate information, it could mislead readers. Furthermore, the proposed technology has the nature of neglecting visual content, which might have an impact on the accessibility of figure captions.

## References

Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. 2022. BEit: BERT pre-training of image transformers. In *International Conference on Learning Representations*.

Regina Barzilay and Mirella Lapata. 2005. Collective content selection for concept-to-text generation. In *Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing*, pages 331–338, Vancouver, British Columbia, Canada. Association for Computational Linguistics.

Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. SciBERT: A pretrained language model for scientific text. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3615–3620, Hong Kong, China. Association for Computational Linguistics.

Constance M. Biegel and Prashant V. Kamat. 2019. Ten tips for capturing figures with captions. *ACS Energy Letters*, 4(3):637–638.

Charles Chen, Ruiyi Zhang, Sungchul Kim, Scott Cohen, Tong Yu, Ryan Rossi, and Razvan Bunescu. 2019. Neural caption generation over figures. In *Adjunct Proceedings of the 2019 ACM International Joint Conference on Pervasive and Ubiquitous Computing and Proceedings of the 2019 ACM International Symposium on Wearable Computers*, UbiComp/ISWC '19 Adjunct, page 482–485, New York, NY, USA. Association for Computing Machinery.

Charles Chen, Ruiyi Zhang, Eunyee Koh, Sungchul Kim, Scott Cohen, and Ryan Rossi. 2020a. Figure captioning with relation maps for reasoning. In *2020 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1526–1534.

Wenhu Chen, Jianshu Chen, Yu Su, Zhiyu Chen, and William Yang Wang. 2020b. Logical natural language generation from open-domain tables. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7929–7942, Online. Association for Computational Linguistics.

Zhiyu Chen, Wenhu Chen, Hanwen Zha, Xiyou Zhou, Yunkai Zhang, Sairam Sundaresan, and William Yang Wang. 2020c. Logic2Text: High-fidelity natural language generation from logical forms. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2096–2111, Online. Association for Computational Linguistics.

Marc Corio and Guy Lapalme. 1999. Generation of texts for information graphics. In *Proceedings of the 7th European Workshop on Natural Language Generation EWNLG'99*, pages 49–58.

Seniz Demir, Sandra Carberry, and Kathleen F McCoy. 2008. Generating textual summaries of bar charts. In *Proceedings of the Fifth International Natural Language Generation Conference*, pages 7–15.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob

Uszkoreit, and Neil Houlsby. 2021. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*.

Zi-Yi Dou and Graham Neubig. 2021. Word alignment by fine-tuning embeddings on parallel corpora. In *Conference of the European Chapter of the Association for Computational Linguistics (EACL)*.

Mortaza Doulaty, Oscar Saz, and Thomas Hain. 2015. Data-selective transfer learning for multi-domain speech recognition. In *INTERSPEECH 2015, 16th Annual Conference of the International Speech Communication Association, Dresden, Germany, September 6-10, 2015*, pages 2897–2901. ISCA.

Charles G Durbin. 2004. Effective use of tables and figures in abstracts, presentations, and papers. *Respiratory care*, 49(10):1233–1237.

Massimo Fasciano and Guy Lapalme. 1996. Postgraphe: a system for the generation of statistical graphics and text. In *Eighth International Natural Language Generation Workshop*.

Andrew Gelman, Cristian Pasarica, and Rahul Dodhia. 2002. Let's practice what we preach: turning tables into graphs. *The American Statistician*, 56(2):121–130.

Li Gong, Josep Crego, and Jean Senellart. 2019. Enhanced transformer model for data-to-text generation. In *Proceedings of the 3rd Workshop on Neural Generation and Translation*, pages 148–156, Hong Kong. Association for Computational Linguistics.

google research. 2022. Python rouge implementation. https://github.com/google-research/google-research/tree/master/rouge.

James Hartley. 2003. Single authors are not alone: Colleagues often help. *Journal of Scholarly Publishing*, 34(2):108–113.

Ting-Yao Hsu, C Lee Giles, and Ting-Hao Huang. 2021. SciCap: Generating captions for scientific figures. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3258–3264, Punta Cana, Dominican Republic. Association for Computational Linguistics.

JaidedAI. 2022. Easyocr. https://github.com/JaidedAI/EasyOCR.

Zhuoxuan Jiang, Ziming Huang, Dong Sheng Li, and Xian-Ling Mao. 2019. Dialogact2vec: Towards end-to-end dialogue agent by multi-task representation learning. *CoRR*, abs/1911.04088.

Kushal Kafle, Brian Price, Scott Cohen, and Christopher Kanan. 2018. Dvqa: Understanding data visualizations via question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5648–5656.

Samira Ebrahimi Kahou, Vincent Michalski, Adam Atkinson, Ákos Kádár, Adam Trischler, and Yoshua Bengio. 2017. Figureqa: An annotated figure dataset for visual reasoning. *Sixth International Conference on Learning Representations Workshop*.

Shankar Kantharaj, Rixie Tiffany Leong, Xiang Lin, Ahmed Masry, Megh Thakkar, Enamul Hoque, and Shafiq Joty. 2022. Chart-to-text: A large-scale benchmark for chart summarization. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4005–4023, Dublin, Ireland. Association for Computational Linguistics.

kermitt2. 2022. Grobid. https://github.com/kermitt2/grobid.

Minghao Li, Tengchao Lv, Lei Cui, Yijuan Lu, Dinei Florencio, Cha Zhang, Zhoujun Li, and Furu Wei. 2021. Trocr: Transformer-based optical character recognition with pre-trained models. *arXiv preprint arXiv:2109.10282*.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Yixin Liu and Pengfei Liu. 2021. SimCLS: A simple framework for contrastive learning of abstractive summarization. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 1065–1072, Online. Association for Computational Linguistics.

Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In *International Conference on Learning Representations*.

Alan Lundgard and Arvind Satyanarayan. 2021. Accessible visualization via natural language descriptions: A four-level model of semantic content. *IEEE transactions on visualization and computer graphics*, 28(1):1073–1083.

A. Mahinpei, Z. Kostic, and C. Tanner. 2022. Linecap: Line charts for data visualization captioning models. In *2022 IEEE Visualization and Visual Analytics (VIS)*, pages 35–39, Los Alamitos, CA, USA. IEEE Computer Society.

Nitesh Methani, Pritha Ganguly, Mitesh M. Khapra, and Pratyush Kumar. 2020. Plotqa: Reasoning over scientific plots. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*.

microsoft. 2022. Blingfire. https://github.com/microsoft/BlingFire.

Vibhu O Mittal, Johanna D Moore, Giuseppe Carenini, and Steven Roth. 1998. Describing complex charts in natural language: A caption generation system. *Computational Linguistics*, 24(3):431–467.

Priscilla Moraes, Gabriel Sina, Kathy McCoy, and Sandra Carberry. 2014. Generating summaries of line graphs. In *Proceedings of the 8th International Natural Language Generation Conference (INLG)*, pages 95–98.

Ramesh Nallapati, Bowen Zhou, Cicero dos Santos, Çağlar Gulçehre, and Bing Xiang. 2016. Abstractive text summarization using sequence-to-sequence RNNs and beyond. In *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*, pages 280–290, Berlin, Germany. Association for Computational Linguistics.

Jason Obeid and Enamul Hoque. 2020. Chart-to-text: Generating natural language descriptions for charts by adapting the transformer model. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 138–147, Dublin, Ireland. Association for Computational Linguistics.

Ankur Parikh, Xuezhi Wang, Sebastian Gehrmann, Manaal Faruqui, Bhuwan Dhingra, Diyi Yang, and Dipanjan Das. 2020. ToTTo: A controlled table-to-text generation dataset. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1173–1186, Online. Association for Computational Linguistics.

Obioma Pelka, Asma Ben Abacha, Alba García Seco de Herrera, Janadhip Jacutprakart, Christoph M Friedrich, and Henning Müller. 2021. Overview of the ImageCLEFmed 2021 concept & caption prediction task. In *CLEF2021 Working Notes*, CEUR Workshop Proceedings, Bucharest, Romania. CEUR-WS.org.

Xin Qian, Eunyee Koh, Fan Du, Sungchul Kim, Joel Chan, Ryan A Rossi, Sana Malik, and Tak Yeon Lee. 2021. Generating accurate caption units for figure captioning. In *Proceedings of the Web Conference 2021*, pages 2792–2804.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Ehud Reiter, Somayajulu Sripada, Jim Hunter, Jin Yu, and Ian Davy. 2005. Choosing words in computer-generated weather forecasts. *Artificial Intelligence*, 167(1-2):137–169.

Joseph Lee Rodgers and W. Alan Nicewander. 1988. Thirteen ways to look at the correlation coefficient. *The American Statistician*, 42(1):59–66.

Nicolas P. Rougier, Michael Droettboom, and Philip E. Bourne. 2014. Ten simple rules for better figures. *PLOS Computational Biology*, 10(9):1–7.

Noah Siegel, Zachary Horvitz, Roie Levin, Santosh Divvala, and Ali Farhadi. 2016. Figureseer: Parsing result-figures in research papers. In *European Conference on Computer Vision*, pages 664–680. Springer.

Chase Stokes and Marti Hearst. 2022. Why more text is (often) better: Themes from reader preferences for integration of charts and text. *The NLVIZ workshop at IEEE Transaction on Visualization and Graphics conference*.

Yixuan Su, David Vandyke, Sihui Wang, Yimai Fang, and Nigel Collier. 2021. Plan-then-generate: Controlled data-to-text generation via planning. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 895–909, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Simeng Sun, Ori Shapira, Ido Dagan, and Ani Nenkova. 2019. How to compare summarizers without target length? pitfalls, solutions and re-examination of the neural summarization literature. In *Proceedings of the Workshop on Methods for Optimizing and Evaluating Neural Language Generation*, pages 21–29, Minneapolis, Minnesota. Association for Computational Linguistics.

Jie Wang and Xinyan Zhao. 2019. Theme-aware generation model for chinese lyrics. *arXiv preprint arXiv:1906.02134*.

Sam Wiseman, Stuart Shieber, and Alexander Rush. 2017. Challenges in data-to-document generation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2253–2263, Copenhagen, Denmark. Association for Computational Linguistics.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Jiacheng Xu, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020. Discourse-aware neural extractive text summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5021–5031, Online. Association for Computational Linguistics.

Zhishen Yang, Raj Dabre, Hideki Tanaka, and Naoaki Okazaki. 2023. Scicap+: A knowledge augmented dataset to study the challenges of scientific figure captioning.

Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter Liu. 2020. Pegasus: Pre-training with extracted gap-sentences for abstractive summarization. In *International Conference on Machine Learning*, pages 11328–11339. PMLR.

Wei Zhao, Maxime Peyrard, Fei Liu, Yang Gao, Christian M. Meyer, and Steffen Eger. 2019. MoverScore: Text generation evaluating with contextualized embeddings and earth mover distance. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 563–578, Hong Kong, China. Association for Computational Linguistics.

Ming Zhong, Pengfei Liu, Yiran Chen, Danqing Wang, Xipeng Qiu, and Xuanjing Huang. 2020. Extractive summarization as text matching. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6197–6208, Online. Association for Computational Linguistics.

Jiawen Zhu, Jinye Ran, Roy Ka-Wei Lee, Zhi Li, and Kenny Choo. 2021. AutoChart: A dataset for chart-to-text generation task. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)*, pages 1636–1644, Held Online. INCOMA Ltd.

# Models of reference production: How do they withstand the test of time?

**Fahime Same**[♡]**, Guanyi Chen**[♠]**,** and **Kees van Deemter**[♠]

[♡]Department of Linguistics, University of Cologne

[♠]Department of Information and Computing Sciences, Utrecht University

`f.same@uni-koeln.de, g.chen@ccnu.edu.cn, c.j.vandeemter@uu.nl`

## Abstract

In recent years, many NLP studies have focused solely on performance improvement. In this work, we focus on the linguistic and scientific aspects of NLP. We use the task of generating referring expressions in context (REG-in-context) as a case study and start our analysis from GREC, a comprehensive set of shared tasks in English that addressed this topic over a decade ago. We ask what the performance of models would be if we assessed them (1) on more realistic datasets, and (2) using more advanced methods. We test the models using different evaluation metrics and feature selection experiments. We conclude that GREC can no longer be regarded as offering a reliable assessment of models' ability to mimic human reference production, because the results are highly impacted by the choice of corpus and evaluation metrics. Our results also suggest that pre-trained language models are less dependent on the choice of corpus than classic Machine Learning models, and therefore make more robust class predictions.

## 1 Introduction

NLP research can have different aims. Some NLP research focuses on developing new algorithms or building practical NLP applications. Another line of NLP work constructs computational models that aim to explain human language and language use; this line of work has been dubbed *NLP-as-Science* (van Deemter, 2023). Among other things, NLP-as-Science demands that we ask ourselves to what extent NLP research findings generalise along a range of dimensions.

In addition to the practical applications of Referring Expression Generation (REG, Reiter, 2017), REG is also one of the typical tasks in NLP-as-Science, where REG algorithms are built to model and explain the reference production of human beings (Krahmer and van Deemter, 2012; van Deemter, 2016). In the computational linguis-tics and cognitive science community, REG can be divided into two distinct tasks: *one-shot REG*, finding a referring expression (RE) to single out a referent from a set, and *REG-in-context*, generating an RE to refer to a referent at a given point in a discourse.

In a classic setup, REG-in-context is often approached in two steps: The first is to decide on the form of an RE at a given point in the discourse, and the second is to decide on its content. Many researchers have been interested in the first sub-task, referential form selection: the task to decide which referential form (e.g., pronoun, proper name, description, etc.) an RE takes (McCoy and Strube, 1999; Henschel et al., 2000; Kibrik et al., 2016). Nearly 15 years ago, Belz et al. (2008) introduced the GREC shared tasks and a number of English REG corpora with two goals: (1) assessing the performance of computational models of reference production (Belz et al., 2009), and (2) understanding the contribution of linguistically-inspired factors to the choice of referential form (Greenbacker and McCoy, 2009b; Kibrik et al., 2016; Same and van Deemter, 2020).

15 years have passed since the GREC challenge was organised, and many new models and corpora have been proposed in the meantime (e.g., Castro Ferreira et al. (2018); Cunha et al. (2020), and Same et al. (2022)). We, therefore, decided that it was time to ask, in the spirit of NLP-as-Science, how well the lessons that GREC once taught our research community hold up when scrutinised in light of all these developments. In other words, we will investigate to what extent the findings from GREC can be *generalised* to other corpora and other models.

To this end, we pursue the following objectives: (1) We extend GREC by testing its REG algorithms not only on the GREC corpora but also on a corpus that was not originally considered and that has a different genre, namely the Wall Street Journal

(WSJ) portion of OntoNotes (Hovy et al., 2006; Weischedel et al., 2013); (2) We fine-tune pre-trained language models on the task of REG-in-context and assess them in the GREC framework.

In Section 2, we detail the GREC shared tasks and introduce the corpora used in GREC. Section 3 spells out our research questions. In Section 4 and Section 5, we introduce the algorithms and corpora that we use. Section 6 reports the performance of each algorithm on each corpus, followed by analyses in Section 7. Section 8 will discuss our findings and draw some lessons.

## 2 The GREC Shared Tasks

In this section, we summarise the GREC task, the corpora used by GREC, and its conclusions.

### 2.1 The GREC Task and its Corpora

According to Belz et al., "*the GREC tasks are about how to generate appropriate references to an entity in the context of a piece of discourse longer than a sentence*" (2009, p. 297). The main task was to predict the referential form, namely whether to use a pronoun, proper name, description or an empty reference at a given point in discourse.

The GREC challenges use two corpora, both created from the introductory sections of Wikipedia articles: (1) GREC-2.0 (henceforth MSR, as it was used in the GREC-MSR shared tasks of 2008 and 2009) consists of 1941 introductory sections of the articles across five domains (people, river, mountain, city, and country); and (2) GREC-People (henceforth NEG as it was used in the GREC-NEG shared task in 2009) contains 1000 introductory sections from Wikipedia articles about composers, chefs, and inventors. Here is an example from NEG:

(1) **David Chang** (born 1977) is a noted American chef. **He** is chef/owner of Momofuku Noodle Bar, Momofuku Ko and Momofuku Ssäm Bar in New York City. **Chang** attended Trinity College, where **he** majored in religious studies. In 2003, **Chang** opened **his** first restaurant, Momofuku Noodle Bar, in the East Village.

A key difference between MSR and NEG lies in their RE annotation practices. In MSR, only those REs that refer to the main topic of the article are annotated, while in NEG, mentions of all *human* referents are annotated. For instance, in a document about David Chang, MSR will only annotate

| Name | GREC ST | ALG | Acc |
|------|---------|-----|-----|
| UDel | MSR '09 | C5.0 | 77.71 |
| ICSI | MSR '09 | CRF | 75.16 |
| CNTS | MSR '08 | MBL | 72.61 |
| IS-G | MSR '08 | MLP | 70.78 |
| OSU | MSR '08 | MaxEnt | 69.82 |
| JUNLG | MSR '09 | Rule | 75.40 |

Table 1: An overview of the algorithms submitted to GREC. The first column contains the name of the respective algorithm. The column GREC ST presents the name of the MSR shared task to which the algorithm was submitted. The third column, ALG, lists the algorithms used, where abbreviations from top to bottom are C5.0 decision tree, conditional random field, memory-based learning, multi-layer perceptron, maximum entropy, and frequency-based rules. The fourth column, Acc, reports the original accuracy of the algorithms, as reported in Belz et al. (2009). Note that UDel, ICSI, and JUNLG were submitted to both the MSR '08 and MSR'09 shared tasks, and we only present the newest results here.

REs referring to David Chang, while NEG will include annotations for all human referents, including David Chang and others.

### 2.2 REG Algorithms Submitted to GREC

Various REG algorithms were submitted to the GREC challenges. These consist of feature-based ML algorithms: CNTS (Hendrickx et al., 2008), ICSI (Favre and Bohnet, 2009), IS-G (Bohnet, 2008), OSU Jamison and Mehay (2008) and UDel Greenbacker and McCoy (2009a), and an algorithm that mixes feature-based ML and rules: JUNLG (Gupta and Bandopadhyay, 2009). Table 1 presents the details of each model, including the ML method, and the original reported accuracy on MSR (cf. Belz et al. (2009) for details).

### 2.3 Feature Selection

The GREC Tasks were designed to find out *what kind of information is useful for making choices between different kinds of referring expressions in context* (Belz et al., 2009, p. 297). However, the original paper does not consider the factors that contributed to the RE choice in the systems submitted to GREC. In a follow-up study, Greenbacker and McCoy (2009b) conducted a feature selection study informed by psycholinguistics. They experimented with various feature subsets derived from their system, known as UDel, which had previously been submitted to the GREC. Additionally, they incorporated selected features from another

REG system, CNTS (Hendrickx et al., 2008), into their study. They show that features motivated by psycholinguistic studies and certain sentence construction features have a positive impact on the performance of REG models. Follow-up feature-selection studies including Kibrik et al. (2016) and Same and van Deemter (2020) also emphasise the contribution of factors such as recency and grammatical role to the choice of RE form.

## 3 Research Questions

15 years after the GREC shared tasks, we were curious to know to what extent the conclusions from GREC still "stand". We, therefore, came up with the following research questions.

In the first place, we are interested in *the impact of the choice of corpus on the performance of REG algorithms* ($\mathcal{R}_1$). GREC uses only the introductory part of Wikipedia articles (see Section 2), which represents only one genre of human language use. Considering that a good REG algorithm needs to model the general use of reference, a better evaluation framework should include texts from multiple genres. Therefore, we also include the WSJ corpus in the study (see Section 5 for more details) and conduct a correlation analysis to quantify how the choice of corpus impacts the evaluation results.

Second, previous studies suggested that classic machine learning (ML) based REG algorithms perform on par with most recent neural methods (Same et al., 2022). However, their study has three limitations: (1) they did not incorporate pre-trained language models (PLMs); (2) they focused on the surface forms of REs, which partly depend on the performance of surface realisation; (3) they did not assess the models based on the intuition that a model with good explanatory power should be less influenced by the choice of corpus. Therefore, we adopt PLMs to the task of REG-in-context (see Section 4 for more details) and investigate *how good is the explanatory power of PLM-based REG models compared to classic ML-based models* ($\mathcal{R}_2$) using the enhanced GREC framework.

Finally, as previously mentioned, one of the primary theoretical objectives of GREC was to computationally explore the contribution of factors that originate from linguistic studies to the choice of referential forms. It is reasonable to expect that such contributions may change depending on the choice of corpus. In this study, we conduct an importance analysis to investigate *whether the impor-*



Figure 1: Illustration of the PLM-based REG Algorithm.

*tance ranking of linguistic factors changes when we use different corpora* ($\mathcal{R}_3$).

## 4 REG Algorithms

In what follows, we introduce the REG algorithms that are considered in this study.

### 4.1 ML-based REG

For this study, we have narrowed our focus to feature-based ML algorithms that predict the type of RE. Consequently, we reconstruct five ML-based REG algorithms, namely UDel, ICSI, CNTS, IS-G, and OSU, along with their respective feature sets, while excluding JUNLG. Note that we implement CNTS slightly differently from Hendrickx et al. (2008). Concretely, Hendrickx et al. (2008) have mentioned that they have used the TiMBL package (Daelemans et al., 2007) for implementing the Memory Based Learning algorithm. Instead, we implemented the k-Nearest Neighbors algorithm. According to Daelemans et al. (2007), Memory Based Learning is the direct descendant of k-Nearest Neighbors. More information on the implementation of these models can be found in Appendix B.

### 4.2 PLM-based REG

Deep learning approaches have been used in many previous works on REG (Castro Ferreira et al., 2019; Cao and Cheung, 2019; Cunha et al., 2020; Chen et al., 2021). Different from previous work[1],

---

[1] Note that Chen et al. (2021, 2023) also leveraged a PLM, but did not fine-tune it. Instead, they used the word representations from the PLM as static inputs to an RNN and made predictions using the RNN.

we fine-tune PLMs on REG corpora in this study.

To fine-tune PLMs on REG corpora, we began by pre-processing each corpus using the same paradigm as described by Cunha et al. (2020). More precisely, each referent in a given document was replaced with its corresponding proper name. For example, all underlined REs in Example (1) were replaced by "David Chang". Subsequently, as depicted in Figure 1, we fed the data into a PLM, and, for each referent (e.g., "David Chang" ), we extracted the representations of its first token and its last token and summed them. The final representations were then sent to a fully connected layer for predicting the RE forms. In this study, we use BERT and RoBERTa (see section 6.1 for more details).

## 5 REG Corpora

In the following, we explain the corpora used in this work. These corpora are English-language corpora.

### 5.1 The MSR and NEG Corpora

In the current study, we only use the articles from the training sets of these corpora (see the number of documents in Table 2). Following the same approach as Castro Ferreira et al. (2018), we created a version of the GREC corpora for the End-to-end (E2E) REG modelling. For the classic ML models, we reproduced the models using the feature sets from the studies mentioned in Section 2.2.

### 5.2 The WSJ Corpus

As mentioned earlier, the WSJ portion of the OntoNotes corpus (Weischedel et al., 2013) is our third data source.[2] We use the version of the corpus that Same et al. (2022) developed for E2E REG modeling.[3] Since empty pronouns are not annotated in WSJ, we decided to also exclude them from the two GREC corpora and focus on a 3-label classification task. The labels considered in this study are *pronoun*, *description*, and *proper name*. Table 2 presents a detailed overview of these corpora.

**Data Splits.** We have made a document-wise split of the data. We split the WSJ data in accordance with the CoNLL 2012 Shared Task (Pradhan et al., 2012). Our WSJ training, development, and test sets contain 20275, 2831, and 2294 samples,

---

[2]We used Ontonotes 5.0 licensed by the Linguistic Data Consortium (LDC) https://catalog.ldc.upenn.edu/LDC2013T19.

[3]Note that WSJ was used in Same et al. (2022), but no corpus analysis or comparison was provided.

|  | MSR | NEG | WSJ |
|---|---|---|---|
| number of documents | 1655 | 808 | 582 |
| word/doc (mean) | 148 | 129 | 530 |
| sent/doc (mean) | 7.1 | 5.8 | 25 |
| par/doc (mean) | 2.3 | 2.2 | 10.8 |
| referent/doc (mean) | 1 | 2.6 | 15 |
| number of RE | 11705 | 8378 | 25400 |
| description % | 13.84% | 4% | 38.29% |
| proper name % | 38.09% | 40.79% | 34.57% |
| pronoun % | 41.79% | 48.75% | 27.14% |
| empty % | 6.28% | 6.47% | - |

Table 2: Comparison of the MSR, NEG, and WSJ corpora in terms of their length-related characteristics and distribution of REs. *Doc*, *sent* and *par* stands for *documents*, *sentences* and *paragraphs*.

respectively. We did an 85-5-10 split of the GREC datasets in accordance with Belz et al. (2009). After excluding empty pronouns, the MSR training, development, and test sets contain 9413, 519, 1038 instances, and the NEG training, development, and test sets contain 6681, 259, 896 instances.

**Proportion of Referring Expressions** As shown in Table 2, pronouns and proper names make up 80% and 89.5% of the referential instances in MSR and NEG, respectively. This implies that the other two referential forms, namely descriptions and empty references, account for approximately 20% of the cases in MSR and about 10% in NEG. Given this imbalance in the frequency of different forms within the two corpora, we question its potential effect on algorithm performance. Specifically, we are wondering if forms with lower frequencies are accurately predicted by the algorithms.

## 6 Evaluation

In this section, we introduce the evaluation protocol and report the performance of the models.

### 6.1 Implementation Details

For BERT and RoBERTa, we used *bert-base-cased* and *roberta-base*, both from Hugging Face. For fine-tuning, we set the batch size to 16, the learning rate to 1e-3, the dropout rate to 0.5, and the size of the output layer to 256. We ran each model for 20 epochs and used the one that achieved the highest F1 score on the development set. The implementation details of the classic ML-based models can be found in Appendix B.

### 6.2 Evaluation Protocol

The main evaluation metric in the GREC-MSR shared tasks was accuracy. In addition to accuracy,

| | MSR | | | NEG | | | WSJ | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Acc. | F1 | wF1 | Acc. | F1 | wF1 | Acc. | F1 | wF1 |
| UDel | 66.86 | 56.76 | 64.3 | **80.80** | 55.45 | 77.9 | 63.74 | 64.23 | 63.2 |
| ICSI | <u>71.19</u> | 64.73 | 70.4 | 80.36 | 64.53 | <u>78.6</u> | 64.62 | 64.15 | 63.4 |
| CNTS | 68.59 | 61.39 | 67.2 | 78.68 | 61.62 | 76.8 | 64.31 | 64.59 | 64.4 |
| OSU | 68.02 | 60.28 | 66.6 | 79.24 | 57.04 | 76.5 | 69.20 | 69.63 | 68.9 |
| IS-G | 67.05 | 58.83 | 65.3 | 77.34 | 59.52 | 75.6 | 69.15 | 69.35 | 69.2 |
| BERT | **71.68** | <u>66.70</u> | **71.4** | 77.79 | <u>72.87</u> | 77.7 | <u>80.95</u> | <u>80.93</u> | <u>80.9</u> |
| RoBERTa | 70.91 | **67.53** | <u>70.7</u> | **80.80** | **77.29** | **80.7** | **82.61** | **82.70** | **82.6** |
| Average | 69.19 | 62.32 | 67.99 | 79.29 | 64.05 | 77.69 | 70.65 | 70.80 | 70.37 |

Table 3: Overall accuracy (Acc.), macro-averaged F1 (F1), and weighted-macro F1 (wF1) scores of the algorithms depicted in Section 4. For instance, MSR-UDel refers to a C5.0 classifier trained on the MSR corpus, using the feature set mentioned in Greenbacker and McCoy (2009a).

we also report macro-F1 and weighted-macro F1. We argue that different metrics evaluate algorithms from different perspectives and provide us with different meaningful insights. For pragmatic tasks like REG, it makes sense to ask how well an algorithm performs on naturally distributed data which is often imbalanced. For these cases, reporting accuracy and weighted F1 are logical. Furthermore, analogous to other classification tasks, minority categories should not be overlooked. Take as an example the class *description* in the NEG corpus, which occurs only 4%. If a model fails to produce this class, the produced document might sound unnatural. Therefore, it is important to ensure that an algorithm is not over- or under-generating certain classes. Looking into accuracy and macro-F1 together provides insights into such cases.

## 6.3 Performance of the Models

The overall accuracy of the models, their macro F1, and their weighted-macro F1 are presented in Table 3. We also present the ranking of the models based on these scores in Appendix A.

**PLM-based Models.** The best-performing models across all corpora and metrics are PLM-based models. In six out of nine rankings, BERT and RoBERTa are ranked as the top two models. The sole exception is NEG, where BERT is the second worst model. The benefit of using PLMs is the largest on the WSJ corpus. For example, RoBERTa improves the macro F1 score from 69.63 (i.e., the performance of the best ML-based model) to 82.70.

**ML-based Models.** In contrast to the robust performance of the PLM models, the performance of the classic ML models is more corpus-dependent. In the case of MSR and NEG, ICSI is the best-performing model, while in the case of WSJ, it

is at the bottom section of the rankings. Another interesting observation is the performance of the UDel models. In terms of accuracy, UDel has the highest performance in NEG, while it has the lowest performance in both MSR and WSJ. In terms of macro-F1 rankings, the NEG UDel model dropped from first to last place, whereas BERT improved from penultimate place to second place. In general, our ML models yielded lower scores than the original models used in the GREC study (Belz et al., 2009). This could be attributed to a variety of factors, including differences in feature engineering and model parameters.

**Comparing Different Metrics.** Upon comparing average scores across the three metrics, we observe that for MSR and NEG, PLMs are clear winners only when macro-F1 is the metric in question. However, for WSJ, PLMs are winners on all three metrics. This may be because the distribution of categories in WSJ is much more balanced than in the other two corpora.

## 7 Analysis

To further compare the different models and investigate the impact of the choice of corpus, we conduct (1) a Bayes Factor (BF) analysis to determine whether the accuracy rates reported in Section 6 come from similar or different distributions, (2) a per-class evaluation of predictions to assess the success of each model in predicting individual classes, (3) a correlation analysis to quantify how the evaluation results change with respect to the choice of a corpus, and (4) a feature selection study to check how the importance of each feature changes as a function of the choice of corpus.

| | | MSR | | | NEG | | | WSJ | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Model | Category | P | R | F | P | R | F | P | R | F |
| Udel | description | 55.36 | 19.38 | 28.71 | 0.00 | 0.00 | 0.00 | 60.29 | 62.95 | 61.59 |
| | name | 72.39 | 62.21 | 66.92 | 76.65 | 80.32 | 78.44 | 60.42 | 49.44 | 54.38 |
| | pronoun | 64.53 | 88.51 | 74.64 | 84.06 | 92.14 | 87.91 | 71.00 | 83.44 | 76.72 |
| ICSI | description | 51.69 | 38.12 | 43.88 | 100.00 | 17.74 | 30.13 | 81.92 | 40.53 | 54.22 |
| | name | 80.33 | 66.82 | 72.95 | 81.85 | 73.14 | 77.25 | 55.12 | 86.40 | 67.37 |
| | pronoun | 69.41 | 87.39 | 77.37 | 79.05 | 94.76 | 86.19 | 72.17 | 69.61 | 70.86 |
| CNTS | description | 53.68 | 31.88 | 40.00 | 75.00 | 14.52 | 24.33 | 64.31 | 63.67 | 63.30 |
| | name | 76.79 | 61.75 | 68.45 | 77.84 | 72.87 | 75.27 | 60.34 | 66.75 | 63.38 |
| | pronoun | 66.16 | 88.51 | 75.72 | 79.32 | 92.14 | 85.25 | 71.90 | 62.54 | 66.89 |
| OSU | description | 53.57 | 28.12 | 36.88 | 100.00 | 4.84 | 9.23 | 72.70 | 56.91 | 63.84 |
| | name | 69.39 | 68.43 | 68.91 | 79.01 | 72.07 | 75.38 | 63.56 | 73.30 | 68.08 |
| | pronoun | 69.20 | 81.98 | 75.05 | 79.27 | 95.20 | 86.51 | 73.43 | 80.87 | 76.97 |
| ISG | description | 57.97 | 25.00 | 34.93 | 77.78 | 11.29 | 19.72 | 73.88 | 63.41 | 68.25 |
| | name | 71.46 | 65.21 | 68.19 | 71.77 | 79.79 | 75.57 | 62.19 | 76.64 | 68.66 |
| | pronoun | 65.10 | 84.01 | 73.36 | 82.30 | 84.28 | 83.28 | 75.36 | 67.36 | 71.14 |
| BERT | description | 52.86 | 46.25 | 49.33 | 62.71 | 59.68 | 61.16 | 82.63 | 79.37 | 80.97 |
| | name | 74.35 | 72.81 | 73.57 | 77.32 | 75.27 | 76.28 | 79.64 | 82.69 | 81.14 |
| | pronoun | 74.84 | 79.73 | 77.21 | 80.04 | 82.31 | 81.16 | 80.48 | 80.87 | 80.67 |
| RoBERTa | description | 56.33 | 55.62 | 55.97 | 76.47 | 62.90 | 69.02 | 86.19 | 77.40 | 81.56 |
| | name | 76.50 | 64.52 | 70.00 | 78.70 | 80.59 | 79.63 | 77.22 | 89.25 | 82.80 |
| | pronoun | 71.40 | 82.66 | 76.62 | 83.04 | 83.41 | 83.22 | 86.47 | 81.19 | 83.75 |

Table 4: Per-class precision, recall and F1 score of each label. The results report on training seven different algorithms on three corpora for predicting three labels, namely description, name, and pronoun.

## 7.1 Bayes Factor Analysis

Given that the accuracy scores are provided for all GREC systems in Belz et al. (2009), we chose to focus our analysis on the raw distributions of these scores. Our aim is to determine if there are significant differences between the accuracies of our models by comparing these distributions. We conduct a Bayes Factor analysis with a beta distribution of 0.01 (henceforth: the threshold). This analysis aims to assess, for each pair of accuracies, how strong the evidence is that they come from a common distribution, or from different ones. A difference below the threshold indicates that accuracy rates come from similar distributions; whereas, a difference above the threshold indicates that they come from different distributions, thus signalling that they differ evidentially. We interpret the strength of the evidence in favour of/against similar/different distributions according to Kass and Raftery (1995). Therefore, based on this approach, we expect that the raw accuracy distributions of the best- and worst-performing models for each corpus differ evidentially.

For MSR, the comparison between the best- and worst-performing models, namely BERT and UDel, provides no evidence that their accuracy rates are

evidentially different from each other (BF = 1.4). The same holds for NEG, where the comparison of the best (UDel and RoBERTa) and worst (IS-G) models appear to have similar probability distributions; therefore, these models are not evidentially different from each other. Conversely, in the case of WSJ, the BF analysis provides strong evidence that the accuracy distributions of the top-performing models, BERT and RoBERTa, are different from those of the classic ML models.

To summarise, we only observed significant differences in the WSJ-based models; the GREC models show more or less the same accuracy distributions. A reason might be that the aggregated calculation of accuracy loses the specificity of the classes being calculated.

## 7.2 Per-class Evaluation

As mentioned earlier, the NEG models demonstrate high accuracy (e.g. the highest average accuracy), but we observe a sharp decline in their macro-F1 values. In this analysis, we want to investigate whether the accuracy scores reported in Table 3 truly reflect the success of these algorithms or if they are merely the by-product of over-generating the dominant label or under-generating the less frequent label. Table 4 presents the *per-class* preci-

sion, recall, and F1 scores of these models.

Upon comparing the F1 scores for the class *description* across the three corpora, we observe that the WSJ models consistently achieve the highest scores, with all algorithms exceeding an F1 score of 50. In contrast, the F1 scores for both MSR and NEG are considerably lower than those of WSJ. The F1 scores for NEG are particularly low, with two notable instances, UDel and OSU, scoring 0 and below 10 respectively. The poor prediction of the class description by the classic ML NEG models is likely due to an insufficient number of instances in the training dataset, thereby hindering the proper training of the algorithms. In contrast, the two PLM models demonstrate acceptable performance in predicting the class description (BERT = 61.16 & RoBERTa = 69.02). This could indicate that pre-trained language models are advantageous where there is a class imbalance.

Another interesting observation concerns the high recall of the "pronoun" prediction in the NEG models. Four of the classic models have a recall of over 92. In the case of OSU, for example, the recall is 95, which means that of all the cases that are pronouns, 95% are labelled correctly. This is possibly an indication that pronouns have been over-generated in this system. In the PLM models, the recall is below 84.

In sum, the results of our per-class evaluation show the difficulties that the classic ML-based NEG models had in predicting the class *description*. The MSR models also had poor performance in predicting descriptions, yet they were more successful than NEG. These results tentatively suggest that feature-based classification models need to be trained on an adequate and relatively balanced number of instances to reliably predict all classes. The results of this study suggest that the PLM models are less dependent on the choice of corpus, and therefore predict classes more robustly.

### 7.3 Correlation Analysis

To quantify how the evaluation results change with respect to corpora, we compute the Spearman correlation coefficient between every pair of corpora, indicating how the rank of the models changes. Table 5 shows the computed coefficients along with the p-values of the tests. It is noteworthy that only the results evaluated by the macro-weighted F1 on MSR and NEG are significantly correlated ($p < .001$).

|  |  | acc | F1 | wF1 |
|---|---|---|---|---|
| MSR/NEG | $r_s$ | -0.1081 | 0.9643 | 0.4643 |
|  | $p$ | 0.8175 | 0.0005 | 0.2939 |
| MSR/WSJ | $r_s$ | 0.2857 | 0.5357 | 0.4643 |
|  | $p$ | 0.5345 | 0.2152 | 0.2939 |
| NEG/WSJ | $r_s$ | -0.1261 | 0.5000 | -0.0357 |
|  | $p$ | 0.7876 | 0.2532 | 0.9394 |

Table 5: Spearman correlation coefficient $r_s$ and the p-value between every pair of corpora in terms of accuracy, macro-averaged F1, and weighted F1.

The lack of correlation between the results on MSR/WSJ and those on NEG/WSJ suggests that using a corpus of a different genre could greatly influence the ranking of the models and, therefore, make the conclusions difficult to generalise. Additionally, these results are in line with the fact that MSR and NEG are from the same source, both being the introductory part of Wikipedia articles, and a higher correlation is to be expected. Also, we may conclude that macro-averaged F1 is a more reliable evaluation metric (see the discussions in Section 6, Section 7.1, and Section 7.2).

### 7.4 Feature Selection Study

We performed a feature importance analysis to check whether the contribution of linguistic factors changes depending on the choice of the corpus. We used XGBoost from the family of Gradient Boosting trees (Chen and Guestrin, 2016) and then computed the permutated variable importance for each model. Data were analysed in two ways: firstly, we used the complete dataset, as outlined in Section 5; secondly, we excluded first-mention REs to concentrate only on subsequent mentions. Considering that the choice of a referent' first mention is less context-dependent, we only report on the latter dataset below:

As expected, the ranking of feature importance varies across different corpora. However, a substantial overlap is observed when considering the most important features across the three corpora. An example is the semantic category of the REs that is used in various MSR and WSJ REG models.[4] In the case of MSR, the REs belong to five semantic categories: human, city, country, river, and mountain. In the case of WSJ, the REs are annotated for a wide

---

[4]Only human referents are annotated in NEG; therefore, this feature is not applicable.
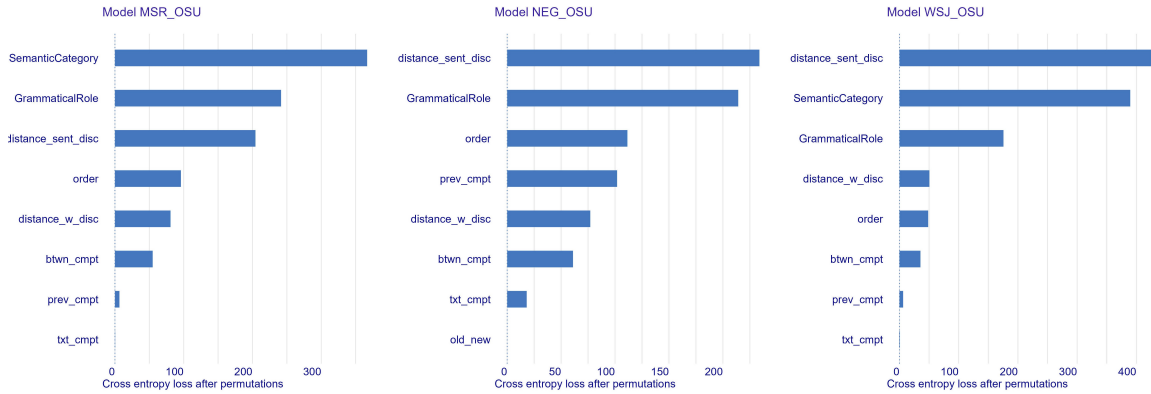
Figure 2: Different rankings of the features in MSR, NEG, and WSJ OSU models.

range of categories including human, city, country, organisation, objects, etc. Notably, in every model that employs semantic category information, this feature has either the highest or second-highest importance ranking. A plausible explanation could be that humans use different referencing strategies to refer to different categories of referents.

In addition to the semantic category, the grammatical role of the RE and the categorical sentential distance to the antecedent consistently have a high importance ranking. The grammatical role marks the distinction between subject, object, and determiner roles. The categorical distance in the number of sentences provides information on how far an RE is to its nearest coreferential antecedent. For instance, whether they are both in the same sentence or are separated by one or more sentences. Figure 2 illustrates the importance rankings of the OSU features in the three corpora. Other importance ranking graphs are available in Appendix C. For a comprehensive description of all features employed in classic ML models and the feature importance analysis, refer to Same and van Deemter (2020).

## 8 Discussion

In this paper, we have conducted a series of re-productions, evaluations, and analyses to check whether the conclusions of GREC are still true after 15 years. Below, we summarise and discuss our findings in accordance with our three research questions in Section 3. We also report our post-hoc observations on the choice of evaluation metric.

**Performance of REG Algorithms.** To answer research question $\mathcal{R}_2$, we extended the GREC by introducing a corpus of a different genre, WSJ, and two pre-trained (PLM-based) REG models. We found that, on MSR, PLM-based and ML-based

models perform similarly, as confirmed by both the BF and per-class analyses. With regards to NEG, PLM-based and ML-based models have similar accuracy scores, as confirmed by the BF analysis, but there are large differences when micro-F1 is used, as confirmed by the per-class evaluation (i.e., ML-based models have difficulty predicting descriptions). On WSJ, PLM-based models are the clear winners.

These results suggest that, in terms of explanatory power, PLM-based models have good performance and good "direct support", i.e., a good ability to generalise to different contexts (see van Deemter (2023) for further discussion). Whether they have good "indirect support" (e.g., whether their predictions are in line with linguistic theories) needs to be investigated in further probing studies.

**Impact of the Choice of Corpus.** As our evaluations and analyses demonstrate, the choice of corpus plays a crucial role in assessing REG algorithms. This role is twofold. Firstly, the choice of corpus strongly influences the evaluation results, pertaining to the research question $\mathcal{R}_1$. Secondly, in addition to the score differences discussed in Section 6, we found that: (1) the difference between PLM-based and ML-based models on WSJ is larger (and evidentially different) than on MSR and NEG models (as evidenced by the BF analysis); (2) the correlations of the evaluation results between WSJ and both MSR and NEG are not significant.

For $\mathcal{R}_3$, we conducted feature selection analyses across the three corpora, discovering that the importance of the features ranks differently for each corpus. This suggests that when investigating the "indirect support" for a model, one needs to aggregate findings from multiple corpora with different genres.

100

**The Use of Evaluation Metrics.** As we discussed in Section 6.2, different metrics evaluate different aspects of a model. This was further ascertained by the inconsistency of the BF analysis and per-class analysis. One lesson we have learned is that it is not enough to report or do analyses on a single metric. Another lesson is that the evaluation results by macro-F1 are more reliable than other metrics because (1) they are consistent across corpora with similar genres (i.e., MSR and NEG; see the Correlation analysis results); (2) the differences identified by using macro-F1 can be confirmed by the per-class evaluation.

## 9   Conclusion

We are now in a position to address the question that we raised in the Introduction: Can the conclusions from the GREC shared tasks still be trusted? By examining a wider class of corpora, models, and evaluation metrics than before, we found that the answer to this question is essentially negative since the GREC conclusions are prone to drastic change once a different corpus or a different metric is employed.

Perhaps this should come as no surprise. According to a widely accepted view of scientific progress (e.g., Jaynes (2002); applied to NLP in (van Deemter, 2023)), theories should be updated again and again in light of new data (i.e., indirect Support), and when new models are proposed, the plausibility of existing models should be compared against the plausibility of these new models (as well as pre-existing ones). New metrics deserve a place in this story as well, even though they are often overlooked. In other words, what we have seen in the present study is nothing more than science in progress – something we are bound to see more of as the enterprise called NLP-as-Science matures.

*Ethics Statement:* Regarding potential biases, in addition to the biases present in text-based datasets, biases can also be introduced by the pre-trained language models (Bender et al., 2021) used in this work. In other words, the REG algorithms we developed in this study may make different predictions with respect to different genders, for instance. In the future, we plan to investigate this phenomenon and find ways to mitigate it.

*Supplementary Materials Availability Statement:* All associated data, source code, output files, scripts, documentation, and other relevant material to this paper are publicly available and can be accessed on our GitHub repository: `https://github.com/fsame/REG_GREC-WSJ`, DOI: 10.5281/zenodo.8182689.

## References

Anja Belz, Eric Kow, Jette Viethen, and Albert Gatt. 2008. The GREC challenge 2008: Overview and evaluation results. In *Proceedings of the Fifth International Natural Language Generation Conference*, pages 183–193, Salt Fork, Ohio, USA. Association for Computational Linguistics.

Anja Belz, Eric Kow, Jette Viethen, and Albert Gatt. 2009. Generating referring expressions in context: The GREC task evaluation challenges. In *Empirical methods in natural language generation*, pages 294–327. Springer.

Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big. *Proceedings of FAccT*.

Bernd Bohnet. 2008. IS-G: The comparison of different learning techniques for the selection of the main subject references. In *Proceedings of the Fifth International Natural Language Generation Conference*, pages 192–193, Salt Fork, Ohio, USA. Association for Computational Linguistics.

Meng Cao and Jackie Chi Kit Cheung. 2019. Referring expression generation using entity profiles. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3163–3172, Hong Kong, China. Association for Computational Linguistics.

Thiago Castro Ferreira, Diego Moussallem, Ákos Kádár, Sander Wubben, and Emiel Krahmer. 2018. NeuralREG: An end-to-end approach to referring expression generation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1959–1969, Melbourne, Australia. Association for Computational Linguistics.

Thiago Castro Ferreira, Chris van der Lee, Emiel van Miltenburg, and Emiel Krahmer. 2019. Neural data-to-text generation: A comparison between pipeline and end-to-end architectures. In *Proceedings of the 2019 Conference on Empirical Methods*

*in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 552–562, Hong Kong, China. Association for Computational Linguistics.

Guanyi Chen, Fahime Same, and Kees van Deemter. 2021. What can neural referential form selectors learn? In *Proceedings of the 14th International Conference on Natural Language Generation*, pages 154–166, Aberdeen, Scotland, UK. Association for Computational Linguistics.

Guanyi Chen, Fahime Same, and Kees van Deemter. 2023. Neural referential form selection: Generalisability and interpretability. *Computer Speech & Language*, 79:101466.

Tianqi Chen and Carlos Guestrin. 2016. XGBoost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '16, page 785–794, New York, NY, USA. Association for Computing Machinery.

Rossana Cunha, Thiago Castro Ferreira, Adriana Pagano, and Fabio Alves. 2020. Referring to what you know and do not know: Making referring expression generation models generalize to unseen entities. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 2261–2272, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Walter Daelemans, Jakub Zavrel, Ko van der Sloot, and Antal van den Bosch. 2007. Timbl: Tilburg memory-based learner.

Benoit Favre and Bernd Bohnet. 2009. ICSI-CRF: The generation of references to the main subject and named entities using conditional random fields. In *Proceedings of the 2009 Workshop on Language Generation and Summarisation (UCNLG+Sum 2009)*, pages 99–100, Suntec, Singapore. Association for Computational Linguistics.

Charles Greenbacker and Kathleen McCoy. 2009a. UDel: Generating referring expressions guided by psycholinguistc findings. In *Proceedings of the 2009 Workshop on Language Generation and Summarisation (UCNLG+Sum 2009)*, pages 101–102, Suntec, Singapore. Association for Computational Linguistics.

Charles F Greenbacker and Kathleen F McCoy. 2009b. Feature selection for reference generation as informed by psycholinguistic research. In *Proceedings of the CogSci 2009 Workshop on Production of Referring Expressions (PRE-Cogsci 2009)*.

Samir Gupta and Sivaji Bandopadhyay. 2009. Junlg-msr: A machine learning approach of main subject reference selection with rule based improvement. In *Proceedings of the 2009 Workshop on Language Generation and Summarisation*, pages 103–104. Association for Computational Linguistics.

Iris Hendrickx, Walter Daelemans, Kim Luyckx, Roser Morante, and Vincent Van Asch. 2008. CNTS: Memory-based learning of generating repeated references. In *Proceedings of the Fifth International Natural Language Generation Conference*, pages 194–95, Salt Fork, Ohio, USA. Association for Computational Linguistics.

Renate Henschel, Hua Cheng, and Massimo Poesio. 2000. Pronominalization revisited. In *Proceedings of the 18th conference on Computational linguistics-Volume 1*, pages 306–312. Association for Computational Linguistics.

Eduard Hovy, Mitchell Marcus, Martha Palmer, Lance Ramshaw, and Ralph Weischedel. 2006. Ontonotes: the 90% solution. In *Proceedings of the human language technology conference of the NAACL, Companion Volume: Short Papers*, pages 57–60.

Emily Jamison and Dennis Mehay. 2008. OSU-2: Generating referring expressions with a maximum entropy classifier. In *Proceedings of the Fifth International Natural Language Generation Conference*, pages 196–197, Salt Fork, Ohio, USA. Association for Computational Linguistics.

E.T. Jaynes. 2002. *Probability Theory: The Logic of Science*. Cambridge University Press, Cambridge, UK.

Robert E Kass and Adrian E Raftery. 1995. Bayes factors. *Journal of the american statistical association*, 90(430):773–795.

Andrej A Kibrik, Mariya V Khudyakova, Grigory B Dobrov, Anastasia Linnik, and Dmitrij A Zalmanov. 2016. Referential choice: Predictability and its limits. *Frontiers in psychology*, 7(1429).

Emiel Krahmer and Kees van Deemter. 2012. Computational generation of referring expressions: A survey. *Computational Linguistics*, 38(1):173–218.

Max Kuhn, Steve Weston, Mark Culp, Nathan Coulter, and Ross Quinlan. 2018. Package 'c50'.

Kathleen E. McCoy and Michael Strube. 1999. Generating anaphoric expressions: Pronoun or definite description? In *The Relation of Discourse/Dialogue Structure and Reference*.

Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Olga Uryupina, and Yuchen Zhang. 2012. CoNLL-2012 shared task: Modeling multilingual unrestricted coreference in OntoNotes. In *Joint Conference on EMNLP and CoNLL - Shared Task*, pages 1–40, Jeju Island, Korea. Association for Computational Linguistics.

Ehud Reiter. 2017. A commercial perspective on reference. In *Proceedings of the 10th International Conference on Natural Language Generation*, pages 134–138, Santiago de Compostela, Spain. Association for Computational Linguistics.

Fahime Same, Guanyi Chen, and Kees Van Deemter. 2022. Non-neural models matter: a re-evaluation of neural referring expression generation systems. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5554–5567, Dublin, Ireland. Association for Computational Linguistics.

Fahime Same and Kees van Deemter. 2020. A linguistic perspective on reference: Choosing a feature set for generating referring expressions in context. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4575–4586, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Kees van Deemter. 2016. *Computational models of referring: a study in cognitive science*. MIT Press.

Kees van Deemter. 2023. Dimensions of explanatory value in nlp models. *Computational Linguistics*.

Ralph Weischedel, Martha Palmer, Mitchell Marcus, Eduard Hovy, Sameer Pradhan, Lance Ramshaw, Nianwen Xue, Ann Taylor, Jeff Kaufman, Michelle Franchini, et al. 2013. Ontonotes release 5.0 ldc2013t19. *Linguistic Data Consortium, Philadelphia, PA*, 23.

## A    Ranking of the Models

### Accuracy-based Ranking

MSR: BERT > ICSI > RoBERTa > CNTS > OSU > IS-G > UDel

NEG: UDel = RoBERTa > ICSI > OSU > CNTS > BERT > IS-G

WSJ: RoBERTa > BERT > OSU > IS-G > ICSI > CNTS > UDel

### Macro-F1 Ranking

MSR: RoBERTa > BERT > ICSI > CNTS > OSU > IS-G > UDel

NEG: RoBERTa > BERT > ICSI > CNTS > IS-G > OSU > UDel

WSJ: RoBERTa > BERT > OSU > IS-G > CNTS > UDel > ICSI

### Macro-weighted F1 Ranking

MSR: BERT > RoBERTa > ICSI > CNTS > OSU > IS-G > UDel

NEG: RoBERTa > ICSI > UDel > BERT > CNTS > OSU > IS-G

WSJ: RoBERTa > BERT > IS-G > OSU > CNTS > ICSI > UDel

## B    Implementation Details for ML-based Models

The R programming language was used mostly for running the classic ML models. The specification of the models can be found below:

**Conditional Random Field [CRF].**   The R Package CRF (https://cran.r-project.org/web/packages/crfsuite/) was used to train these models. The iterations are set to 3000, and the learning method is Stochastic Gradient Descent with L2 regularization term (l2sgd).

**Decision Tree [C5.0].**   The R Package C5.0 (Kuhn et al., 2018) was used to build the decision trees. The number of boosting iterations (trials) is set to 3, and the splitting criterion is information gain (entropy).

**Memory-Based Learning [MBL].**   As mentioned before, we implemented the k-Nearest Neighbors [KNN] algorithm instead of MBL. The R package caret with the method KNN was used to implement this model.

**Maximum Entropy [MaxEnt].** The multinom algorithm from the nnet R package was used to implement this model.

**Multi-Layer Perceptron [MLP].** The Keras package was used to implement MLP. The model consists of two hidden layers with 16 and 8 units, respectively. The hidden layers use the rectified linear activation function (ReLU), and the output layer uses the Sigmoid activation function. The model is fitted for 50 training epochs. In addition, 50 samples (batch size) are propagated through the network.

**eXtreme Gradient Boosting [XGBoost].** XG-Boost was used for the feature selection experiments. We used the R packages xgboost and DALEXtra for the analysis. We set the learning rate to 0.05, the minimum split loss to 0.01, the maximum depth of a tree to 5, and the sub-sample ratio of the training instances to 0.5.

## C Feature Importance Rankings

The graphs in Figure 3 show the rankings across MSR, WSJ, and WSJ. A maximum number of eight features is depicted in the graphs.
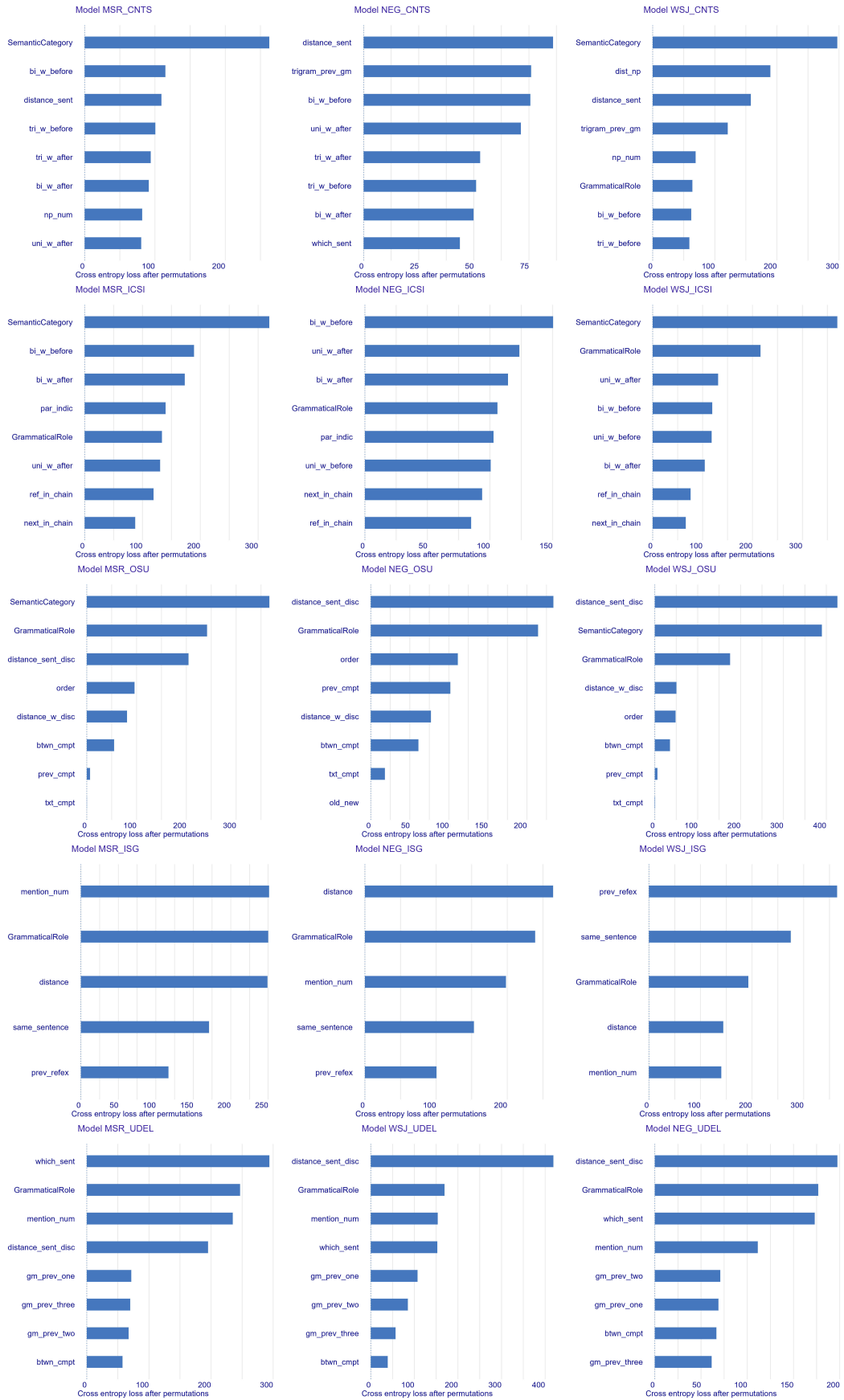
Figure 3: Importance ranking of the features in MSR, NEG, and WSJ models.

# Generating Faithful Text From a Knowledge Graph
# with Noisy Reference Text

**Tahsina Hashem[1], Weiqing Wang[1], Derry Tanti Wijaya[2],**
**Mohammed Eunus Ali[3], Yuan-Fang Li[1]**
[1]Department of Data Science & AI, Monash University, Australia
[2]Department of Data Science, Monash University, Indonesia
[3]Department of CSE, Bangladesh University of Engineering and Technology, Bangladesh
{tahsina.hashem, Teresa.Wang, derry.wijaya, yuanfang.li}@monash.edu;
eunus@cse.buet.ac.bd

## Abstract

Knowledge Graph (KG)-to-Text generation aims at generating fluent natural-language text that accurately represents the information of a given knowledge graph. While significant progress has been made in this task by exploiting the power of pre-trained language models (PLMs) with appropriate graph structure-aware modules, existing models still fall short of generating faithful text, especially when the ground-truth natural-language text contains additional information that is not present in the graph. In this paper, we develop a KG-to-text generation model that can generate faithful natural-language text from a given graph, in the presence of noisy reference text. Our framework incorporates two core ideas: Firstly, we utilize contrastive learning to enhance the model's ability to differentiate between faithful and hallucinated information in the text, thereby encouraging the decoder to generate text that aligns with the input graph. Secondly, we empower the decoder to control the level of hallucination in the generated text by employing a controllable text generation technique. We evaluate our model's performance through the standard quantitative metrics as well as a ChatGPT-based quantitative and qualitative analysis. Our evaluation demonstrates the superior performance of our model over state-of-the-art KG-to-text models on faithfulness.

## 1 Introduction

A knowledge graph (KG) is a structured representation of information as a network of interconnected real-world entities, and relationships. The task of KG-to-text generation has been proposed (Ribeiro et al., 2020a; Koncel-Kedziorski et al., 2019) to make this structured information more accessible to humans, aiming to generate fluent, informative, and faithful natural-language sentences that should describe the contents of an input KG. Recently,

this task plays a significant role in a variety of applications such as knowledge-grounded dialogue generation (Zhou et al., 2018; Zhao et al., 2020), story generation (Guan et al., 2019; Ji et al., 2020), event narration (Colas et al., 2021a), and question-answering (Agarwal et al., 2021; Chen et al., 2023; Saxena et al., 2020).

Significant progress has been made in the KG-to-text generation task by utilizing a set of Transformer-based (Vaswani et al., 2017) pre-trained language models (PLMs) such as BART (Lewis et al., 2019), T5 (Raffel et al., 2020) or GPT (Radford et al., 2019) with appropriate graph structure-aware modules (Ke et al., 2021; Colas et al., 2022; Han and Shareghi, 2022). However, ensuring the faithfulness of KG-to-text generation, i.e. reducing hallucinations (Ji et al., 2022; Wang et al., 2022; Raunak et al., 2021; Rebuffel et al., 2022), is an under-explored problem, and existing KG-to-text models fall short of generating faithful text when the ground-truth text of the training dataset contains wrong or extra information that is not consistent with the input.

Figure 1 shows an example of a small KG about a house, which contains information on its internal features and neighborhood, and the corresponding ground-truth reference text, from a real-world real-estate KG (Das et al., 2021). The ground-truth text, while summarizing the features of the house accurately, also mentions some information that is not available in the input KG (i.e. extrinsic hallucination, highlighted in red).

When a KG-to-text model is trained with such hallucinated reference text, it is likely to produce text that is also hallucinated. This hallucination problem significantly reduces the faithfulness and thus trustworthiness of the generated text. Thus, the ability to reduce hallucination in the presence of noisy reference text is important for the practi-

cal application of KG-to-text and other NLG techniques, especially in mission- and safety-critical domains such as medical diagnostics and scientific research.

A number of techniques have been proposed (Ji et al., 2022) to control this hallucination problem in abstractive summarization, table-to-text generation, generative question-answering, neural machine translation, and knowledge-grounded dialogue generation (Wang et al., 2022; Tang et al., 2022; Rebuffel et al., 2022; Krishna et al., 2021; Zhou et al., 2021; Zhang et al., 2022). However, to the best of our knowledge, controlling hallucination in graph-to-text generation with noisy reference text has not been investigated.

In this paper, we propose a novel framework to address this important and practical problem. Our framework combines contrastive learning technique and controllable text generation. Contrastive learning enables the model to distinguish between faithful and hallucinated text and guides the decoder to generate faithful text instead of hallucinated text. The controllable text generation technique learns the level of hallucination from noisy training text and controls (i.e. minimizes) the level of hallucinated information in the generated text. Our framework can be employed in any KG-to-Text encoder-decoder model to generate faithful natural language text from a given KG, in the presence of noisy reference text.

Our contributions are as follows:

- We propose a framework to deal with the hallucination problem in KG-to-text generation task. Our framework comprises two core ideas: (i) Employing contrastive learning to enable the KG-to-text generation model to better differentiate between faithful and hallucinated information in the reference text and guide the decoder to generate text that is faithful to KG. (ii) Controlling the level of hallucination while generating text from KG using a controllable text generation technique.
- We conduct experiments and evaluate performance using a standard quantitative analysis with automatic metrics. Our comprehensive evaluation on two noisy datasets demonstrates the superior performance of our proposed model over the state-of-art KG-to-text generation models on faithfulness metrics.
- We further propose and perform novel ChatGPT-based quantitative and qualitative

evaluations to assess the performance of our model more comprehensively. The evaluation also shows our model's effectiveness in generating faithful text over existing KG-to-text generation models.

## 2 Related Work

### 2.1 Knowledge Graph-to-Text Generation

KG-to-text generation techniques (Koncel-Kedziorski et al., 2019; Guo et al., 2020; Ribeiro et al., 2020b; Chen et al., 2020) utilize graph neural networks (Veličković et al.) and graph Transformers (Vaswani et al., 2017) to effectively encode a graph's structural information. With the rapid advancement of pre-trained language models (PLMs) (Lewis et al., 2019; Raffel et al., 2020; Radford et al., 2019), researchers have started adapting and fine-tuning these models to KG-to-text generation tasks and obtained better results compared to previous models (Ribeiro et al., 2021; Chen et al., 2020; Kale and Rastogi, 2020). Recently, researchers further improved the KG-to-text models' performance by integrating pre-trained language models with appropriate graph-structure-aware modules (Ke et al., 2021; Colas et al., 2022) and employing some graph masking pre-training tasks (Ke et al., 2021; Han and Shareghi, 2022).

However, we have empirically observed that although these state-of-art KG-to-text generation models (Ke et al., 2021; Colas et al., 2022; Han and Shareghi, 2022) introduce graph aware encoders and/or apply graph masking pre-training strategies to enhance graph-text alignments, still these models are struggling with hallucination problems when trained with noisy input ground-truth text.

### 2.2 Controlling Hallucinations in Text Generation

This hallucination problem is well explored in other natural language generation tasks such as in table-to-text generation, summarization, dialogue generation, question-answering, and neural machine translation. Planning (Su et al., 2021) or skeleton-based method (Wang et al., 2021), joint learning strategy (Xu et al., 2021), Bayes training framework (Tian et al., 2019), table-text optimal-transport matching strategy (Wang et al., 2020), control token approach (Filippova, 2020) are widely used in controlling hallucinations in table-to-text generation tasks. Most recently, Rebuffel et

(a) House knowledge graph

(b) Ground-Truth Text

Figure 1: A sample knowledge graph for the House dataset with its ground-truth text. The red colored text in the ground-truth text represents extrinsic hallucination information.

al. (Rebuffel et al., 2022) proposed a multi-branch decoder approach to control hallucination at decoding time in this area.

Prior works have also focused on minimizing hallucinations in summarization, dialogue generation, question-answering and neural machine translation areas. Some of the recent hallucination mitigation techniques are based on control token approach (Filippova, 2020; Rashkin et al., 2021; Wang et al., 2022), contrastive learning approach (Cao and Wang, 2021; Tang et al., 2022), generate then-refine strategy (Dziri et al., 2021), a routing transformer based approach (Krishna et al., 2021) and self-training of neural machine translation based approach (Zhou et al., 2021). To the best of our knowledge, no work has been done in graph-to-text generation tasks with hallucinated ground-truth text.

## 2.3 Evaluation using ChatGPT

Large language models such as ChatGPT have recently been employed for evaluating the quality and factual consistency of the generated text in NLP tasks with respect to the source input through ranking, rating, and entailment inference (Kocmi and Federmann, 2023; Wang et al., 2023; Luo et al., 2023). Luo et al. (2023) closely investigated Chat-GPT's ability under a zero-shot setting with three factual consistency evaluation tasks: binary entailment inference, summary ranking, and consistency rating. Experimental findings show that ChatGPT generally performs better than previous evaluation metrics across the three tasks, demonstrating its significant potential for factual consistency evaluation. However, they also point out some limitations of ChatGPT such as its preference on lexical similarity instead of semantic entailment, false

reasoning, and poor understanding of instructions. Moreover, while these approaches can compute an overall faithfulness score of the output text, they fall short in terms of explaining the score e.g., by quantifying the amount of hallucination (out of all the output facts, how many are hallucinated?), precision (out of all the output facts, how many are input facts?) and recall (out of all the input facts, how many appear in the output?). In this work, we use ChatGPT to quantify each of these values and obtain a finer-grained explanation of what a faithfulness score entails.

## 3 Proposed Model

### 3.1 Problem Formulation

Let $G = (V, E)$ represent a knowledge graph, where $V = \{e_1, e_2, \ldots, e_{|V|}\}$ represents the entity set and $E = \{r_{ij}\} \subseteq V \times V$ represents the relations connecting the entities, the task of KG-to-text aims to generate a passage of text $\hat{Y} = (y_1, y_2, \ldots, y_n)$, that faithfully represents the information contained in a graph $G$. The model is given a training set $\mathcal{D} = \{(G_i, Y_i)\}$, in which the reference text $Y_i$ may contain *hallucinated* information.

### 3.2 Our Framework

Standard fine-tuning approaches use a cross-entropy loss to maximize the similarity between the ground-truth text and the output text. Thus, if the ground-truth text contains hallucination, the model trained through fine-tuning also learns to generate hallucinated text. To overcome this hallucination problem, we introduce an effective fine-tuning approach that combines a contrastive loss function and a controllable text generation technique with the cross-entropy loss function. As a result, our

method can train a KG-to-text generation model to generate faithful text from a KG.

Figure 2 depicts the overall architecture of our proposed model. The following two subsections illustrate our two proposed techniques in detail.
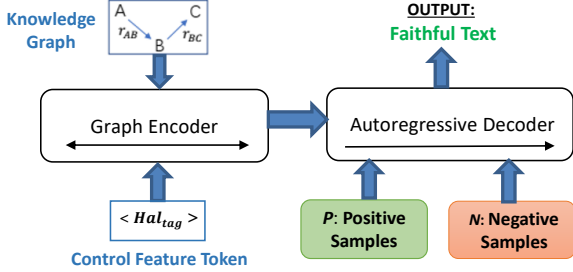


Figure 2: The overall framework of our KG-to-text model.

## 3.3 Minimizing Hallucinations with Contrastive Learning

Contrastive learning is a popular and effective representation learning method. Originally proposed for computer vision tasks (Khosla et al., 2020; Yang et al., 2022), contrastive learning has been successfully applied to learn representations for sentences/documents (Gao et al., 2021; Zhang et al., 2021), abstractive summarization (Liu and Liu, 2021; Cao and Wang, 2021; Wan and Bansal, 2022) and dialogue generation (Tang et al., 2022; Dziri et al., 2022; Geng et al., 2022). Inspired by them, we have utilized this learning framework to reduce hallucinations while generating text from knowledge graphs. It enables the model to differentiate between faithful information and hallucinated information in the text, which then assists the decoder in generating text that should be free of hallucinations.

For an input pair of a graph and an anchor reference text $(G_i, Y_i)$ from the training data $\mathcal{D}$, $P_i$ represents a set of positive samples and $N_i$ represents a set of hallucinated summaries (i.e. negative samples). The contrastive learning objective function is formulated as follows in Equation 1:

$$L_{CL} = - \sum_{(G_i, Y_i) \in \mathcal{D}} \sum_{Y_j \in P_i} \log \frac{\exp(\cos(h_i, h_j))}{\sum_{Y_k \in N_i} \exp(\cos(h_i, h_k))}$$

(1)

Here, $Y_j$ is a positive sample from the set $P_i$, $Y_k$ is a negative sample from the set $N_i$, and $h_i$, $h_j$, $h_k$ are the BART decoder representations of $Y_i$, $Y_j$, and $Y_k$ respectively.

This contrastive objective function encourages the model to learn a preference for positive (faithful) summaries over negative (hallucinated) ones. While the ground-truth text in the training data $\mathcal{D}$ is noisy, it is reasonable to assume that each reference text is more faithful to the paired graph than a randomly sampled text from $\mathcal{D}$. Based on this observation, we carefully select the positive and negative samples to ensure the effectiveness of our contrastive learning technique.

**Positive sample construction.** Back-translation (Mallinson et al., 2017) is an effective approach for preserving meanings and providing linguistic diversity. Hence, we use NLPAug (Ma, 2019) to translate each anchor text to German and back to English and take the translated text as a positive sample for the anchor passage.

**Negative sample construction.** For the anchor text of a given graph, we treat the text of any other graph in $\mathcal{D}$ as a potential negative sample. We randomly select four such text to construct $N$ for each anchor text. Dataset-specific knowledge can be easily incorporated in this approach to improve the quality of contrastive learning. For the House dataset, we adopt a simple heuristic for constructing the negative sample set. Here, we give more importance to the six major features of a house graph: (1) house location (2) house address (3) number of bedrooms (4) number of bathrooms (5) number of parking spaces, and (6) house property type. If all of these major features of a house differ from the anchor house, then the house's paired text is selected as the negative sample for the anchor house. We choose these six features as major features because information of these features is available in almost every house (91%) in the training set.

## 3.4 Controlling Hallucinations with Control Feature Token

In contrastive learning, we use the ground-truth reference text as a positive sample. As the ground-truth text contains hallucinations, when training with contrastive learning for generating text, the output text still contains some hallucinations. Thus, we employ a controllable text generation approach to further enhance the faithfulness of our model. Specifically, we append controllable features to the input graph in training in order to control the level of hallucination in the generated text.

**Control feature token.** Control feature token is a hallucination measure that quantifies how much

the given ground-truth text is faithful to the source graph. We linearized the knowledge graph (Chen et al., 2020) into a list of verbalized triples and employ BARTScore (Yuan et al., 2021) as the measure of faithfulness between the linearized graph and the corresponding ground-truth text, as it has been shown that it is closely associated with human evaluations of faithfulness (Yuan et al., 2021).



Figure 3: Controllable Text generation with Control Feature Token

**Controllable generation.** According to the BARTScore of the training samples, we split the samples into three buckets, where each bucket contains a list of training samples at a specific range of BARTScore. This range is chosen in a manner that ensures each bucket contains approximately an equal number of samples. These three buckets are represented using the following hallucination tags, $Hal_{tag}$={$Hal_{low}$, $Hal_{medium}$ and $Hal_{high}$} following existing work (Filippova, 2020; Zhang et al., 2022). At training time, we append the corresponding hallucination tag to the input sample according to its BARTScore. These three hallucination tags represent the three control feature tokens that act as a special input to control the level of hallucination during text generation.

Figure 3 illustrates the fine-tuning process with the control tokens. Let $G$ and $Y = (y_1, y_2, \ldots, y_n)$ be the input sample graph and its corresponding reference text, and $H$ be the hallucination tag (i.e. control feature token) for this input sample. Formally, we define the objective function of our fine-tuning strategy with the control token as follows:

$$L_{CE\_CtrlTok} = -\sum_{i=1}^{n} log P(y_i|y_{<i}, G, H) \quad (2)$$

Thus, during training, the model learns the mapping between the graph-text pair ($G$, $Y$) and its corresponding control token $H$. The model then becomes an expert at evaluating samples according to the control token. At inference time, the control

token is set to the desired hallucinated value i.e., low ($Hal_{low}$) to generate faithful text from the KG.

**The overall training objective** of our proposed model is the sum of the contrastive loss and the cross-entropy loss with the control token:

$$L = L_{CL} + L_{CE\_CtrlTok} \quad (3)$$

Thus, during training, instead of blindly following the ground-truth text, the model gives more focus on the faithful parts of the text instead of the hallucinated ones. Moreover, the decoder is encouraged to generate text by minimizing hallucinations through controlled measures.

## 4 Experiments

### 4.1 Dataset

We conduct experiments and evaluation on two KG-to-text generation datasets: the House dataset (Das et al., 2021) about real-estate house listing and the GenWiki dataset (Jin et al., 2020). In both datasets, the ground-truth text contains a significant amount of hallucinated information, making the task of generating faithful text especially challenging. Thus, these datasets are the most appropriate to evaluate the performance of our proposed model. Table 1 shows the statistics of these two datasets in detail. Note that we use the "FINE" version (Jin et al., 2020) of GenWiki.

| Dataset | #Relations | #KG-Text Pairs (Train / Valid / Test) |
|---|---|---|
| House | 68 | 33K / 10K / 10, 219 |
| GenWiki$_{FINE}$ | 287 | 750K / 7, 152 / 1, 000 |

Table 1: Statistics of the datasets, including the total number relations and the data split

**House.** The dataset is prepared from the large real-estate and POI datasets of Melbourne, Australia (Das et al., 2021). It includes 53, 220 records of house sales transactions from 2013 to 2015. It consists of three types of point-of-interests (POIs), namely regions, schools, and train stations, along with their corresponding features. Every sample in the dataset includes a ground-truth advertisement text describing the features of the house. However, the given ground-truth text contains a significant level of hallucinated information.

**GenWiki.** It is a large-scale non-parallel (Colas et al., 2021b) dataset prepared by matching Wikipedia articles with DBpedia entities (Jin et al., 2020).

| House Dataset | | | | | |
| --- | --- | --- | --- | --- | --- |
| **Model** | Comparison with ground-truth text | | | Comparison with linearized graph | |
| | BLEU ↑ | METEOR ↑ | ROUGE-L ↑ | BARTScore ↑ | FactCC ↑ |
| Ground-truth text (5K samples) | - | - | - | -4.564 | 48.48 |
| JointGT (Ke et al., 2021) | **3.61** | 11.96 | **18.62** | -3.685 | 49.53 |
| GAP (Colas et al., 2022) | 3.47 | **12.05** | 18.16 | -3.666 | 52.71 |
| GMP (Han and Shareghi, 2022) | 3.09 | 10.73 | 16.23 | -3.941 | 48.47 |
| **Our Full Model** | 2.54 | 11.06 | 16.86 | **-3.245** | **63.61** |
| Control token only | 2.88 | 11.2 | 17.35 | -3.567 | 52.97 |
| Contrastive learning only | 2.56 | 11.04 | 16.89 | -3.247 | 63.04 |
| GenWiki Dataset | | | | | |
| **Model** | Comparison with ground-truth text | | | Comparison with linearized graph | |
| | BLEU ↑ | METEOR ↑ | ROUGE-L ↑ | BARTScore ↑ | FactCC ↑ |
| Ground-truth text (5K samples) | - | - | - | -3.464 | 53.80 |
| CycleGT (Guo et al., 2020) | **41.59** | **35.72** | **63.31** | -3.276 | 76.86 |
| JointGT (Ke et al., 2021) | 37.93 | 32.60 | 59.06 | -2.299 | 79.94 |
| GMP (Han and Shareghi, 2022) | 35.43 | 32.68 | 57.63 | **-1.601** | 76.62 |
| **Our Full Model** | 37.48 | 32.70 | 60.40 | -2.182 | **82.85** |
| Control token only | 37.01 | 32.38 | 59.57 | -2.268 | 81.98 |
| Contrastive learning only | 35.19 | 31.33 | 57.89 | -2.309 | 81.48 |

Table 2: Results on the **House** and **GenWiki** datasets. We have used BART-base and T5-base for House dataset and Genwiki dataset respectively. **Bold** fonts denote the best results.

## 4.2 Baseline Models

We evaluate the performance of our proposed model against graph-to-text generation models that are based on an encoder-decoder architecture. On the House dataset, we choose three state-of-the-art models: JointGT model (Ke et al., 2021) that jointly learns the graph structure and text; GAP (Colas et al., 2022) that is aware of the graph structure; and GMP (Han and Shareghi, 2022), a self-supervised graph masking pre-training model. On the Gen-Wiki dataset, we compare the results of the following models: the state-of-the-art unsupervised model CycleGT (Guo et al., 2020) for Genwiki dataset, JointGT (T5) model (Ke et al., 2021) and GMP (Han and Shareghi, 2022). Note that in addition to the existing state-of-the art model, GMP, we also include CycleGT as it has the best reported performance on GenWiki dataset.

## 4.3 Experimental Settings

We adopt JointGT (Ke et al., 2021) as our base model for fine-tuning. JointGT is initialized with the Hugging Face's pre-trained BART-base checkpoint[1] for House Dataset. For GenWiki dataset the model is initialized with the Hugging Face's pre-trained T5-base checkpoint[2]. We select the pre-

[1] https://huggingface.co/facebook/bart-base
[2] https://huggingface.co/t5-base

trained LM BART-base or T5-base in order to do a fair comparison with the baseline models.

JointGT is pre-trained with a KGTEXT dataset (Chen et al., 2020). For contrastive learning, we use two positive samples and four negative samples for each training sample. For the House dataset, we fine-tune our model for 5 epochs; for the Gen-Wiki dataset, we fine-tune our model for 4000 steps. The batch size is set to 32. The maximum length of linearized input graphs is 600 and the maximum length of text sequences is set to 128 tokens. We adopt Adam (Kingma and Ba, 2015) as the optimizer and set the learning rate to be 3e-5. We used one A40 48GB GPU and one A10 24GB GPU for the experiments

## 4.4 Main Results

We use automatic metrics to measure both fluency and faithfulness of generated text. Following existing KG-to-text work, we employ standard metrics BLEU (Papineni et al., 2002), ME-TEOR (Banerjee and Lavie, 2005), and ROUGE-L (Lin, 2004). These metrics are usually used to measure accuracy and fluency of the generated text with respect to the ground-truth text. However, as the ground-truth text contains hallucinations, we cannot verify the faithfulness of the generated text by comparing with these metrics. Thus, we use BARTScore (Yuan et al., 2021) and

FactCC (Kryściński et al., 2020) for comparing the generated text with the linearized input graph for measuring faithfulness. These two metrics have been widely used for measuring faithfulness in other NLP tasks (Tang et al., 2022; Gao and Wan, 2022; Cao and Wang, 2021; van der Poel et al., 2022).

The faithfulness of the reference text of the House dataset and the GenWiki dataset is also reported in Table 2, as measured by BARTScore and FactCC score. As can be seen, the reference text of both datasets contains significant amounts of hallucination (low BARTScore and FactCC scores).

Table 2 presents the results on the House and GenWiki datasets. From the results on the House dataset, we can observe that our full model achieves best results on faithfulness measures (i.e. when compared with the linearized graph), outperforming the best baseline models on BARTScore and FactCC score by 0.421 and 10.9 absolute points respectively. The performance delta on the GenWiki dataset is smaller, where our model achieves the best performance on FactCC of 1.55 points and second best performance on BARTScore. We posit the larger performance delta on the House dataset is due to it being significantly more noisy evidenced by lower BARTScore and FactCC scores.

For BLEU, METEOR and ROUGE-L, the baseline models perform modestly better than our model when comparing with the ground-truth text. This result is expected and reasonable as compared with our model, the other models tend to generate text with higher similarity with the ground-truth text, resulting in higher values as measured by these metrics. At the same time, due to the noisy nature of the reference text, a high similarity also indicates high hallucination, as discussed above.

In Section 4.5 below, we further measure the faithfulness and fluency of generated text with ChatGPT as the oracle, where we demonstrate that our model achieves superior faithfulness while maintaining fluency.

Table 3 shows a sample ground-truth text and the text generated by different models, where correct facts are highlighted in blue and hallucinated text is highlighted in red. More examples can be found in Appendix C.

## 4.5 ChatGPT-based Evaluation

We propose to utilize ChatGPT to further measure the factual consistency and fluency of the generated



Figure 4: ChatGPT-based evaluation on 50 samples from the House test set.

text with respect to the input graph. We randomly sample 50 houses from the House test set, and perform evaluation on the text generated by different models.

To measure **fluency**, similar to (Wang et al., 2023), we prompt ChatGPT to score the fluency of the generated text. To measure **factual consistency**, we carefully design prompts to instruct ChatGPT to enumerate facts in the (linearized) graph (*# input facts*), the common facts between the graph and generated text (*# common facts*), and the hallucinated facts in the generated text (*# hallucinated facts*), respectively. By enumerating facts that are correctly generated, missing, or hallucinated, our ChatGPT-based evaluation provides better explainability of models' faithfulness. Details and examples of our prompts and ChatGPT's responses can be found in Appendix A.

In addition to enumerating the facts, ChatGPT-based evaluation provides a way to measure quantitative metrics such as precision, recall, and hallucination rates. We randomly sample 50 graph-text pairs from the test House dataset, and measure the precision (P), recall (R) and amount of hallucination (H) in generated text of these samples, which are formulated as follows: $P = \frac{\text{\# common facts}}{\text{\# output facts}}$, $R = \frac{\text{\# common facts}}{\text{\# input facts}}$, and $H = \frac{\text{\# hallucinated facts}}{\text{\# output facts}}$.

The number of output facts (*# output facts*) is computed by summing up the number of hallucinated facts (*# hallucinated facts*) and the number of common facts (*# common facts*).

Figure 4 shows the results of this analysis. It can be seen that our model outperforms all the baseline KG-to-text generation models on precision, recall and faithfulness (i.e. low hallucination) and achieves competitive scores in terms of fluency.

To determine the gap between our model and the

| |
|---|
| **Ground-truth text:** |
| **Exceptional potential with elevated outlook:**Set on the high side in a wonderfully quiet and convenient cul-de-sac location, this clinker-brick home provides an ideal opportunity to renovate and extend, rebuild or redevelop (STCA) and capture the city skyline views. Immediately comfortable, the home currently offers five bedrooms (two with built-in robes) and central bathroom, two living areas and well-appointed kitchen. Other features include heating and cooling, expansive backyard, plus garage and carport via full-length side driveway. Close to Rosehill Road/Hoffmans Road buses, shops and cafes, walking distance to St Bernard???s College and Rosehill Secondary College. Currently tenanted at $1653.00pm. Land size 15.8m x 42.1m ??? 665sqm (approx). This is the one! |
| **JointGT model:** |
| AUCTION THIS SATURDAY AT 11:30AM!:In one of Niddrie's most sought-after pockets, this solid-brick home presents an exciting opportunity to renovate and extend (STCA) on a deep block of land (approx 690sqm). Current floorplan includes five bedrooms (all with built-in robes) and central bathroom, formal lounge through to the kitchen/meals area, as well as a separate toilet. Other features of the home include gas heating and split-system air conditioning, plus a north-facing backyard complete with double garage (remote-controlled |
| **GAP model:** |
| **Spacious family home in a sought-after Niddrie location!**:Set on a large block of land (approx 697sqm), this solid-brick home presents an exciting opportunity to renovate, extend or redevelop (STCA). Current floorplan includes five bedrooms (two with BIRs) and central bathroom, open-plan living/meals area and updated kitchen with stainless steel appliances. Other features include split-system heating/cooling, polished floorboards, polished hardwood floors and a double garage with internal access. Close to schools, shops and public transport, it's a must-see |
| **GMP model:** |
| **Exceptional family living in a prime Niddrie location!**:In a quiet cul-de-sac close to schools, shops and transport, this impressive home offers an immediately comfortable family lifestyle with plenty of scope to further improve and extend (STCA). Beyond its attractive facade, the home's inviting interior comprises five bedrooms (all with built-in robes) and stylish central bathroom, spacious lounge through to the kitchen/meals area and adjoining family room. Other highlights include ducted heating and evaporative cooling, polished hardwood floors, plus a double garage (with |
| **Our model:** |
| SOLD by Frank Dowling Real Estate: Located in a quiet pocket of Niddrie, this immaculately presented family home is sure to impress. Comprising of 5 good sized bedrooms with BIRs, master with WIR and walk in robe, central bathroom, open plan living and dining area, central kitchen with stainless steel appliances, family bathroom and separate toilet. Other features include gas heating, evaporative cooling, split system air conditioner, double garage with internal access and a large rear yard with rear access. Close to schools, shops and transport. |

Table 3: An example of ground-truth and generated text on the House dataset. Here red colored text represents hallucinated information and blue colored text represents the faithful information.

most capable language models, we also compare our model with ChatGPT on a set of 1,000 random samples from the House dataset in different settings. A comprehensive analysis of this experiment is presented in Appendix B. As can be expected, ChatGPT achieves significantly better performance in faithfulness in zero-shot setting. However, when given noisy ground-truth text as few-shot examples, ChatGPT generates hallucinated text similar to the ground-truth text, showing that it is also prone to noise in the reference text. Our model outperforms ChatGPT in this (3-shot) setting in terms of precision and hallucination (i.e., lower hallucination).

### 4.6 Ablation Studies

To investigate the effect of contrastive learning and control token techniques individually, we experiment on both datasets with two configurations of our full model: one with control token only and the other one with contrastive learning only.

As we see in Table 2, both model components

contribute to our model's better faithfulness, with contrastive learning making a larger impact in House dataset.

## 5 Conclusion

In this paper, we have proposed a novel approach to generate faithful text from a knowledge graph having noisy ground-truth text. To ensure faithful text generation, we have introduced two key ideas: (i) contrastive learning to better differentiate between faithful and hallucinated information, (ii) control token to regulate the level of hallucination in the generated text. Experimental results on two noisy KG-to-text datasets demonstrates that KG-to-text model with our framework outperforms all the baseline models in terms of faithfulness metrics. Moreover, we have proposed a novel Chat-GPT based evaluation technique for an in-depth quantitative and qualitative analysis, which further verifies the superior performance of our model on

precision, recall and faithfulness.

**Limitation and Future work** We have applied our proposed framework only in PLM based KG-to-text encoder-decoder model. In future, we plan to explore the hallucination problem in AMR (Abstract Meaning Representations) graph datasets, which can also preserve a number of meaningful semantic relations and widely used in NLP areas.

## Ethical Considerations

Our model utilizes existing pre-trained language model based KG-to-text generation model, thus the ethical concerns associated with these models would also be applicable to our proposed framework.

## Acknowledgments

## References

Oshin Agarwal, Heming Ge, Siamak Shakeri, and Rami Al-Rfou. 2021. Knowledge graph based synthetic corpus generation for knowledge-enhanced language model pre-training. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3554–3565.

Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72.

Shuyang Cao and Lu Wang. 2021. Cliff: Contrastive learning for improving faithfulness and factuality in abstractive summarization. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6633–6649.

Wenhu Chen, Yu Su, Xifeng Yan, and William Yang Wang. 2020. Kgpt: Knowledge-grounded pre-training for data-to-text generation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8635–8648.

Yu Chen, Lingfei Wu, and Mohammed J Zaki. 2023. Toward subgraph-guided knowledge graph question generation with graph neural networks. *IEEE Transactions on Neural Networks and Learning Systems*.

Anthony Colas, Mehrdad Alvandipour, and Daisy Zhe Wang. 2022. Gap: A graph-aware language model framework for knowledge graph-to-text generation. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 5755–5769.

Anthony Colas, Ali Sadeghian, Yue Wang, and Daisy Zhe Wang. 2021a. Eventnarrative: A large-scale event-centric dataset for knowledge graph-to-text generation. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual.*

Anthony Colas, Ali Sadeghian, Yue Wang, and Daisy Zhe Wang. 2021b. Eventnarrative: A large-scale event-centric dataset for knowledge graph-to-text generation. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual.*

Sarkar Snigdha Sarathi Das, Mohammed Eunus Ali, Yuan-Fang Li, Yong-Bin Kang, and Timos Sellis. 2021. Boosting house price predictions using geo-spatial network embedding. *Data Mining and Knowledge Discovery*, 35:2221–2250.

Nouha Dziri, Ehsan Kamalloo, Sivan Milton, Osmar Zaiane, Mo Yu, Edoardo M Ponti, and Siva Reddy. 2022. Faithdial: A faithful benchmark for information-seeking dialogue. *Transactions of the Association for Computational Linguistics*, 10:1473–1490.

Nouha Dziri, Andrea Madotto, Osmar R Zaiane, and Avishek Joey Bose. 2021. Neural path hunter: Reducing hallucination in dialogue systems via path grounding. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2197–2214.

Katja Filippova. 2020. Controlled hallucinations: Learning to generate faithfully from noisy data. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 864–870.

Mingqi Gao and Xiaojun Wan. 2022. Dialsummeval: Revisiting summarization evaluation for dialogues. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5693–5709.

Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. Simcse: Simple contrastive learning of sentence embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6894–6910.

Zhichao Geng, Ming Zhong, Zhangyue Yin, Xipeng Qiu, and Xuan-Jing Huang. 2022. Improving abstractive dialogue summarization with speaker-aware supervised contrastive learning. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 6540–6546.

Jian Guan, Yansen Wang, and Minlie Huang. 2019. Story ending generation with incremental encoding and commonsense knowledge. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6473–6480.

Qipeng Guo, Zhijing Jin, Xipeng Qiu, Weinan Zhang, David Wipf, and Zheng Zhang. 2020. Cyclegt: Unsupervised graph-to-text and text-to-graph generation via cycle training. In *Proceedings of the 3rd International Workshop on Natural Language Generation from the Semantic Web (WebNLG+)*, pages 77–88.

Jiuzhou Han and Ehsan Shareghi. 2022. Self-supervised graph masking pre-training for graph-to-text generation. In *Empirical Methods in Natural Language Processing 2022*, pages 4845–4853. Association for Computational Linguistics (ACL).

Haozhe Ji, Pei Ke, Shaohan Huang, Furu Wei, Xiaoyan Zhu, and Minlie Huang. 2020. Language generation with multi-hop reasoning on commonsense knowledge graph. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 725–736.

Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2022. Survey of hallucination in natural language generation. *ACM Computing Surveys*.

Zhijing Jin, Qipeng Guo, Xipeng Qiu, and Zheng Zhang. 2020. Genwiki: A dataset of 1.3 million content-sharing text and graphs for unsupervised graph-to-text generation. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 2398–2409.

Mihir Kale and Abhinav Rastogi. 2020. Text-to-text pre-training for data-to-text tasks. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 97–102.

Pei Ke, Haozhe Ji, Yu Ran, Xin Cui, Liwei Wang, Linfeng Song, Xiaoyan Zhu, and Minlie Huang. 2021. Jointgt: Graph-text joint representation learning for text generation from knowledge graphs. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 2526–2538.

Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. 2020. Supervised contrastive learning. *Advances in neural information processing systems*, 33:18661–18673.

Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.

Tom Kocmi and Christian Federmann. 2023. Large language models are state-of-the-art evaluators of translation quality. *arXiv e-prints*, pages arXiv–2302.

Rik Koncel-Kedziorski, Dhanush Bekal, Yi Luan, Mirella Lapata, and Hannaneh Hajishirzi. 2019. Text generation from knowledge graphs with graph transformers. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2284–2293.

Kalpesh Krishna, Aurko Roy, and Mohit Iyyer. 2021. Hurdles to progress in long-form question answering. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4940–4957.

Wojciech Kryściński, Bryan McCann, Caiming Xiong, and Richard Socher. 2020. Evaluating the factual consistency of abstractive text summarization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9332–9346.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*.

Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.

Yixin Liu and Pengfei Liu. 2021. Simcls: A simple framework for contrastive learning of abstractive summarization. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 1065–1072.

Zheheng Luo, Qianqian Xie, and Sophia Ananiadou. 2023. Chatgpt as a factual inconsistency evaluator for abstractive text summarization. *arXiv preprint arXiv:2303.15621*.

Edward Ma. 2019. Nlp augmentation. https://github.com/makcedward/nlpaug.

Jonathan Mallinson, Rico Sennrich, and Mirella Lapata. 2017. Paraphrasing revisited with neural machine translation. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 881–893.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

Liam van der Poel, Ryan Cotterell, and Clara Meister. 2022. Mutual information alleviates hallucinations in abstractive summarization. In *EMNLP 2022*. arXiv.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551.

Hannah Rashkin, David Reitter, Gaurav Singh Tomar, and Dipanjan Das. 2021. Increasing faithfulness in knowledge-grounded dialogue with controllable features. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 704–718.

Vikas Raunak, Arul Menezes, and Marcin Junczys-Dowmunt. 2021. The curious case of hallucinations in neural machine translation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1172–1183.

Clément Rebuffel, Marco Roberti, Laure Soulier, Geoffrey Scoutheeten, Rossella Cancelliere, and Patrick Gallinari. 2022. Controlling hallucinations at word level in data-to-text generation. *Data Mining and Knowledge Discovery*, pages 1–37.

Leonardo F. R. Ribeiro, Yue Zhang, Claire Gardent, and Iryna Gurevych. 2020a. Modeling global and local node contexts for text generation from knowledge graphs. *Transactions of the Association for Computational Linguistics*, 8:589–604.

Leonardo FR Ribeiro, Martin Schmitt, Hinrich Schütze, and Iryna Gurevych. 2021. Investigating pretrained language models for graph-to-text generation. In *Proceedings of the 3rd Workshop on Natural Language Processing for Conversational AI*, pages 211–227.

Leonardo FR Ribeiro, Yue Zhang, Claire Gardent, and Iryna Gurevych. 2020b. Modeling global and local node contexts for text generation from knowledge graphs. *Transactions of the Association for Computational Linguistics*, 8:589–604.

Apoorv Saxena, Aditay Tripathi, and Partha Talukdar. 2020. Improving multi-hop question answering over knowledge graphs using knowledge base embeddings.

In *Proceedings of the 58th annual meeting of the association for computational linguistics*, pages 4498–4507.

Yixuan Su, David Vandyke, Sihui Wang, Yimai Fang, and Nigel Collier. 2021. Plan-then-generate: Controlled data-to-text generation via planning. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 895–909.

Xiangru Tang, Arjun Nair, Borui Wang, Bingyao Wang, Jai Desai, Aaron Wade, Haoran Li, Asli Celikyilmaz, Yashar Mehdad, and Dragomir Radev. 2022. Confit: Toward faithful dialogue summarization with linguistically-informed contrastive fine-tuning. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5657–5668.

Ran Tian, Shashi Narayan, Thibault Sellam, and Ankur P. Parikh. 2019. Sticking to the facts: Confident decoding for faithful data-to-text generation. *CoRR*.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *International Conference on Learning Representations*.

David Wan and Mohit Bansal. 2022. Factpegasus: Factuality-aware pre-training and fine-tuning for abstractive summarization. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1010–1028.

Jiaan Wang, Yunlong Liang, Fandong Meng, Haoxiang Shi, Zhixu Li, Jinan Xu, Jianfeng Qu, and Jie Zhou. 2023. Is chatgpt a good nlg evaluator? a preliminary study. *arXiv preprint arXiv:2303.04048*.

Peng Wang, Junyang Lin, An Yang, Chang Zhou, Yichang Zhang, Jingren Zhou, and Hongxia Yang. 2021. Sketch and refine: Towards faithful and informative table-to-text generation. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4831–4843.

Tianshu Wang, Faisal Ladhak, Esin Durmus, and He He. 2022. Improving faithfulness by augmenting negative summaries from fake documents. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11913–11921.

Zhenyi Wang, Xiaoyang Wang, Bang An, Dong Yu, and Changyou Chen. 2020. Towards faithful neural table-to-text generation with content-matching constraints. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1072–1086.

Xinnuo Xu, Ondřej Dušek, Verena Rieser, and Ioannis Konstas. 2021. Agggen: Ordering and aggregating while generating. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1419–1434.

Jinyu Yang, Jiali Duan, Son Tran, Yi Xu, Sampath Chanda, Liqun Chen, Belinda Zeng, Trishul Chilimbi, and Junzhou Huang. 2022. Vision-language pre-training with triple contrastive learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15671–15680.

Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021. Bartscore: Evaluating generated text as text generation. *Advances in Neural Information Processing Systems*, 34:27263–27277.

Dejiao Zhang, Feng Nan, Xiaokai Wei, Shang-Wen Li, Henghui Zhu, Kathleen Mckeown, Ramesh Nallapati, Andrew O Arnold, and Bing Xiang. 2021. Supporting clustering with contrastive learning. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5419–5430.

Haopeng Zhang, Semih Yavuz, Wojciech Kryściński, Kazuma Hashimoto, and Yingbo Zhou. 2022. Improving the faithfulness of abstractive summarization via entity coverage control. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 528–535.

Xueliang Zhao, Wei Wu, Can Xu, Chongyang Tao, Dongyan Zhao, and Rui Yan. 2020. Knowledge-grounded dialogue generation with pre-trained language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3377–3390.

Chunting Zhou, Graham Neubig, Jiatao Gu, Mona Diab, Francisco Guzmán, Luke Zettlemoyer, and Marjan Ghazvininejad. 2021. Detecting hallucinated content in conditional neural sequence generation. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1393–1404.

Hao Zhou, Tom Young, Minlie Huang, Haizhou Zhao, Jingfang Xu, and Xiaoyan Zhu. 2018. Commonsense knowledge aware conversation generation with graph attention. In *IJCAI*, pages 4623–4629.

## A    Prompt Design for ChatGPT-based Evaluation

The prompt templates are shown in Figure 5.

**Listing the facts of a graph:**   Here we give ChatGPT an input linearized graph and ask it to "list the features one by one from the INPUT" (Figure 5-Left). Figure 6 shows an example of this prompt to ChatGPT and its response for a sample from the House test set. ChatGPT has made no error in all 50 test samples of House data.

**Listing the common facts:**   ChatGPT was unable to correctly list the common facts between the linearized input graph and the generated text. Hence, we prompt ChatGPT for each fact listed in the input, whether that fact is included in the output. Here, each fact (or "feature") represents a single triple of the input linearized graph (Figure 5-Middle). Then, we count the answer with a "yes" response from ChatGPT. On average, ChatGPT makes 2-3 mistakes per sample. Figure 7 shows an example of this prompt and and ChatGPT's response. The red colored text indicates the mistakes done by ChatGPT.

**Listing the hallucinated facts:**   Here, we prompt ChatGPT to list both the extrinsic and intrinsic hallucination facts in the generated text by providing ChatGPT with an input (linearized graph) and an output (generated text). Firstly, to list the extrinsic hallucination facts we instruct ChatGPT to "List the features one by one from the OUTPUT that is not mentioned in the INPUT". Secondly, to list the intrinsic hallucination facts we instruct ChatGPT to "List the features one by one from the OUTPUT that is contradictory to the INPUT" (Figure 5-Right). Here, ChatGPT makes no mistakes in the 50 House test samples. Figure 8 illustrates the steps with an example and ChatGPT's response.

## B    Comparing Our Result with ChatGPT

We randomly take 1000 sample graphs from the House dataset. Our experiments are conducted using the API of Chat- GPT (gpt-3.5-turbo) model. We input ChatGPT the sample graphs in a linearized format and asked to summarize the linearized graphs in a real-estate advertising format. We experiment with ChatGPT-ZeroShot (without giving any reference text) , ChatGPT-$k$-FewShot, (where $k$ represents the number of noisy ground-truth text sample is given to ChatGPT as a reference in addition to the input linearized graph) and compare these with our full model.

Table 4 shows that in terms of faithfulness metrics (BARTScore), ChatGPT-ZeroShot has the best performance. This is because, ChatGPT is a large model and ChatGPT-ZeroShot generates text without taking any noisy ground-truth text as a reference. Whereas, our model is a small (BART-base/T5-base) language model and the model is trained with the full noisy training House dataset. We also notice that the performance of ChatGPT-$k$-FewShot drops with the increase of number of noisy reference text samples. Thus, the more we increase the number of noisy ground-truth texts as a reference to ChatGPT, the more ChatGPT generates hallucinated text similar to ground-truth text. That's why the BLEU, METEOR and ROUGE-L scores increase and BARTscore, FactCC scores decrease with the increase of few shot samples.

We also compare the results using ChatGPT-based evaluation. Table 5 shows the average of precision, recall and hallucinations which we compute using ChatGPT. The results also show that ChatGPT-ZeroShot performs best in all metrics as usual. Our model outperforms ChatGPT-3-FewShot in terms of precision (higher precision) and hallucination (lower hallucination).

**Performance Based on Salient Facts:**   We rank in descending order the features (type-wise) of the house graph based on their frequency of occurrence in the House trainining dataset. We take top ten features as *salient* facts. The salient facts are: 1) house_location, 2) house_property-type, 3) num. of bedrooms, 4) num. of bathrooms, 5) num of parking spaces, 6) has_ac, 7) has_dining, 8) has_heating, 9) has garage_spaces and 10) nearest_train_station. Using ChatGPT, we enumerate the presence of these facts and measure salient precision, $P_{salient}$ and salient recall, $R_{salient}$ as follows.

$$P_{salient} = \frac{\#\ salient\ common\ facts}{\#\ output\ facts} \tag{4}$$

$$R_{salient} = \frac{\#\ salient\ common\ facts}{\#\ salient\ input\ facts} \tag{5}$$

The results from Table: 6 shows that our model achieves the best average salient precision, $P_{salient}$, and ChatGPT-ZeroShot achieves the best average salient recall. The reason behind this result is that ChatGPT-ZeroShot generated output text contains mostly all the facts from the input graph, whereas

Figure 5: Prompt templates for enumerating facts using ChatGPT. Template-1 (left) is to enumerate facts in the input (linearized graph). Template-2 (middle) is to enumerate common facts between the input (linearized graph) and the output (generated text). Template-3 (right) is to enumerate hallucinated facts in the output (generated text).



Figure 6: An Example of Template-1



Figure 7: An Example of Template-2

Figure 8: An Example of Template-3

| Generation Model | Comparison with ground-truth text | | | Comparison with linearized graph | |
|---|---|---|---|---|---|
| | BLEU ↑ | METEOR ↑ | ROUGE-L ↑ | BARTScore ↑ | FactCC ↑ |
| ChatGPT-ZeroShot | 1.21 | 11.86 | 12.91 | **-2.389** | 71.02 |
| ChatGPT-1-Shot | 1.95 | 12.73 | 15.02 | -2.872 | **76.34** |
| ChatGPT-2-Shot | 2.06 | 12.67 | 15.58 | -2.937 | 72.02 |
| ChatGPT-3-Shot | 2.25 | **13.31** | 15.76 | -3.036 | 73.88 |
| **Our Full Model** | **2.68** | 11.21 | **17.10** | -3.246 | 62.84 |

Table 4: Results on 1000 test samples from the House dataset. **Bold** fonts denote the best results.

| Generation Model | Avg. Precision | Avg. Recall | Avg. Hallucination |
|---|---|---|---|
| ChatGPT-ZeroShot | **73.28** | **88.21** | **26.71** |
| ChatGPT-3-Shot | 65.45 | 64.39 | 34.55 |
| **Our Full Model** | 67.06 | 58.81 | 32.94 |

Table 5: ChatGPT Evaluation Results based on 50 samples from the House Dataset. **Bold** fonts denote the best results.

our model generated output text gives more focus on the salient facts.

# C   Generated Samples

Figure 9 and Figure 10 show qualitative examples of sample graphs, the ground-truth texts and the texts generated by different models on House dataset and Genwiki dataset, respectively.

| Generation Model | Avg. Salient Precision | Avg. Salient Recall |
|---|---|---|
| ChatGPT-ZeroShot | 26.75 | **92.66** |
| ChatGPT-3-FewShot | 30.27 | 86.36 |
| **Our Full Model** | **31.64** | 77.16 |

Table 6: ChatGPT Evaluation Results based on 50 samples from the House dataset considering salient features. **Bold** fonts denote the best results.

**House Knowledge Graph:**



**Ground-truth text:**

AUCTION THIS SATURDAY 10TH OCTOBER AT 1PM: Superbly blending charm with a modern sensibility and fresh infusion of light, this fully renovated three bedroom Californian Bungalow offers stylish family living in a sought after locale.An expansive open plan lounge and dining room provides spacious living while showcasing a marvellous kitchen boasting 700mm stainless steel appliances, abundant storage and feature island bench, ensuring effortless entertaining.A serene second living room, boasting leafy garden views offers retreat from the bustle of family life.Three spacious bedrooms are serviced by a modern bathroom, which cleverly incorporates an oversized European laundry.Brilliantly positioned to transport, schools, cafes, vibrant Sydney Road, CityLink and featuring ducted heating, alarm, split system air-conditioning, BIR's, ceiling fans and carport with OSP, stylish family living starts here.

**JointGT model:**

AUCTION THIS SATURDAY AT 11:00AM: This renovated three bedroom brick veneer home is sure to appeal to first home buyers, downsizers and investors alike.The home comprises of three double bedrooms (all with built-in robes), spacious lounge and dining room, renovated kitchen with stainless steel appliances, separate laundry, ducted heating, evaporative cooling, alarm system and a single garage with internal access.The property also features a large rear garden with established gardens and a bungalow.

**GAP model:**

Renovated And Ready To Move In: This beautifully renovated brick veneer home is ready to move in and enjoy. Featuring a spacious living/dining area, renovated kitchen with stainless steel appliances and adjoining meals area, three double bedrooms with built-in robes, central bathroom and separate laundry. Other features include ducted heating, evaporative cooling, alarm system, double garage with internal access, polished floorboards, polished hardwood floors, security shutters and much more. Close to Sydney Road shops, schools, public transport and easy access to CBD.

**GMP model:**

SOLD BY HARRY KONTOSSIS: Sale by SET DATE 10/10/2015 (unless sold prior) Situated in a quiet cul-de-sac in the heart of Brunswick West, this immaculately presented three bedroom, two bathroom home is sure to appeal to first home buyers, investors and downsizers alike.Accommodation includes a formal lounge and dining room, well-appointed kitchen with stainless steel appliances, adjoining family/meals area, three double bedrooms with built-in robes, central bathroom, separate toilet.

**Our framework:**

Renovated Family Home in the Heart of Brunswick West: This beautifully renovated family home is set on a generous allotment of 444 sqm (approx) and is ideal for first home buyers, downsizers or investors looking for a low maintenance lifestyle. Comprising of 3 spacious bedrooms (master with WIR and walk in robe), formal lounge and dining room, updated kitchen with stainless steel appliances, central bathroom and laundry. Other features include ducted heating, evaporative cooling, split system air conditioner, ducted cooling, alarm system, double garage and much more.

Figure 9: Example of generated text on the House dataset. Here red colored text represents hallucinated information and blue colored text represents the faithful information

121

**Genwiki Knowledge Graph:**

**Ground-truth text:**

Daniel Monte McLain ( May 11 , 1955 – November 8 , 1995 ) , known by the stage name Country Dick Montana , was a musician best known as a member of The Beat Farmers . Montana was born in Carmel , California .

**JointGT model:**

Montana was born on May 11 , 1955 in Carmel , California .

**CycleGT model:**

Daniel Monte McLain ( May 11 , 1955 in Carmel , Montana – November 8 , 1995 in Carmel , California ) was a musician , best known as the founder of the band Country Dick Montana .

**GMP model:**

Daniel Monte McLain ( May 11, 1955 – November 8, 1995 ), known professionally as Country Dick Montana, was an American singer, songwriter, and musician.

**Our framework:**

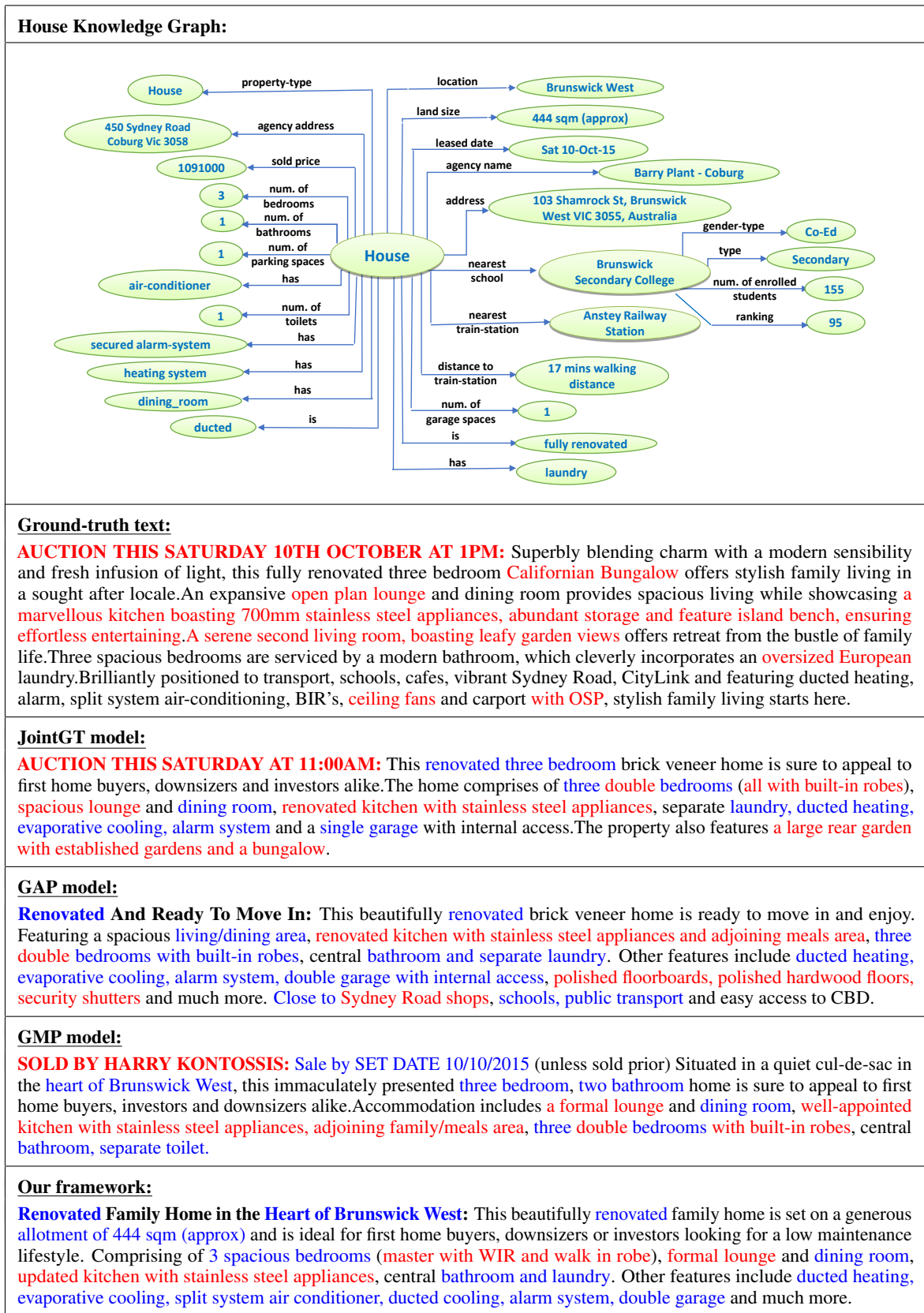Daniel Monte McLain ( May 11 , 1955 – November 8 , 1995 ) was an American musician .

Figure 10: Example of generated text on the Genwiki dataset. Here red colored text represents hallucinated information and blue colored text represents faithful information.

# Entropy-based Sampling for Abstractive Multi-document Summarization in Low-resource Settings

**Laura Mascarell** and **Ribin Chalumattu** and **Julien Heitmann**
ETH Zurich
{lmascarell,cribin,julien.heitmann}@inf.ethz.ch

## Abstract

Research in Multi-document Summarization (MDS) mostly focuses on the English language and depends on large MDS datasets that are not available for other languages. Some of these approaches concatenate the source documents, resulting in overlong model inputs. Existing transformer architectures are unable to process such long inputs entirely, omitting documents in the summarization process. Other solutions address this issue by implementing multi-stage approaches that also require changes in the model architecture. In this paper, we introduce various sampling approaches based on information entropy that allow us to perform MDS in a single stage. These approaches also consider all source documents without using MDS training data nor changing the model's architecture. Besides, we build a MDS test set of German news articles to assess the performance of our methods on abstractive multi-document summaries. Experimental results show that our entropy-based approaches outperform previous state-of-the-art on German MDS, while still remaining primarily abstractive. We release our code[1] and MDS test set[2] to encourage further research in German abstractive MDS.

## 1 Introduction

In light of the ever-growing volume of available information, it becomes essential to be able to automatically summarize information from several sources. Multi-document Summarization (MDS) aims at condensing the most important information from different documents. Despite the advances in single-document summarization (Zhang et al., 2020), summarizing multiple related documents remains a greater challenge due to its input length and the presence of redundant information (Fan et al., 2019; Song et al., 2022). Therefore, some research focuses on implementing multi-stage approaches

that first identify the relevant information to then feed it into a summarization model (Lebanoff et al., 2018; Liu and Lapata, 2019a). More recent works utilize pre-trained language models (Lewis et al., 2020; Raffel et al., 2020; Xiao et al., 2022) fine-tuned for the summarization task and feed them with the source documents concatenated (Johner et al., 2021; Xiao et al., 2022). However, these approaches pose two major issues. First, concatenated inputs exceeding the length limit of the model are truncated, which might lead to the omission of entire documents. Second, they rely on multi-document datasets that are scarce or unavailable in languages other than English.

Hokamp et al. (2020) introduce a decoding strategy that adapts single- to multi-document summarization without using additional training data nor applying changes to the single-input model architecture. At every decoding timestep, it averages the output probabilities of a single-document summarization model for each individual document, combining them into a single output. Instead of averaging all log-probabilities, which favours highly frequent tokens, we propose to make a more informed decision. In particular, we leverage entropy to measure the model confidence in the next token prediction and thus select the most informative output. We implement different entropy-based approaches and evaluate their performance on MDS of German text. Our main contributions are:

- We present different entropy-based sampling approaches for the MDS task. These are specially well-suited for languages like German that have limited or unavailable MDS data.

- We build and release a new German MDS test set in the news domain that is more suitable for evaluating abstractive summarization than the existing MDS German dataset auto-hMDS (Zopf, 2018). We expect our dataset to foster research on German abstractive MDS.

---

[1] Link to GitHub repository.
[2] Link to Multi-GeNews repository.

- The experimental results demonstrate that our method achieves the state-of-the-art performance in German abstractive MDS in terms of ROUGE scores and manual evaluation.

## 2 Related Work

**Multi-document Summarization** Some prior work approaches MDS as a multi-stage process (Liu et al., 2018; Zhu et al., 2021) that first extracts salient sentences from the source documents to then distill them into a summary using different methods such as graph-based modeling (Li et al., 2020; Chen et al., 2021) or modifying the attention mechanism (Perez-Beltrachini and Lapata, 2021). In contrast, Lebanoff et al. (2018) highlights the importance of adapting Single-document Summarization (SDS) models to summarize multiple documents and propose an approach that adapts their attention weights. Similarly, other works propose various changes in the model architecture (Liu and Lapata, 2019a; Elsahar et al., 2021). The main disadvantage of these approaches is that they are tailored to specific model architectures.

More recently, Xiao et al. (2022) introduce PRIMERA, a pre-trained model for MDS that can be applied in zero- or few-shot settings. The source documents are concatenated and fed into the Longformer Transformer model, which can handle long inputs up to 4,096 or even 16k tokens with current GPUs. Nevertheless, PRIMERA is only available for English and there is no alternative for other languages. Similarly, Johner et al. (2021) performs MDS on German text using the pre-trained language model BART (Lewis et al., 2020) and concatenating the source documents as input.[3] However, BART input length is restricted to 1,024 tokens, which may end up excluding entire documents from the summarization process.

Overall, our entropy-based approaches present the following advantages over prior work: (a) they do not require a pre-step to extract salient information (b) nor changes in the SDS model architecture with (c) no need for additional MDS training data, and (d) still considering all source documents in the summarization process. This work is built upon the dynamic ensemble approach from Hokamp et al. (2020), improving the decoding strategy by sampling on more informative predictions.

---

[3]To the best of our knowledge, Johner et al. (2021) is the only work that tackles MDS in German besides Zopf (2018) with the auto-hMDS dataset.

**Entropy in Summarization** Xu et al. (2020) leverage entropy to analyze the performance of Transformer-based models in the SDS task. Later, van der Poel et al. (2022) use entropy to determine when the model is uncertain about the next token prediction and apply Pointwise Mutual Information (PMI) instead to alleviate hallucination. Similarly, we apply the conditional PMI approach to MDS. Instead of finding the conditional entropy threshold through hyperparameter search as in van der Poel et al. (2022), we apply maximum probabilistic information entropy (Li et al., 2021). This novel entropy definition has been successfully used to reduce the size of image datasets by selecting the most informative samples.

## 3 Entropy Background

In information theory, the entropy of a random variable denotes the amount of information, or lack thereof (i.e. uncertainty), associated with its possible outcomes. Thus, given a probability distribution $p$ over all possible outcomes $x_1, \ldots, x_n$ of a random variable $X$, we quantify the entropy of $X$ using the standard Shannon entropy equation:

$$H(X) = -\sum_{i=1}^{n} p(x_i) \, \log p(x_i) \qquad (1)$$

The entropy is then maximum for uniform distributions, where all outcomes are equally likely, indicating high uncertainty.

In the context of automatic text generation, we can leverage entropy to quantify the confidence of probabilistic models in their predictions (Xu et al., 2020; van der Poel et al., 2022). More specifically, summarization models aim at generating a summary string $\mathbf{y}^*$ of a given source document $\mathbf{x}$ that maximizes the scoring function:

$$\mathbf{y}^* = \underset{\mathbf{y} \in \mathcal{Y}}{\mathrm{argmax}} \; \log p(\mathbf{y} \mid \mathbf{x}), \qquad (2)$$

where $\mathbf{y}$ is the sequence of tokens $y_0, \ldots, y_T$ from the model vocabulary $\mathcal{V}$ generated at every timestep $t$, $0 < t < T$. During decoding, that is, the prediction of each sequence token $y_t \in \mathcal{V}$, the model provides a probability distribution $p(\cdot \mid \mathbf{y}_{<t}, \mathbf{x})$ over $\mathcal{V}$ that also takes into account the context of the previous tokens. According to Equation 1, we can then use such distribution to measure the model's confidence in the prediction:

$$H(p(\cdot \mid \mathbf{y}_{<t}, \mathbf{x})) = -\sum_{y \in \mathcal{V}} \Big( p(y \mid \mathbf{y}_{<t}, \mathbf{x}) \qquad (3)$$
$$\times \log p(y \mid \mathbf{y}_{<t}, \mathbf{x}) \Big)$$

## 4 Entropy-based MDS

Given a set of documents $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, the dynamic-ensemble approach (*DynE*) described in Hokamp et al. (2020) adapts single- to multi-document summarization as follows: at every decoding timestep $t$, it computes the output probabilites for each individual source document using a single-document summarization model; next, it averages these outputs to obtain a single log-probability distribution assigned to the token $y$:

$$p(y \mid \mathcal{X}) = \frac{1}{|\mathcal{X}|} \sum_{\mathbf{x} \in \mathcal{X}} p(y \mid \mathbf{y}_{<t}, \mathbf{x}) \qquad (4)$$

We leverage entropy information to adapt the *DynE* approach and implement various sampling strategies that select the most informative output at each decoding timestep $t$.

**Minimum Entropy** ($H_{min}$)  Based on the hypothesis that low entropy indicates a higher confidence in the prediction, this approach picks the token prediction of the model instance with the lowest entropy $\min_{1 \le i \le |\mathcal{X}|} H(p(\cdot \mid \mathbf{y}_{<t}, \mathbf{x}_i))$. Note that this approach does not guarantee certainty in all token predictions. In those cases where all model instances exhibit high uncertainty, the selected instance could still have high entropy and thus provide an arbitrary prediction.

**Max-predicted Probability Threshold** ($H_{th}$)  Li et al. (2021) focus on the maximum-predicted probability $p_{max}$ to determine the model's confidence and reduce redundancy in datasets of images. Specifically, the authors state that a low maximum probability indicates high entropy and consequently, low confidence in the prediction. Therefore, they propose to measure entropy as:

$$H(X) = -p_{\max} \log p_{\max} \qquad (5)$$

where $p_{\max} = \max_{x \in X} p(x)$. Figure 1 plots Equation 5, showing the correlation between information entropy and $p_{max}$. Note that the entropy is highest when $p_{max}$ is 0.35, with a positive correlation for probabilities below this threshold and a negative correlation for probabilities above it.



Figure 1: Plot of the maximum probabilistic information entropy (Equation 5), which illustrates the correlation between the maximum-predicted probability $p_{max}$ and information entropy, when $0 \le p_{max} \le 1$.

Inspired by Li et al. (2021) approach, we apply maximum probabilistic information entropy in MDS, assuming that values of $p_{max}$ below the threshold 0.35 indicate that the model is essentially guessing. At each decoding step, we obtain the maximum-predicted probability for each input document in $\mathcal{X}$ and proceed as follows:

a) we choose the prediction with the highest $p_{max}$ among those above the threshold. The higher the probability, the lower the entropy.

b) if all probabilities are below the threshold, we conclude that there is not enough information for the current prediction and we average their log-probabilities as in Equation 4.

**Mutual Information Decoding** ($H_{pmi}$)  Several works apply mutual information approaches during decoding to favor more specific and informative outputs (Li et al., 2016; Takayama and Arase, 2019). Later, van der Poel et al. (2022) observe that highly frequent tokens often indicate hallucinated content and implement a decoding strategy that integrates mutual information to mitigate hallucination in single-document summarization. In particular, their approach optimizes for Pointwise Mutual Information (PMI) when the model is uncertain about its prediction, generating summaries that are more faithful to the source document:

$$p(y|\mathbf{x}) = \log p(y|\mathbf{y}_{<t}, \mathbf{x}) - \lambda \log p(y|\mathbf{y}_{<t}), \quad (6)$$

where $0 < \lambda < 1$ to avoid excessively penalizing high-frequent tokens, which could lead to

ungrammatical outputs (Li et al., 2016). Based on these findings, we propose an additional variation of our $H_{th}$ approach, which applies PMI when there is no certainty in any of the predictions, that is, all probabilities are below the 0.35 threshold.[4]

## 5 Datasets

This section describes the datasets used to train and evaluate our MDS approaches. Specifically, we consider three pre-existing German datasets that are suitable for single-document—GeWiki (Frefel, 2020) and 20m (Rios et al., 2021)—and multi-document summarization—auto-hMDS (Zopf, 2018). Moreover, we build Multi-GeNews, a MDS test set in the news domain that is specifically tailored for abstractive MDS.

### 5.1 Single-document Summarization

**GeWiki** This is the largest dataset available for single-document abstractive summarization in German, consisting of 240k summary-article pairs. Here, the lead text of Wikipedia articles are extracted as summaries of the rest of the article.

**20m** A single-document summarization dataset with 18,305 news articles and their corresponding manually-written summaries collected from the Swiss newspaper *20 Minuten* ('20 Minutes').

### 5.2 Multi-document Summarization

**auto-hMDS** This multi-document summarization dataset consists of 2,210 summaries from Wikipedia leads as in GeWiki and 10,454 source documents. The documents were obtained by automatically querying the internet with summary sentences, resulting in a highly extractive dataset. Nonetheless, we consider it in our experiments for comparison with the related work. Despite being the largest MDS dataset in German, auto-hMDS is significantly smaller than its English counterpart Multi-News (Fabbri et al., 2019).[5]

**Multi-GeNews** Due to the lack of abstractive MDS datasets in German, we built a MDS test set to assess the performance of the proposed approaches. The data comes from the news portal of the Swiss media company SRF[6] and consists of news articles published between January and March 2020.

The articles published on the SRF website are often followed by a *Mehr zum Thema* ('More on the topic') section with related articles on the subject. To build our test set, we first utilize this section to obtain clusters of related articles. Specifically, we collect the related article suggestions and filter those published within one day of each other to ensure that they cover the same news. Next, we generate the reference summaries, which will be used to compute the automatic scores, concatenating the lead paragraphs of the articles in each cluster.[7] Hence, the reference summaries are a combination of lead texts. We finally filter salient sentences and remove duplicated information from the reference summaries using a pretrained extractive summarization model for German text. To build this model, we adapted the BertExt architecture (Liu and Lapata, 2019b) for the German language.[8] The adaption involved initializing the Bert component of the BertExt architecture using a German Bert checkpoint[9] and subsequently fine-tuning the entire model on the newswire 20m dataset.

The resulting dataset consists of 754 unique articles grouped into 402 clusters. Each cluster contains two to six articles with a median of four articles and the corresponding generated reference summary.[10] The average length of the articles and summaries are 593 and 61 tokens, respectively.

## 6 Experiments

We evaluate the performance of the entropy-based sampling approaches on our Multi-GeNews and the auto-hMDS datasets in terms of automatic ROUGE scores (Lin, 2004) and extractive fragment density $\rho$ (Grusky et al., 2018). Since we focus on abstractive summarization, the latter allows us to measure the degree of extractiveness of the summaries, and in turn, abstractiveness—higher $\rho$ values indicate that the summary is more extractive and contains larger text chunks from the source article. Furthermore, we collect human annotations on a subset of the Multi-GeNews to assess the faithfulness of the generated summaries and get a deeper understanding on their quality (Section 7).

---

[4]We use a $\lambda$ of 0.25 in our experiments, which we manually selected based on the impact of various values on the output.

[5]over 56k summaries and 250k source documents.

[6]https://www.srf.ch/news

[7]Similarly to the GeWiki dataset, we consider the lead paragraph of an article as its summary.

[8]https://github.com/nlpyang/BertSum

[9]https://huggingface.co/dbmdz/bert-base-german-uncased

[10]Although an article can belong to different clusters, there are no identical clusters with the same articles.

| Method | 100 words | | | | 200 words | | | |
|--------|-----------|---|---|---|-----------|---|---|---|
| | **R1**↑ | **R2**↑ | **RL**↑ | $\rho$↓ | **R1**↑ | **R2**↑ | **RL**↑ | $\rho$↓ |
| mBART concat | 18.4 | 6.2 | 12.5 | 27.9 | 24.5 | 7.7 | 15.0 | 35.6 |
| mBART + $DynE$ | **23.4** | 6.9 | 15.1 | 2.2 | 26.8 | 7.0 | 16.0 | 1.9 |
| mBART + $H_{min}$ | 20.7 | 8.6 | 14.7 | 17.9 | 26.9 | 10.4 | 17.4 | 16.6 |
| mBART + $H_{th}$ | 21.5 | **9.0** | **15.3** | 16.1 | **27.8** | **10.8** | **18.0** | 14.7 |
| mBART + $H_{pmi}$ | 16.5 | 6.9 | 12.3 | 12.5 | 21.0 | 7.9 | 14.5 | 10.0 |

Table 1: Performance of the entropy-based approaches and the baseline models on the **auto-hMDS** dataset in terms of ROUGE scores and extractive fragment density $\rho$. The mBART model is fine-tuned on the auto-hMDS dataset by concatenating the source articles into a single input. Similarly, the mBART baseline is fed with the concatenated source articles. Overall, $H_{th}$ achives the highest performance among the various methods evaluated.

| Method | 100 words | | | |
|--------|-----------|---|---|---|
| | **R1**↑ | **R2**↑ | **RL**↑ | $\rho$↓ |
| mBART concat | 23.0 | 6.0 | 14.8 | 9.23 |
| mBART + $DynE$ | 22.2 | 4.8 | 14.9 | 1.5 |
| mBART + $H_{min}$ | 23.4 | 5.6 | 15.0 | 2.46 |
| mBART + $H_{th}$ | **24.5** | 6.2 | 15.6 | 2.72 |
| mBART + $H_{pmi}$ | 23.9 | **7.2** | **16.1** | 2.78 |

Table 2: Performance of the entropy-based approaches and the baselines on our **Multi-GeNews** test set. The mBART model is fine-tuned on the 20m dataset as described in Section 6.1. The mBART baseline receives as input the source articles concatenated.

## 6.1 Models

This section describes the implementation details to build the models used in our experiments. Namely, the two summarization models, individually fine-tuned on the newswire 20m and the auto-hMDS datasets, and the language model used by the pointwise mutual information decoding approach $H_{pmi}$.

**Summarization Models** We evaluate the performance of our MDS approaches using two summarization models fine-tuned on the news domain dataset 20m[11] and the MDS dataset auto-hMDS, respectively. The latter allows us to compare the performance of our approaches against prior work on German MDS. The models are based on mBART, a multilingual sequence-to-sequence transformer-based model that effectively handles multiple lan-

guages including German (Liu et al., 2020) and initialized with the `facebook/mbart-large-cc25` checkpoint available at the Hugging Face Hub.[12]

In particular, we fine-tune the model on the 20m dataset for 10 epochs and batch size of 2 using the AdamW optimizer (Loshchilov and Hutter, 2019) with a learning rate of $3e - 5$. The gradient accumulation steps is 16, resulting in a total effective batch size of 32. To fine-tune the single-input mBART model with the multi-document summarization dataset auto-hMDS, we follow the work in Johner et al. (2021) and concatenate the source articles in a single input. We train the model with $3e - 5$ learning rate and batch size of 2 for 5 epochs.

**Language Model** We build a language model to apply the mutual information decoding approach. Specifically, we use the GPT-2 (Radford et al., 2019) checkpoint from Hugging Face[13] and fine-tune it on the same in-domain data as the corresponding mBART summarization model. To ensure that both mBART and GPT-2 models share the same vocabulary, we train GPT-2 using the same tokenizer as mBART. We then fine-tune it for 3 epochs using the AdamW optimizer with learning rate of $5e - 4$ and batch size of 16. We set the maximum context length of the model to 256 tokens, since we do not generate longer summaries than that. The gradient accumulation steps is set to 8, resulting in a total effective batch size of 128.

## 6.2 Results

Table 2 compares the performance of our entropy-based methods against the $DynE$ (Hokamp et al., 2020) and mBART baselines on our Multi-GeWiki

---

[11]Since the GeWiki dataset is significantly larger than the in-domain 20m, we also considered to build a model using both datasets through behavioral fine-tuning. However, the performance on the single-document summarization task was inferior than simply fine-tuning on 20m. Several factors could contribute to this results such as a domain shift or a discrepancy in summary length distribution.

[12]https://huggingface.co/facebook/mbart-large-cc25
[13]https://huggingface.co/gpt2

| | | |
|---|---|---|
| Ref. | | Das Coronavirus beherrscht die Medien – doch das sei derzeit kaum angebracht, sagt Wissenschaftsredaktor Häusler. Innerhalb eines Tages kletterte die Zahl der bestätigten Infektionen in China auf 2744. In den umliegenden Ländern gibt es bereits Erkrankte.<br><br>[The coronavirus dominates the media - but that is hardly appropriate at the moment, says science editor Häusler. Within one day, the number of confirmed infections in China climbed to 2744, and there are already people ill in surrounding countries.] |
| Rank 1 | $H_{th}$ | Die Zahl der Corona-Toten in China ist auf 80 gestiegen. Die Zahl der Touristen, die sich mit dem Virus infiziert haben, ist auf überschaubar. Die Tourismusbranche rechnet nicht damit, dass das Virus sich auf der ganzen Welt ausbreitet. SRF-Wissenschaftler Thomas Häusler und Daniel Sulzer von Interlaken erklären, wie gefährlich die Situation für die chinesische Bevölkerung ist. |
| | | [The number of Covid deaths in China has risen to 80. The number of tourists infected with the virus is on manageable. The tourism industry does not expect the virus to spread around the world. SRF scientist Thomas Häusler and Daniel Sulzer from Interlaken explain how dangerous the situation is for the Chinese population.] |
| Rank 2 | $H_{pmi}$ | In China ist die Zahl der Corona-Fälle stark angestiegen. In der Volksrepublik gibt es bereits über 700 Fälle. Die Zahl der Touristen, die sich mit dem Virus infiziert haben, ist doppelt so hoch wie bei der saisonalen Grippe. Ein Tourismusexperte erklärt, wie gefährlich die Situation für die chinesische Bevölkerung ist. |
| | | [In China, the number of Covid cases has risen sharply. There are already over 700 cases in the People's Republic. The number of tourists infected with the virus is twice as high as the seasonal flu. A tourism expert explains how dangerous the situation is for the Chinese population.] |
| Rank 3 | $DynE$ | Die Zahl der Corona-Fälle in der Schweiz steigt, die Zahl der Infizierten nimmt zu. Die Gefahr, sich mit dem Virus anzustecken, ist noch nicht bedroht. Trotzdem steigen die Zahlen wieder an. Laut Experten ist die Gefahr grösser, dass sich das Virus noch weiter ausbreitet. |
| | | [The number of Covid cases in Switzerland is rising, the number of infected persons is increasing. The danger of contracting the virus is not yet threatened. Nevertheless, the numbers are rising again. According to experts, there is a greater risk that the virus will spread even further.] |
| Rank 4 | $H_{min}$ | Der Tourismusverband Chinas rechnet mit einem Rückgang der Touristenzahlen. In Interlaken und Luzern gibt es nur noch wenige Berichte über das Coronavirus. In der Schweiz gibt es aber Hoffnung: Vermehrt Japaner und Chinesen berichten von Infektionen mit dem Virus. Ein Tourismusdirektor glaubt, dass der Tourismus in der Volksrepublik eine globale Pandemie auslösen könnte. |
| | | [The Tourism Association of China expects a decline in tourist numbers. In Interlaken and Lucerne there are only few reports of the coronavirus. In Switzerland, however, there is hope: Increasing numbers of Japanese and Chinese report infections with the virus. One tourism director believes that tourism in the People's Republic could trigger a global pandemic.] |

Table 3: Example of the summary ranking task for the input articles 18126230, 18127577, and 18130289, where at least two annotators agreed on the ranking position for each summary. In contrast to the entropy-based approaches, *DynE* is susceptible to generate overly general summaries.

test set. We use the single-document summarization mBART model fine-tuned on the in-domain dataset 20min—see details in Section 6.1. To feed the mBART baseline with multiple documents, we concatenate them in a single input as in Johner et al. (2021). Overall, the automatic ROUGE scores indicate that $H_{th}$ and $H_{pmi}$ achieve the highest performance. Similarly, $H_{th}$ outperforms the other approaches on the auto-hMDS dataset (see Table 1).

**Abstractiveness of the Summaries**  Table 1 and Table 2 reveal that *DynE* summaries are the most abstractive (lowest $\rho$ scores). In contrast, concatenating the source articles as input results in highly extractive summaries,[14] and the gap is even more

significant with mBART fine-tuned on auto-hMDS, since the dataset is highly extractive (Table 1). Although we aim at generating abstractive summaries, the *DynE* approach is prone to generate highly frequent tokens,[15] resulting in general summaries that fail to consider relevant and specific information from the source articles (see example in Table 3). Instead, our entropy-based approaches generate summaries with a moderate level of abstractiveness that also include concrete information.

## 7 Human Evaluation

We recruited three native German speakers to perform a manual evaluation on the Multi-GeNews test

---

[14]The results on the extractiveness of mBART summaries are also supported in Johner et al. (2021).

[15]Since the *DynE* approach averages the log-probability outputs at each decoding step, common tokens obtain higher probabilities and are more likely to be predicted.

(a) Heatmap illustrating the distribution of relative preference among the approaches. The x-axis indicates the preferred approach over the y-axis. $H_{min}$ and $H_{th}$ are the most favoured approaches, whereas $H_{pmi}$ ranks as the least preferred.



(b) Percentage of instances where each approach ranked at the top and the bottom positions, according to the annotators. While $H_{th}$ summaries were consistently rated among the top positions, the annotators rated $H_{pmi}$ summaries low.

Figure 2: Evaluation of the quality of the summaries among the different approaches. We only consider those instances where the majority of the annotators agreed on the (a) relative or (b) absolute ranking position.

set.[16] This evaluation task is twofold: (1) assess the relative quality of the summaries, ranking them accordingly (Goyal et al., 2022; Ermakova et al., 2019) and (2) the faithfulness of the generated summaries to the source articles (Krishna et al., 2023), that is, whether the information presented in the summaries is supported by the articles.

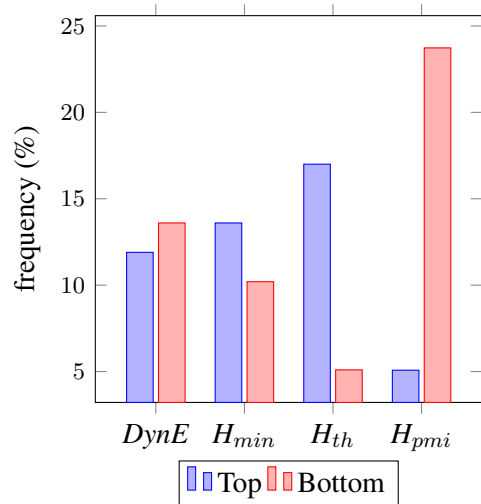Since the task requires to read a considerable amount of text, we ask the participants to annotate a sample of the MDS test set. This sample comprises 20 randomly selected instances that meet the following criteria: (a) each instance consists of three source articles, (b) the generated summaries end with a punctuation mark to avoid incomplete sentences, and (c) the token-level edit distance among the summaries is above five to ensure lexical differences. For each participant, we randomly shuffle the evaluation instances and the required annotations to avoid any biases. Additionally, we do not provide them with any information about which specific approach generated each summary.

**Summary Ranking Task** The objective of this task is to gain insights into human preferences of the generated summaries. For each instance (i.e. a set of related articles), we ask the participants to rank the generated summaries according to the

informativeness of the summaries and their preference. That is, they must evaluate how effectively the summaries are at capturing the essential information from the three source articles. Note that we are not evaluating other linguistic aspects such as cohesion or fluency, since we do not implement any specific methods to improve those. Table 3 provides an example of the annotation task.

This task is specially challenging when multiple summaries either contain similar information or suffer from hallucinations, or both. In fact, the annotators reported that it was hard to decide the rank of a summary between two consecutive positions. This negatively impacted on the inter-annotator agreement, resulting in a final Kendall's tau coefficient of 0.22. In the analysis of this task, we concentrate on the relative performance of the approaches and only consider instances with a majority agreement among annotators.

Figure 2a illustrates the relative preference among the different approaches in this ranking task. The results demonstrate that the summaries from the approaches $H_{min}$ and $H_{th}$ are consistently rated higher than the others. In contrast, the $H_{pmi}$ summaries are the least preferred. These results are also supported in Figure 2b, where we compare the frequency with which each approach was ranked within the top two and the bottom two positions. Furthermore, Figure 2b shows that while the base-

---

[16]The participants received a voucher worth CHF 75.- as compensation for their participation.

| | $H_{th}$ |
|---|---|
| F | **Donald Trump hielt sich in** der Nacht auf Mittwoch in den beiden Kammern des US-Kongresses seine dritte Rede ab. [Donald Trump delivered his third speech to both chambers of the U.S. Congress on Wednesday night.] |
| F | Die Rede ist von einem Triumphgehabe gegen die Demokraten. [There is talk of triumphant action against the Democrats.] |
| F | Das Verfahren gegen Trump ist nach wie vor im Gange. [The case against Trump is still ongoing.] |

| | $H_{pmi}$ |
|---|---|
| F | **Donald Trump hielt sich in** den USA nicht an die Corona-Regeln. [Donald Trump did not follow the Covid rules in the USA.] |
| F | Die demokratische Mehrheit im Kongress hielt sich dagegen und sprach Trump ab. [The Democratic majority in Congress held against this and absolved Trump.] |
| T | Die Rede ist von einem Triumph für Trump. [The talk is of a triumph for Trump.] |

Table 4: Example of the faithfulness annotation task. The boolean in the first column represents whether the text span is factual (T=True) or not (F=False) based on the majority agreement among the annotators. Here, the $H_{th}$ summary was ranked at the top positions of the ranking and $H_{pmi}$ at the bottom positions, even though the latter has a text span annotated as factual. The highlighted text indicate the common tokens between $H_{th}$ and $H_{pmi}$ until $H_{pmi}$ applies PMI, hallucinating on the Covid virus, although Covid is not even mentioned in the source articles.

line summaries *DynE* receive mixed ratings, the $H_{th}$ summaries are consistently ranked in the top positions. This indicates a consistent preference for the $H_{th}$ summaries over the baseline.

**Faithfulness Annotation Task** van der Poel et al. (2022) leverage PMI to improve faithfulness and evaluate it in terms of automatic metrics (Section 4). The goal of this annotation task is to manually evaluate the faithfulness of our $H_{pmi}$ approach, which applies PMI to MDS, and compare it to the other proposed approaches that do not specifically address faithfulness. Specifically, we follow the guidelines described in Krishna et al. (2023) and split the summaries into text spans to ensure lower inter-annotator variance.[17] We then ask the annotators to judge whether each span is faithful to the source articles, that is, the statements can be verified against the articles. The final Fleiss' $\kappa$ (Fleiss, 1971) inter-annotator agreement is 0.62.

Overall, the annotations indicate that hallucination is a general issue in all generated summaries. To evaluate the impact of $H_{pmi}$ on hallucination, we only consider those annotations where at least two annotators agree on the factuality label. The results show that $H_{min}$ and $H_{th}$ obtain a factuality rate of 36% and 33.3%, respectively, while $H_{pmi}$ achieves a slightly higher factuality rate of 36.2%. Given the small size of the evaluation sample, we con-

clude that there is no significant improvement of factuality with the $H_{pmi}$ approach on this task.

Since $H_{pmi}$ is an enhanced version of $H_{th}$, and $H_{th}$ is consistently preferred over $H_{pmi}$ (Figure 2), we delve deeper into cases where $H_{pmi}$ shows an improvement in factuality, yet it receives a lower rating than $H_{th}$. The results indicate that $H_{pmi}$ indeed redirects the prediction of the rest of the summary, specially when applied early on as stated in Li et al. (2016). However, it does not necessarily address the issue of hallucination. For example, the first text span of $H_{th}$ in Table 4 hallucinates the moment when the speech occurs *Nacht auf Mittwoch* ('Wednesday night'), and it is therefore annotated as not factual. In contrast, the $H_{pmi}$ generates a sentence about the Covid rules. However, none of the source articles refer to this topic,[18] which results in a more severe hallucination.

## 8 Conclusion

In this work, we tackle Multi-document Summarization (MDS) in low-resource settings where there is a lack of MDS training data. We therefore present various sampling approaches built upon prior works that use single-document summarization models for the MDS task. Specifically, we leverage information entropy as a metric to measure the model certainty in each token prediction. The experimental results on German MDS show that

---

[17] Although the guidelines mainly refer to long summaries of at least 150 words, we found them also useful in our setting.

[18] Source articles ids: 18163721, 18160037, and 18160205.

our $H_{th}$ approach, which specifically applies maximum probabilitic information entropy, achieves the state-of-art in German abstractive MDS. In our experiments, we also assessed an extended version of the $H_{th}$ approach that applies Pointwise Mutual Information (PMI) when all predictions exhibit uncertainty. Although PMI has been used in prior work to address hallucination, we observe in the manual evaluation that PMI changes the prediction of the rest of summary, but it does not inherently tackle hallucination. Future work should focus on addressing the issue of hallucination in automatic summarization, including further research on the efficacy of PMI to mitigate hallucinations. Additionally, it would be interesting to explore alternative approaches to enhance the $H_{th}$ approach when there is uncertainty in the prediction. Finally, we built a MDS test set of German news articles that will help the research community to evaluate abstractive MDS on German text.

## Ethics Statement

**Human Annotation**    We recruited the annotators for the manual evaluation task on a voluntary basis and provided them with information about the goals and scope of the task. The data was collected anonymously such that no conclusion can be drawn about any particular annotator. This human evaluation obtained the corresponding ethical approval from the Ethics Commission of ETH Zurich university (EK-2023-N-37).

**Text Generation Models**    Ethical considerations documented for natural language generation systems (Smiley et al., 2017; Kreps et al., 2022) also apply to our work. We do not anticipate any additional concerns.

*Supplementary Materials Availability Statement:* Source code for the presented entropy-based sampling approaches[19] in Section 4 and the Multi-GeNews dataset[20] described in Section 5.2 are available from GitHub.

## Acknowledgements

---

[19]Link to GitHub repository.
[20]Link to Multi-GeNews repository.

## References

Moye Chen, Wei Li, Jiachen Liu, Xinyan Xiao, Hua Wu, and Haifeng Wang. 2021. SgSum:transforming multi-document summarization into sub-graph selection. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4063–4074, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Hady Elsahar, Maximin Coavoux, Jos Rozen, and Matthias Gallé. 2021. Self-supervised and controlled multi-document opinion summarization. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1646–1662, Online. Association for Computational Linguistics.

Liana Ermakova, Jean Valère Cossu, and Josiane Mothe. 2019. A survey on evaluation of summarization methods. *Information processing & management*, 56(5):1794–1814.

Alexander Fabbri, Irene Li, Tianwei She, Suyi Li, and Dragomir Radev. 2019. Multi-news: A large-scale multi-document summarization dataset and abstractive hierarchical model. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1074–1084, Florence, Italy. Association for Computational Linguistics.

Angela Fan, Claire Gardent, Chloé Braud, and Antoine Bordes. 2019. Using local knowledge graph construction to scale Seq2Seq models to multi-document inputs. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4186–4196, Hong Kong, China. Association for Computational Linguistics.

J.L. Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological Bulletin*, 76(5):378–382.

Dominik Frefel. 2020. Summarization corpora of Wikipedia articles. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 6651–6655, Marseille, France. European Language Resources Association.

Tanya Goyal, Junyi Jessy Li, and Greg Durrett. 2022. News summarization and evaluation in the era of gpt-3. *arXiv preprint arXiv:2209.12356*.

Max Grusky, Mor Naaman, and Yoav Artzi. 2018. Newsroom: A dataset of 1.3 million summaries with diverse extractive strategies. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 708–719, New Orleans, Louisiana. Association for Computational Linguistics.

Chris Hokamp, Demian Gholipour Ghalandari, Nghia The Pham, and John Glover. 2020. DynE: Dynamic ensemble decoding for multi-document summarization. *arXiv preprint arXiv:2006.08748.*

Timo Johner, Abhik Jana, and Chris Biemann. 2021. Error analysis of using BART for multi-document summarization: A study for English and German language. In *Proceedings of the 23rd Nordic Conference on Computational Linguistics (NoDaLiDa)*, pages 391–397, Reykjavik, Iceland (Online). Linköping University Electronic Press, Sweden.

Sarah Kreps, R. Miles McCain, and Miles Brundage. 2022. All the news that's fit to fabricate: AI-generated text as a tool of media misinformation. *Journal of Experimental Political Science*, 9(1):104–117.

Kalpesh Krishna, Erin Bransom, Bailey Kuehl, Mohit Iyyer, Pradeep Dasigi, Arman Cohan, and Kyle Lo. 2023. LongEval: Guidelines for human evaluation of faithfulness in long-form summarization. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1650–1669, Dubrovnik, Croatia. Association for Computational Linguistics.

Logan Lebanoff, Kaiqiang Song, and Fei Liu. 2018. Adapting the neural encoder-decoder framework from single to multi-document summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4131–4141, Brussels, Belgium. Association for Computational Linguistics.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. A diversity-promoting objective function for neural conversation models. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 110–119, San Diego, California. Association for Computational Linguistics.

Wei Li, Xinyan Xiao, Jiachen Liu, Hua Wu, Haifeng Wang, and Junping Du. 2020. Leveraging graph to improve abstractive multi-document summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6232–6243, Online. Association for Computational Linguistics.

Yang Li, Jiachen Yang, and Jiabao Wen. 2021. Entropy-based redundancy analysis and information screening. *Digital Communications and Networks*.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Peter J. Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam Shazeer. 2018. Generating wikipedia by summarizing long sequences. In *International Conference on Learning Representations*.

Yang Liu and Mirella Lapata. 2019a. Hierarchical transformers for multi-document summarization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5070–5081, Florence, Italy. Association for Computational Linguistics.

Yang Liu and Mirella Lapata. 2019b. Text summarization with pretrained encoders. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3730–3740, Hong Kong, China. Association for Computational Linguistics.

Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pre-training for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:726–742.

Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In *International Conference on Learning Representations*.

Laura Perez-Beltrachini and Mirella Lapata. 2021. Multi-document summarization with determinantal point process attention. *Journal of Artificial Intelligence Research*, 71:371–399.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551.

Annette Rios, Nicolas Spring, Tannon Kew, Marek Kostrzewa, Andreas Säuberli, Mathias Müller, and Sarah Ebling. 2021. A new dataset and efficient baselines for document-level text simplification in German. In *Proceedings of the Third Workshop on New Frontiers in Summarization*, pages 152–161, Online and in Dominican Republic. Association for Computational Linguistics.

Charese Smiley, Frank Schilder, Vassilis Plachouras, and Jochen L. Leidner. 2017. Say the right thing right: Ethics issues in natural language generation

systems. In *Proceedings of the First ACL Workshop on Ethics in Natural Language Processing*, pages 103–108, Valencia, Spain. Association for Computational Linguistics.

Yun-Zhu Song, Yi-Syuan Chen, and Hong-Han Shuai. 2022. Improving multi-document summarization through referenced flexible extraction with credit-awareness. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1667–1681, Seattle, United States. Association for Computational Linguistics.

Junya Takayama and Yuki Arase. 2019. Relevant and informative response generation using pointwise mutual information. In *Proceedings of the First Workshop on NLP for Conversational AI*, pages 133–138, Florence, Italy. Association for Computational Linguistics.

Liam van der Poel, Ryan Cotterell, and Clara Meister. 2022. Mutual information alleviates hallucinations in abstractive summarization. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5956–5965, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Wen Xiao, Iz Beltagy, Giuseppe Carenini, and Arman Cohan. 2022. PRIMERA: Pyramid-based masked sentence pre-training for multi-document summarization. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5245–5263, Dublin, Ireland. Association for Computational Linguistics.

Jiacheng Xu, Shrey Desai, and Greg Durrett. 2020. Understanding neural abstractive summarization models via uncertainty. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6275–6281, Online. Association for Computational Linguistics.

Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter Liu. 2020. PEGASUS: Pre-training with extracted gap-sentences for abstractive summarization. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pages 11328–11339. Proceedings of Machine Learning Research (PMLR).

Fangwei Zhu, Shangqing Tu, Jiaxin Shi, Juanzi Li, Lei Hou, and Tong Cui. 2021. TWAG: A topic-guided Wikipedia abstract generator. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4623–4635, Online. Association for Computational Linguistics.

Markus Zopf. 2018. Auto-hMDS: Automatic construction of a large heterogeneous multilingual multi-document summarization corpus. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).

# Claim Optimization in Computational Argumentation

**Gabriella Skitalinskaya[1,2], Maximilian Spliethöver[1], and Henning Wachsmuth[1]**

[1]Leibniz University Hannover, Institute of Artificial Intelligence

[2]University of Bremen, Department of Computer Science

{g.skitalinska,m.spliethoever,h.wachsmuth}@ai.uni-hannover.de

## Abstract

An optimal delivery of arguments is key to persuasion in any debate, both for humans and for AI systems. This requires the use of clear and fluent claims relevant to the given debate. Prior work has studied the automatic assessment of argument quality extensively. Yet, no approach actually improves the quality so far. To fill this gap, this paper proposes the task of *claim optimization*: to rewrite argumentative claims in order to optimize their delivery. As multiple types of optimization are possible, we approach this task by first generating a diverse set of candidate claims using a large language model, such as BART, taking into account contextual information. Then, the best candidate is selected using various quality metrics. In automatic and human evaluation on an English-language corpus, our quality-based candidate selection outperforms several baselines, improving 60% of all claims (worsening 16% only). Follow-up analyses reveal that, beyond copy editing, our approach often specifies claims with details, whereas it adds less evidence than humans do. Moreover, its capabilities generalize well to other domains, such as instructional texts.

## 1 Introduction

The delivery of arguments in clear and appropriate language is a decisive factor in achieving persuasion in any debating situation, known as *elocutio* in Aristotle's rhetoric (El Baff et al., 2019). Accordingly, the claims composed in an argument should not only be grammatically fluent and relevant to the given debate topic, but also unambiguous, self-contained, and more. Written arguments therefore often undergo multiple revisions in which various aspects are optimized (Zhang and Litman, 2015).

Extensive research has been done on the automatic assessment of argument quality and the use of large language models on various text editing tasks. Yet, no work so far has studied how to ac-



Figure 1: Examples of different optimized versions of an *original claim* found on the debate platform Kialo. All optimizations were generated by the approach proposed in this paper, using the *debate topic* as context.

tually improve argumentative texts. However, developing respective approaches is a critical step towards building effective writing assistants, which could help learners write better argumentative texts (Wambsganss et al., 2021) or rephrase arguments made by an AI debater (Slonim et al., 2021). In this work, we close the outlined gap by studying how to employ language models for rewriting argumentative text to optimize its delivery.

We start by defining the task of *claim optimization* in Section 3, and adjust the English-language claim revision dataset of Skitalinskaya et al. (2021) for evaluation. The new task requires complementary abilities: On the one hand, different types of quality issues inside a claim must be detected, from grammatical errors to missing details. If not all quality aspects can be improved simultaneously, specific ones must be targeted. On the other hand, improved claim parts need to be integrated with the context of the surrounding discussion, while preserving the original meaning as far as possible. Figure 1 shows three exemplary optimizations of a claim from the debate platform *Kialo*. The first elaborates what the consequence of weaponization

is, whereas the second rephrases the claim to clarify what weaponizing means, employing knowledge about the debate topic. The third renders the stance of the claim explicit. We observe that different ways to optimize a claim exist, yet the level of improvement differs as well.

To account for the multiplicity of claim optimization, we propose a controlled generation approach that combines the capabilities of large language models with quality assessment (Section 4). First, a fine-tuned generation model produces several candidate optimizations of a given claim. To optimize claims, we condition the model on discourse context, namely the debate topic and the previous claim in the debate. The key to selecting the best optimization is to then score candidates using three quality metrics: *grammatical fluency*, *meaning preservation*, and *argument quality*. Such candidate selection remains understudied in many generative tasks, particularly within computational argumentation.

In automatic and manual evaluation (Section 5), we demonstrate the effectiveness of our approach, employing fine-tuned BART (Lewis et al., 2020) for candidate generation. Our results stress the benefits of quality assessment (Section 6). Incorporating context turns out especially helpful for making shorter claims—where the topic of the debate is difficult to infer—more self-contained. According to human annotators, our approach improves 60% of all claims and harms only 16%, clearly outperforming standard fine-tuned generation.

To gain further insights, we carry out a manual annotation of 600 claim optimizations and identify eight types typically found in online debate communities, such as *elaboration* and *disambiguation* (Section 7). Intriguingly, our approach covers similar optimization types as in human revisions, but we also observe limitations (Section 7). To explore to what extent it generalizes to other revision domains, we also carry out experiments on instructional texts (Anthonio and Roth, 2020) and formal texts (Du et al., 2022), finding that it outperforms strong baselines and state-of-the-art approaches.

In summary, the contributions of this paper are:

1. *a new task*, claim optimization, along with a manual analysis of typical optimization types;

2. *a computational approach* that selects the best generated candidate claim in terms of quality;

3. *empirical insights* into the impact and challenges of optimizing claims computationally.[1]

## 2 Related Work

Quality assessment has become a key topic in computational argumentation research (Lapesa et al., 2023). Various quality dimensions exist in argumentation theory, as surveyed by Wachsmuth et al. (2017) and assessed computationally in various works (Lauscher et al., 2020; Marro et al., 2022). Many of them relate to quality aspects we consider in this work, from clarity and organization (Wachsmuth et al., 2016) to the general evaluability of arguments (Park and Cardie, 2018), potential fallacies in their reasoning (Goffredo et al., 2022), and the appropriateness of the language used (Ziegenbein et al., 2023). Recently, (Skitalinskaya and Wachsmuth, 2023) tackled the question whether an argumentative claim is in need of revision, whereas Jundi et al. (2023) investigated where to best elaborate a discussion. While Gurcke et al. (2021) leverage claim generation for a refined assessment of argument quality, we are not aware of any prior work that actually optimizes arguments or their components in order to improve quality.

As shown in Figure 1, there can be several ways to optimize a given text. Our key idea is to select the best optimization among diverse candidates generated by a language model. Prior generation work on candidate selection hints at the potential benefits of such setup, albeit in other tasks and domains. In early work on rule-based conversational systems, Walker et al. (2001) introduced dialogue quality metrics to optimize template-based systems towards user satisfaction. Kondadadi et al. (2013) and Cao et al. (2018) chose the best templates for generation, and Mizumoto and Matsumoto (2016) used syntactic features to rank candidates in grammar correction. Recently, Yoshimura et al. (2020) proposed a reference-less metric trained on manual evaluations of grammar correction system outputs to assess generated candidates, while Suzgun et al. (2022) utilize pre-trained language models to select the best candidate in textual style transfer tasks.

In generation research on computational argumentation, candidate selection remains largely understudied. Most relevant in this regard is the approach of Chakrabarty et al. (2021) which reframes arguments to be more trustworthy (e.g., less partisan). It generates multiple candidates and selects

---

[1]Data, code, and models from our experiments are found at https://github.com/GabriellaSky/claim_optimization

one based on the entailment relation scores to the input. Extending this idea, we select candidates based on various properties, including argument quality.

Understanding the editing process of arguments is crucial, as it reveals what quality dimensions are considered important. For Wikipedia, Daxenberger and Gurevych (2013) proposed a fine-grained taxonomy as a result of their multi-label edit categorization of revisions (Daxenberger and Gurevych, 2012). The taxonomy focuses solely on the editing actions performed, such as inserting, deleting, and paraphrasing. In contrast, Yang et al. (2017) identified various semantic intentions behind Wikipedia revisions, from *copy editing* to *content clarifications* and *fact updates*. Their taxonomy defines a starting point for our research. Not all covered intentions generalize beyond Wiki scenarios, though.

Wikipedia-based corpora have often been used in the study of editing and rewriting, including paraphrasing (Max and Wisniewski, 2010), grammar correction (Lichtarge et al., 2019), bias neutralization (Pryzant et al., 2020), and controllable text editing (Faltings et al., 2021; Du et al., 2022). Similarly, WikiHow enabled summarization (Koupaee and Wang, 2018) and knowledge acquisition (Zhou et al., 2019). However, neither of these includes *argumentative* texts. Instead, we thus rely on the corpus of Skitalinskaya et al. (2021), which consists of revision histories of argumentative claims from online debates. Whereas the authors *compare* claims in terms of quality, we propose and study the new task of automatically *optimizing* claim quality. Moreoever, we see the revision types they distinguish (clarification, grammar correction, linking to external sources) as too coarse-grained to represent the diversity of claim optimizations. We refine them manually into eight optimization types, allowing for a more systematic analysis. Skitalinskaya et al. (2021) also found low correlations between the revision types and 15 common argument quality dimensions (Wachsmuth et al., 2017), suggesting that they are rather complementary. Primarily, they target the general form a well-phrased claim should have and its relevance to the debate.

For the analysis of argumentative text rewriting, Zhang and Litman (2015) incorporated both argumentative writing features and surface changes. To explore the classification of essay revisions, they defined a two-dimensional schema, combining the revision operation (e.g., modify, add, or delete) with the component being revised (e.g., reasoning or evidence). Moreover, Afrin and Litman (2018) created a small corpus of between-draft revisions of 60 student essays to study whether revision improves quality. However, these works do not uncover the reasoning behind a revision operation and are more geared towards analysis at the essay level.

## 3 Task and Data

This section introduces the proposed task and presents the data used for development and evaluation.

### 3.1 Claim Optimization

We define the claim optimization task as follows:

**Task** Given as input an argumentative claim $c$, potentially along with context information on the debate, rewrite $c$ into an output claim $\tilde{c}$ such that

(a) $\tilde{c}$ improves upon $c$ in terms of text quality and/or argument quality, and

(b) $\tilde{c}$ preserves the meaning of $c$ as far as possible.

While we conceptually assume that $c$ consists of one or more sentences and has at least one quality flaw, our approaches do not model this explicitly. Moreover, note that $c$ might have multiple flaws, resulting in $n \geq 2$ candidate optimizations $\tilde{C} = \{\tilde{c}_1, \ldots, \tilde{c}_n\}$. In this case, the goal is to identify the candidate $c^* \in \tilde{C}$ that maximizes overall quality.

### 3.2 Data for Development and Evaluation

We start from the ClaimRev dataset (Skitalinskaya et al., 2021), consisting of 124,312 claim revision histories from the debate platform *Kialo*. Each history defines a chain $(c_1, ..., c_m)$, in which claim $c_i$ is a revised version of the previous claim, $c_{i-1}$ with $1 < i \leq m$, improving upon its quality. According to the authors, this holds in 93% of all cases.

From each revision chain, we derived all possible optimization pairs $(c, \tilde{c}) := (c_{i-1}, c_i)$, in total 210,222. Most revisions are labeled with their intention by the users who performed them, rendering them suitable for learning to optimize claims automatically.[2] Overall, 95% of all pairs refer to three intention labels: *clarification*, *typo/grammar correction*, and *corrected/added links*. To avoid noise from the few remaining labels, we condensed the data to 198,089 instances of the three main labels.[3]

---

[2] As 26% of all pairs were unlabeled, we trained a BERT model to assign such pairs one of the 6 most prominent labels.

[3] The labels of the removed instances denote changes to the meaning of $c$ and statements from which no action or intention can be derived (e.g., "see comments", "moved as pro").

For the final task dataset, we associated each remaining pair $(c, \tilde{c})$ to its context: the *debate topic* $\tau$ (i.e., the thesis on Kialo) as well as the *previous claim* $\hat{c}$ (the parent on Kialo), which is supported or opposed by $c$ (see Figure 1). We sampled 600 revision pairs pseudo-randomly as a test set (200 per intention label), and split remaining pairs into training (90%) and validation set (10%). As the given labels are rather coarse-grained, we look into the optimizations in more detail in Section 7.

## 4 Approach

We now present the first approach to automatic claim optimization. To account for the variety of possible optimizations, multiple candidate claims are generated that are pertinent to the context given and preserve the claim's meaning. Then, the best candidate is selected based on quality metrics. Both steps are detailed below and illustrated in Figure 2.

### 4.1 Seq2Seq-based Candidate Generation

To generate candidates, we fine-tune a Seq2Seq model on pairs $(c, \tilde{c})$, by treating the original claim $c$ as encoder source and revised claim $\tilde{c}$ as the decoder target. In a separate experiment, we condition the model on context information, the debate topic $\tau$ and the previous claim $\hat{c}$, during fine-tuning to further optimize the relevance of generated candidates. The context is separated from $c$ by delimiter tokens (Keskar et al., 2019; Schiller et al., 2021).

Multiple ways to improve $c$ exist, especially if it suffers from multiple flaws, since not all flaws may be fixed in a single revision. Therefore, we first generate $n$ suitable candidates, $\tilde{c}_1, \ldots, \tilde{c}_n$, among which the best one is to be found later ($n$ is set to 10 in Section 5). However, the top candidates created by language models often tend to be very similar. To increase the diversity of candidates, we perform top-$k$ sampling (Fan et al., 2018), where we first generate the most probable claim (top-1) and then vary $k$ with in steps of 5 (e.g. top-5, top-10, etc).

### 4.2 Quality-based Candidate Selection

Among the $n$ candidates, we aim to find the optimal claim, $c^*$, that most improves the delivery of $c$ in terms of text and argument quality. Similar to Yoshimura et al. (2020), we tackle this task as a candidate selection problem. In our proposed strategy, *AutoScore*, we integrate three metrics: (1) grammatical fluency, (2) meaning preservation, and (3) argument quality. This way, we can *explicitly* favor



Figure 2: Proposed claim optimization approach: First, we generate $n$ candidates from the *original claim*, possibly conditioned on context information. Then, the *optimized claim* is selected using three quality metrics.

specific quality dimensions via respective models:

**Grammatical Fluency** We learn to assess fluency on the MSR corpus (Toutanova et al., 2016) where the grammaticality of abstractive compressions is scored by 3–5 annotators from 1 (disfluent) to 3 (fluent). We chose this corpus, since multiple compressions per input make a trained model sensitive to the differences in variants of a text. For training, we average all annotator scores and make the task binary, namely, a text is seen as disfluent unless all annotators gave score 3. Then, we train BERT on the data to obtain fluency probabilities (details found in Appendix A). The accuracy of our model on the suggested data split is 77.4.

**Meaning Preservation** To quantify to what extent a generated candidate maintains the meaning of the original claim, we compute their semantic similarity as the cosine similarity of the SBERT sentence embeddings (Reimers and Gurevych, 2019).

**Argument Quality** Finally, to examine whether the generated candidates are better than the original claim from an argumentation perspective, we fine-tune a BERT model on the task of pairwise

137

argument classification using the ClaimRev dataset. Since this corpus is also used to fine-tune the Seq2Seq model, we apply the same training and validation split as described in Section 3.2 to avoid data leakage, and obtain 75.5 accuracy. We then use its probability scores to determine relative quality improvement (for more details see Appendix A).

Given the three quality metrics, we calculate the final evaluation score, $AutoScore$, as the weighted linear sum of all three individual scores as

$$\alpha \cdot fluency + \beta \cdot meaning + \gamma \cdot argument,$$

where $fluency$, $meaning$, and $argument$ are normalized scores of the three outlined quality metrics. The non-negative weights satisfy $\alpha + \beta + \gamma = 1$.

It should be noted that depending on the domain or writing skills of the users, there may be other more suitable datasets or approaches to capturing the outlined quality aspects, which could potentially lead to further performance improvements. While we do explore how well the suggested approaches transfer to certain other domains of text (see Section 7.3), identifying the optimal model for each quality dimension falls beyond the scope of this paper.

## 5 Experiments

This section describes our experimental setup to study how well the claims from Section 3 can be improved using our approach from Section 4. We focus on the impact of candidate selection.

### 5.1 Seq2Seq-based Candidate Generation

For candidate generation, we employ the pre-trained conditional language model BART (Lewis et al., 2020), using the *bart-large* checkpoint. However, other Seq2Seq architectures can also be considered within our approach (see Appendices A, B).

### 5.2 Quality-based Candidate Selection

We evaluate our candidate selection approach in comparison to three ablations and four baselines:

**Approach** To utilize AutoScore for choosing candidates, the optimal weighting of its metrics must be determined. We follow Yoshimura et al. (2020), performing a grid search in increments of 0.01 in the range of 0.01 to 0.98 for each weight to maximize the Pearson's correlation coefficient between AutoScore and the original order of the revisions

from revision histories in the validation set. Similar has been done for counterargument retrieval by Wachsmuth et al. (2018). The best weights found are $\alpha = 0.43$, $\beta = 0.01$, and $\gamma = 0.56$, suggesting that meaning preservation is of low importance and potentially may be omitted. We suppose this is due to the general similarity of the generated candidates, so a strong meaning deviation is unlikely.

**Ablations** To assess the impact of each considered quality metric used in AutoScore, we perform an ablation study, where optimal candidates are chosen based on the individual metric scores:

- *Max Fluency.* Highest grammatical fluency
- *Max Argument.* Highest argument quality
- *Max Meaning.* Highest semantic similarity

**Baselines** We test four selection strategies for 10 candidates generated via top-$k$ sampling:

- *Unedited.* Return the original input as output.
- *Top-1.* Return the most likely candidate (obtained by appending the most probable token generated by the model at each time step).
- *Random.* Return candidate pseudo-randomly.
- *SVMRank.* Rerank candidates with SVMRank (Joachims, 2006). Using sentence embeddings we decide which of the claim versions is better, by fine-tuning SBERT (*bert-base-cased*) on the corpus of Skitalinskaya et al. (2021).

### 5.3 Evaluation

We explore claim optimization on all 600 test cases, both automatically and manually:

**Automatic Evaluation** We compare all content selection strategies against the reference revisions using the precision-oriented *BLEU* (Papineni et al., 2002), recall-oriented *Rouge-L* (Lin, 2004), *SARI* (Xu et al., 2016), which computes the average $F_1$-scores of the added, kept, and deleted $n$-grams in comparison to the ground truth revision output, and the *exact match accuracy*. We also compute the semantic similarity of the optimized claim and the context information to capture whether conditioning claims on context affects their topic relevance.

**Manual Evaluation** As we fine-tune existing generation models rather than proposing new ones, we focus on the *candidate selection* in two manual annotation studies. For each instance, we acquired five independent crowdworkers via *MTurk*.

138

In the first study, the annotators scored all candidates with respect to the three considered quality metrics. We used the following Likert scales:

- *Fluency.* 1 (major errors, disfluent), 2 (minor errors), and 3 (fluent)

- *Meaning Preservation.* 1 (entirely different), 2 (substantial differences), 3 (moderate differences), 4 (minor differences), and 5 (identical)

- *Argument Quality.* 1 (notably worse than original), 2 (slightly worse), 3 (same as original), 4 (slightly improved), and 5 (notably improved)

A challenge of crowdsourcing is to ensure good results (Sabou et al., 2014). To account for this, we obtained the fina fluency, argument quality and meaning preservation scores using MACE (Hovy et al., 2013), a Bayesian model that gives more weight to reliable workers. In the given case, 39% of the 46 annotators had a MACE competence value $> 0.3$, which can be seen as reasonable in MTurk studies.

In the second study, we asked annotators to rank four candidates, returned by the content selection strategies, by perceived overall quality. If multiple candidates were identical, we showed each only once. While Krippendorff's $\alpha$ agreement was only 0.20 and percent agreement was 0.36% (majority voting), such values are common in subjective tasks (Wachsmuth et al., 2017; Alshomary et al., 2021).

## 6 Results and Discussion

Apart from evaluating the applicability of large generative language models to the task of argumentative claim optimization in general, our experiments focus on two questions: (1) Does the use of explicit knowledge about text and argument quality lead to the selection of better candidates? (2) Does the use of contextual information make the generated candidates more accurate and relevant to the debate?

### 6.1 Overall Claim Optimization Performance

**Automatic Evaluation**  Table 1 shows the automatic scores of all considered candidate selection strategies. The high scores of the baseline *Unedited* on metrics such as BLEU and ROUGE-L indicate that many claim revisions change little only. In contrast, *Unedited* is worst on SARI, a measure taking into account words that are added, deleted, and kept in changes, making it more suitable for evaluation. Here, *BART+AutoScore* performs best on SARI (43.7) and exact match accuracy (8.3%).

| Approach | BLEU | RouL | SARI | NoEd↓ | ExM |
|---|---|---|---|---|---|
| **Baselines** | | | | | |
| Unedited | **69.4** | **0.87** | 27.9 | 1.00 | 0.0% |
| BART + Top-1 | 64.0 | 0.83 | 39.7 | 0.31 | 7.8% |
| BART + Random | 62.6 | 0.83 | 38.7 | 0.28 | 6.8% |
| BART + SVMRank | 55.7 | 0.76 | 38.8 | 0.03 | 4.5% |
| **Approach** | | | | | |
| BART + AutoScore | 59.4 | 0.80 | **43.7** | **0.02** | **8.3%** |
| **Ablation** | | | | | |
| BART + Max Fluency | 57.6 | 0.78 | 41.5 | 0.09 | 5.8% |
| BART + Max Argument | 60.9 | 0.81 | 43.6 | 0.02 | 8.0% |
| BART + Max Meaning | 69.0 | 0.87 | 33.8 | 0.72 | 5.2% |

Table 1: Automatic evaluation: Performance of each candidate selection strategy on 600 test cases in terms of BLEU, Rouge-L, SARI, ratio of unedited cases, and ratio of exact matches to target reference.

| Model | Strategy | Fluency | Argument | Meaning | Rank |
|---|---|---|---|---|---|
| BART | Top-1 | 2.29 | 3.61 | 3.65 | 2.16 |
| | Random | 2.26 | 3.50 | 3.53 | 2.06 |
| | SVMRank | **2.33** | **3.69** | **3.66** | 1.95 |
| | **AutoScore** | **2.33** | 3.61 | 3.57 | **1.92** |

Table 2: Manual evaluation: Scores on the 600 test cases generated by BART using our candidate selection strategy *AutoScore* or the baselines: fluency (1–3), argument quality and meaning (1–5), mean rank (1–4, lower better). AutoScore ranks significantly better than *Top-1* ($p < .005$), *Random* ($p < .05$), and *SVMRank* ($p < .1$).

The *BART+Max Meaning* ablation supports the intuition that the candidates with highest meaning preservation scores are those with minimal changes, if any (72% of the candidates remain identical to the input). Such identical outputs are undesirable, as the claims are not optimized successfully, which is also corroborated by the low weight parameter ($\beta = 0.01$) found for the meaning preservation metric when optimizing AutoScore (see Section 5).

**Manual Evaluation**  Table 2 shows that human annotators prefer optimized candidates selected by *AutoScore*, with an average rank of 1.92. The difference to *Top-1* and *Random* is statistically significant ($p < .05$ in both cases) according to a Wilcoxon signed-rank test, whereas the gain over the second-best algorithm, *SVMRank*, is limited. Also, candidates of AutoScore and SVMRank are deemed more fluent than those of Top-1 and Random (2.33 vs. 2.29 and 2.26). In terms of argument quality, the results deviate from the automatic evaluation (Table 1), showing marginally higher scores for SVMRank and Top-1. Further analysis revealed that AutoScore and SVMRank agreed on the optimal candidate in 35% of the cases, partially

| Context | BLEU | Original | Previous | Topic |
|---|---|---|---|---|
| Claim only | 59.4 | 0.95 | 0.55 | 0.55 |
| + Previous Claim | **60.3** | 0.95 | **0.57** | **0.57** |
| + Debate Topic | 60.0 | 0.95 | 0.55 | 0.55 |
| Human-Baseline | 100.0 | 0.94 | 0.55 | 0.55 |

Table 3: BLEU and semantic similarity score with respect to the *original* claim, the debate's *previous* claim, and its *topic* of BART+AutoScore, depending on the context given for the 600 test samples.

explaining their close scores. Although SVMRank achieved high scores across the three quality metrics, we note that the annotators preferred candidates scores generated by AutoScore, highlighting the importance of more diverse revision changes reflected by lower meaning preservation scores.

Overall, our findings suggest that using candidate selection approaches that incorporate quality assessments (i.e., AutoScore and SVMRank) leads to candidates of higher fluency and argument quality while preserving the meaning of the original claim. In addition to Figure 1, examples of automatically-generated optimized claims can be found in the appendix.

## 6.2 Performance with Context Integration

Table 3 shows the semantic similarity of claims optimized by our approach and context information, depending on the context given. The results reveal slight improvements when conditioning the model on the previous claim (e.g., 60.3 vs. 59.4 BLEU). To check whether this led to improved claims, two authors of the paper compared 600 claims generated with and without the use of the previous claim in terms of (a) which claim seems better overall and (b) which seems more grounded. We found that using the previous claim as context improved quality in 12% of the cases and lowered it in 1% only, while leading to more grounded claims in 36%.

**Qualitative Analysis** Our manual inspection of a claim sample revealed the following insights:

First, conditioning on context reduces the number of erroneous specifications, particularly for very short claims with up to 10 words. This seems intuitive, as such claims often convey little information about the topic of the debate, making inaccurate changes without additional context likely.

Next, Kialo revisions often adhere to the following form: A claim introduces a statement and/or supporting facts, followed by a conclusion. This

pattern was frequently mimicked by our approach. Yet, in some cases, it added a follow-up sentence repeating the original claim in different wording or generated conclusions containing fallacious or unsound phrases contradicting the original claim in others. Modeling context mitigated this issue.

Finally, we found that models conditioned on different contexts sometimes generated candidates optimized in different regards, whereas a truly optimal candidate would be a fusion of both suggestions.

## 7 Analysis

To explore the nature of claim optimization and the capabilities of our approach, this section reports on (a) what types of optimizations exist, (b) how well our approach can operationalize these, and (c) how well it generalizes to non-argumentative domains.

## 7.1 Taxonomy of Optimization Types

To understand the relationship between optimizations found in the data and the underlying revision intentions, two authors of this paper inspected 600 claim revisions of the test set. Opposed to actions, intentions describe the goal of an edit (e.g., making a text easier to read) rather than referring to specific changes(e.g., paraphrasing or adding punctuation). We build on ideas of Yang et al. (2017) who provide a taxonomy of revision intentions in Wikipedia texts. Claims usually do not come from encyclopedias, but from debate types or from monological arguments, as in essays (Persing and Ng, 2015). Therefore, we adapt the terminology of Yang et al. (2017) to gear it more towards argumentative texts.

As a result of a joint discussion of various sample pairs, we decided to distinguish eight optimization types, as presented in Table 4. Both authors then annotated all 600 test pairs for these types, which led to only 29 disagreement cases, meaning a high agreement of 0.89 in terms of Cohen's $\kappa$. These cases were resolved by both annotators together.[4]

Table 4 also shows cooccurrences of the types and intention labels. *Typo/grammar correction* and *correcting/adding links* align well with *copy editing* and *corroboration* respectively. In contrast, clarification is broken into more fine-grained types, where *specification* seems most common with 58

---

[4]We acknowledge that there is potential bias inherent in self-annotation. However, we would like to point out that no knowledge about the test set was used to develop the approach presented in Section 4.

| # | Optimization | Description of the Type | Clarification | Grammar | Links |
|---|---|---|---|---|---|
| 1 | Specification | Specifying or explaining a given fact or meaning (of the argument) by adding an example or discussion without adding new information. | 58 | 1 | – |
| 2 | Simplification | Removing information or simplifying the sentence structure, e.g., with the intent to reduce the complexity or breadth of the claim. | 43 | – | – |
| 3 | Reframing | Paraphrasing or rephrasing a claim, e.g., with the intent to specify or generalize the claim, or to add clarity. | 29 | – | – |
| 4 | Elaboration | Extending the claim by more information or adding a fact with the intent to make the claim more self-contained, sound, or stronger. | 23 | – | – |
| 5 | Corroboration | Adding, editing, or removing evidence in the form of links that provide supporting information or external resources to the claim. | 8 | – | 153 |
| 6 | Neutralization | Rewriting a claim using a more encyclopedic or neutral tone, e.g., with the intent to remove bias or biased language. | 7 | – | – |
| 7 | Disambiguation | Reducing ambiguity, e.g., replacing pronouns by concepts mentioned before in the debate, or replacing acronyms with what they stand for. | 7 | – | 1 |
| 8 | Copy editing | Improving the grammar, spelling, tone, or punctuation of a claim, without changing the main point or meaning. | 41 | 200 | 52 |

Table 4: Descriptions of the eight claim optimization types identified in the 600 test pairs. The right columns show the count of claims per type for each of the three intention labels from Skitalinskaya et al. (2021): *clarification*, typo/*grammar* correction, and correcting/adding *links*. Note, that a revision may be assigned to multiple categories.

| Type | Human | Approach | Better | Same | Worse |
|---|---|---|---|---|---|
| Specification | 59 | 152 | 65% | 19% | 16% |
| Simplification | 43 | 18 | 61% | 28% | 11% |
| Reframing | 29 | 21 | 62% | 33% | 5% |
| Elaboration | 23 | 55 | 62% | 18% | 20% |
| Corroboration | 161 | 38 | 53% | 23% | 24% |
| Neutralization | 7 | 0 | – | – | – |
| Disambiguation | 8 | 8 | 63% | 25% | 12% |
| Copy editing | 293 | 301 | 59% | 26% | 15% |
| **Overall** | 623 | 593 | 60% | 24% | 16% |

Table 5: Manual analysis: Comparison of the human-optimized claims of all 600 test cases (some have multiple) and of the claims optimized by BART+AutoScore (15 claims were unchanged). The three right columns show the ratio of optimized claims judged *better*, *same*, or *worse* than the original in terms of overall quality.

cases, followed by *simplification* and *reframing*. Examples of each type are found in the appendix.

We point out that the eight types are not exhaustive for all possible claim quality optimizations, but rather provide insights into the semantic and discourse-related phenomena observed in the data. We see them as complementary to the argument quality taxonomy of Wachsmuth et al. (2017) as ways to improve the delivery-related quality dimensions: *clarity*, *appropriateness*, and *arrangement*.

## 7.2 Performance across Optimization Types

To enable comparison between the human optimizations and automatically generated outputs, two authors of the paper labeled 600 optimized claims with the types defined in Table 4. Due to resource constrains only the best performing ap-

proach, BART+AutoScore, was considered. Overall, our approach generates better claims in 60% of the cases, while 84% remain at least of similar quality.

Most noteworthily, we observe that our approach performs optimizations of the type *specification* 2.5 times as often as humans, and more than double as many *elaboration* revisions (55 vs. 23). In contrast, it adds, edits, or removes evidence in the form of links (*corroboration*) four times less often than humans. The model also made fewer *simplifications* (18 vs. 43) and no *neutralization* edits, which may be due to data imbalance regarding such types.

In terms of average quality, *specification* (65%) and *disambiguation* edits (63%) most often lead to improvements, but the eight types appear rather balanced in this regard. The Jaccard similarity score between optimizations performed by humans and our approach is 0.37, mostly agreeing on copy edits (178 cases) and corroboration (22 cases). Given such low overlap, future work should consider conditioning models to generate specific optimizations.

## 7.3 Performance across Revision Domains

Lastly, we examine whether our approach, along with the chosen text quality metrics, applies to texts from other domains. We consider two datasets: *WikiHow* (Anthonio and Roth, 2020), containing revisions of instructional texts, and *IteraTeR* (Du et al., 2022), containing revisions of various formal texts, such as encyclopedia entries, news, and scientific papers. For our experiments, we use the provided document-level splits, and sample 1000

| Approach | BLEU | RouL | SARI | NoEd↓ | ExM |
|---|---|---|---|---|---|
| **WikiHow Dataset** | | | | | |
| Unedited | 65.7 | 0.85 | 28.4 | 1.00 | 0.00% |
| BART + Top-1 | **64.7** | **0.83** | 41.3 | 0.50 | 13.0% |
| BART + AutoScore | 61.8 | 0.80 | **48.5** | **0.08** | **16.0%** |
| **IteraTeR Dataset** | | | | | |
| Unedited | 74.0 | 0.86 | 28.6 | 1.00 | 0.00% |
| BART + Top-1 | **68.9** | **0.83** | 37.0 | 0.07 | 0.00% |
| BART + AutoScore | 64.8 | 0.80 | **38.6** | **0.02** | 0.00% |

Table 6: Automatic evaluation: Performance of candidate selection strategies on data from other domains, in terms of BLEU, Rouge-L, SARI, ratio of unedited samples, and ratio of exact matches to target reference.

revision pairs pseudo-randomly as a final test set.

Table 6 shows automatic evaluation results. In both cases, *BART+Autoscore* leads to higher SARI scores (48.5 vs. 41.3 for WikiHow, 38.6 vs. 37.0 for IteraTeR), and notably reduces the number of cases where the models failed to revise the input (0.08 vs. 0.50 for WikiHow). The reported *BART+Top1* model represents the approach of Du et al. (2022), indicating that our approach and its text quality metrics achieve state-of-the-art performance with systematic improvements across domains, when generating optimized content. However, as different domains of text have different goals, different notions of quality, and, subsequently, different revision types performed, integrating domain-specific quality metrics may further improve performance. We leave this for future work.

## 8 Conclusion

With this paper, we work towards the next level of computational argument quality research, namely, to not only *assess* but also to *optimize* argumentative text. Applications include suggesting improvements in writing support and automatic phrasing in debating systems. We presented an approach that generates multiple candidate claim optimizations and then selects the best one using various quality metrics. In experiments, combining fine-tuned BART with such candidate selection improved 60% of the claims from online debates, outperforming several baseline models and candidate selection strategies. We showcased generalization capabilities on two out-of-domain datasets, but we also found some claim optimization types hard to automate.

In future work, we seek to examine whether recent large language models (e.g., Alpaca) and end-to-end models (where generation and candidate se-

lection are learned jointly) can further optimize the quality of claims. As our approach so far relies on the availability of large claim revision corpora and language models, techniques for low-resource scenarios and languages should be explored to make claim optimization more widely applicable.

## Acknowledgments

## Ethical Considerations

This work contributes to the task of argumentative text editing, namely we explore how to revise claims automatically in order to optimize their quality. While our work may also improve downstream task performance on other tasks, it is mainly intended to support humans in scenarios, such as the creation and moderation of content on online debate platforms as well as the improvement of arguments generated or retrieved by other systems. In particular, the presented approach is meant to help users by showing examples of how to further optimize their claims in relation to a certain debate topic, so they can deliver their messages effectively and hone their writing skills.

However, our generation approach still comes with limitations and may favor revision patterns over others in unpredictable ways, both of which might raise ethical concerns. For example, it may occasionally produce false claims based on untrue or non-existent facts. We think, humans should be able to identify such cases in light of the available context though, as long as the improvements remain suggestions and do not happen fully automatically, as intended.

The presented technology might further be subject to intentional misuse. A word processing software, for example, could be conditioned to automatically detect and adapt claims made by the user in subtle ways that favors political or social views of the software provider. Such misuse might then not only change the intended message of the text, but also influence or even change the views of the user (Jakesch et al., 2023).

In a different scenario, online services, such as social media platforms or review portals, might change posted claims (e.g. social media posts, online reviews) to personalize them and increase user

engagement or revenue. These changes might not only negatively affect the posting, but also the visiting user.

While it is hard to prevent such misuse, we think that the described scenarios are fairly unlikely, as such changes tend to be noticed by the online community quickly. Furthermore, the presented architecture and training procedure would require notable adaptations to produce such high-quality revisions.

An aspect that remains unexplored in this work is the ability of the presented approaches to work with variations of the English language, such as African-American English, mainly due to the lack of available data. In this regard, the approach might unfairly disadvantage or favor particular language varieties and dialects, potentially inducing social bias and harm if applied in public scenarios. We encourage researchers and practitioners to stay alert for such cases and to choose training data with care for various social groups.

Finally, our work included the labeling of generated candidate claims on a crowdsourcing platform. As detailed in Section 5, we compensated MTurk workers $13 per hour, complying with minimum wage standards in most countries at the time of conducting the experiment.

## References

Tazin Afrin and Diane Litman. 2018. Annotation and classification of sentence-level revision improvement. In *Proceedings of the Thirteenth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 240–246, New Orleans, Louisiana. Association for Computational Linguistics.

Milad Alshomary, Wei-Fan Chen, Timon Gurcke, and Henning Wachsmuth. 2021. Belief-based generation of argumentative claims. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 224–233, Online. Association for Computational Linguistics.

Talita Anthonio and Michael Roth. 2020. What can we learn from noun substitutions in revision histories? In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 1359–1370, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Ziqiang Cao, Wenjie Li, Sujian Li, and Furu Wei. 2018. Retrieve, rerank and rewrite: Soft template based neural summarization. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 152–161,

Melbourne, Australia. Association for Computational Linguistics.

Tuhin Chakrabarty, Christopher Hidey, and Smaranda Muresan. 2021. ENTRUST: Argument reframing with language models and entailment. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4958–4971, Online. Association for Computational Linguistics.

Johannes Daxenberger and Iryna Gurevych. 2012. A corpus-based study of edit categories in featured and non-featured Wikipedia articles. In *Proceedings of COLING 2012*, pages 711–726, Mumbai, India. The COLING 2012 Organizing Committee.

Johannes Daxenberger and Iryna Gurevych. 2013. Automatically classifying edit categories in Wikipedia revisions. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 578–589, Seattle, Washington, USA. Association for Computational Linguistics.

Wanyu Du, Vipul Raheja, Dhruv Kumar, Zae Myung Kim, Melissa Lopez, and Dongyeop Kang. 2022. Understanding iterative revision from human-written text. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3573–3590, Dublin, Ireland. Association for Computational Linguistics.

Roxanne El Baff, Henning Wachsmuth, Khalid Al Khatib, Manfred Stede, and Benno Stein. 2019. Computational argumentation synthesis as a language modeling task. In *Proceedings of the 12th International Conference on Natural Language Generation*, pages 54–64, Tokyo, Japan. Association for Computational Linguistics.

Felix Faltings, Michel Galley, Gerold Hintz, Chris Brockett, Chris Quirk, Jianfeng Gao, and Bill Dolan. 2021. Text editing by command. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5259–5274, Online. Association for Computational Linguistics.

Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 889–898, Melbourne, Australia. Association for Computational Linguistics.

Pierpaolo Goffredo, Shohreh Haddadan, Vorakit Vorakitphan, Elena Cabrio, and Serena Villata. 2022. Fallacious argument classification in political debates. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 4143–4149. International Joint Conferences on Artificial Intelligence Organization. Main Track.

Timon Gurcke, Milad Alshomary, and Henning Wachsmuth. 2021. Assessing the sufficiency of arguments through conclusion generation. In *Proceedings of the 8th Workshop on Argument Mining*, pages 67–77, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Dirk Hovy, Taylor Berg-Kirkpatrick, Ashish Vaswani, and Eduard Hovy. 2013. Learning whom to trust with MACE. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1120–1130, Atlanta, Georgia. Association for Computational Linguistics.

Maurice Jakesch, Advait Bhat, Daniel Buschek, Lior Zalmanson, and Mor Naaman. 2023. Co-Writing with Opinionated Language Models Affects Users' Views. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, CHI '23, pages 1–15, New York, NY, USA. Association for Computing Machinery.

Thorsten Joachims. 2006. Training linear svms in linear time. In *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '06, page 217–226, New York, NY, USA. Association for Computing Machinery.

Iman Jundi, Neele Falk, Eva Maria Vecchi, and Gabriella Lapesa. 2023. Node placement in argument maps: Modeling unidirectional relations in high & low-resource scenarios. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5854–5876, Toronto, Canada. Association for Computational Linguistics.

Nitish Shirish Keskar, Bryan McCann, Lav R. Varshney, Caiming Xiong, and Richard Socher. 2019. CTRL: A conditional transformer language model for controllable generation. *CoRR*, abs/1909.05858.

Ravi Kondadadi, Blake Howald, and Frank Schilder. 2013. A statistical NLG framework for aggregated planning and realization. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1406–1415, Sofia, Bulgaria. Association for Computational Linguistics.

Mahnaz Koupaee and William Yang Wang. 2018. Wikihow: A large scale text summarization dataset. *CoRR*, abs/1810.09305.

Gabriella Lapesa, Eva Maria Vecchi, Serena Villata, and Henning Wachsmuth. 2023. Mining, assessing, and improving arguments in NLP and the social sciences. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics: Tutorial Abstracts*, pages 1–6, Dubrovnik, Croatia. Association for Computational Linguistics.

Anne Lauscher, Lily Ng, Courtney Napoles, and Joel Tetreault. 2020. Rhetoric, logic, and dialectic: Advancing theory-based argument quality assessment in natural language processing. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4563–4574, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Jared Lichtarge, Chris Alberti, Shankar Kumar, Noam Shazeer, Niki Parmar, and Simon Tong. 2019. Corpora generation for grammatical error correction. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3291–3301, Minneapolis, Minnesota. Association for Computational Linguistics.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Santiago Marro, Elena Cabrio, and Serena Villata. 2022. Graph embeddings for argumentation quality assessment. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 4154–4164, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Aurélien Max and Guillaume Wisniewski. 2010. Mining naturally-occurring corrections and paraphrases from Wikipedia's revision history. In *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10)*, Valletta, Malta. European Language Resources Association (ELRA).

Tomoya Mizumoto and Yuji Matsumoto. 2016. Discriminative reranking for grammatical error correction with statistical machine translation. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1133–1138, San Diego, California. Association for Computational Linguistics.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

144

Joonsuk Park and Claire Cardie. 2018. A corpus of eRulemaking user comments for measuring evaluability of arguments. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).

Isaac Persing and Vincent Ng. 2015. Modeling argument strength in student essays. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 543–552, Beijing, China. Association for Computational Linguistics.

Reid Pryzant, Richard Diehl Martinez, Nathan Dass, Sadao Kurohashi, Dan Jurafsky, and Diyi Yang. 2020. Automatically neutralizing subjective bias in text. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(01):480–489.

Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing*, pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.

Marta Sabou, Kalina Bontcheva, Leon Derczynski, and Arno Scharl. 2014. Corpus annotation through crowdsourcing: Towards best practice guidelines. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 859–866, Reykjavik, Iceland. European Language Resources Association (ELRA).

Benjamin Schiller, Johannes Daxenberger, and Iryna Gurevych. 2021. Aspect-controlled neural argument generation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 380–396, Online. Association for Computational Linguistics.

Gabriella Skitalinskaya, Jonas Klaff, and Henning Wachsmuth. 2021. Learning from revisions: Quality assessment of claims in argumentation at scale. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1718–1729, Online. Association for Computational Linguistics.

Gabriella Skitalinskaya and Henning Wachsmuth. 2023. To revise or not to revise: Learning to detect improvable claims for argumentative writing support. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15799–15816, Toronto, Canada. Association for Computational Linguistics.

Noam Slonim, Yonatan Bilu, Carlos Alzate, Roy Bar-Haim, Ben Bogin, Francesca Bonin, Leshem Choshen, Edo Cohen-Karlik, Lena Dankin, Lilach

Edelstein, et al. 2021. An autonomous debating system. *Nature*, 591(7850):379–384.

Mirac Suzgun, Luke Melas-Kyriazi, and Dan Jurafsky. 2022. Prompt-and-rerank: A method for zero-shot and few-shot arbitrary textual style transfer with small language models.

Kristina Toutanova, Chris Brockett, Ke M. Tran, and Saleema Amershi. 2016. A dataset and evaluation metrics for abstractive compression of sentences and short paragraphs. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 340–350, Austin, Texas. Association for Computational Linguistics.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Henning Wachsmuth, Khalid Al Khatib, and Benno Stein. 2016. Using argument mining to assess the argumentation quality of essays. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 1680–1691. The COLING 2016 Organizing Committee.

Henning Wachsmuth, Nona Naderi, Yufang Hou, Yonatan Bilu, Vinodkumar Prabhakaran, Tim Alberdingk Thijm, Graeme Hirst, and Benno Stein. 2017. Computational argumentation quality assessment in natural language. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 176–187, Valencia, Spain. Association for Computational Linguistics.

Henning Wachsmuth, Shahbaz Syed, and Benno Stein. 2018. Retrieval of the best counterargument without prior topic knowledge. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 241–251. Association for Computational Linguistics.

Marilyn A. Walker, Rebecca Passonneau, and Julie E. Boland. 2001. Quantitative and qualitative evaluation of darpa communicator spoken dialogue systems. In *Proceedings of the 39th Annual Meeting of the Association for Computational Linguistics*, pages 515–522, Toulouse, France. Association for Computational Linguistics.

Thiemo Wambsganss, Christina Niklaus, Matthias Söllner, Siegfried Handschuh, and Jan Marco Leimeister. 2021. Supporting cognitive and emotional empathic writing of students. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4063–4077, Online. Association for Computational Linguistics.

Sam Wiseman and Alexander M. Rush. 2016. Sequence-to-sequence learning as beam-search optimization. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1296–1306, Austin, Texas. Association for Computational Linguistics.

Wei Xu, Courtney Napoles, Ellie Pavlick, Quanze Chen, and Chris Callison-Burch. 2016. Optimizing statistical machine translation for text simplification. *Transactions of the Association for Computational Linguistics*, 4:401–415.

Diyi Yang, Aaron Halfaker, Robert Kraut, and Eduard Hovy. 2017. Identifying Semantic Edit Intentions from Revisions in Wikipedia. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2000–2010, Copenhagen, Denmark. Association for Computational Linguistics.

Ryoma Yoshimura, Masahiro Kaneko, Tomoyuki Kajiwara, and Mamoru Komachi. 2020. SOME: Reference-less sub-metrics optimized for manual evaluations of grammatical error correction. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6516–6522, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Fan Zhang and Diane Litman. 2015. Annotation and classification of argumentative writing revisions. In *Proceedings of the Tenth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 133–143, Denver, Colorado. Association for Computational Linguistics.

Yilun Zhou, Julie Shah, and Steven Schockaert. 2019. Learning household task knowledge from WikiHow descriptions. In *Proceedings of the 5th Workshop on Semantic Deep Learning (SemDeep-5)*, pages 50–56, Macau, China. Association for Computational Linguistics.

Timon Ziegenbein, Shahbaz Syed, Felix Lange, Martin Potthast, and Henning Wachsmuth. 2023. Modeling appropriate language in argumentation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4344–4363, Toronto, Canada. Association for Computational Linguistics.

# A  Implementation and Training Details

## A.1  Candidate Generation Models

For generation, we use the pre-trained BART model implemented in the fairseq library. The library and pre-trained models are BSD-licensed. We use the BART-large checkpoint (400M parameters) and further finetune the model for 10 epochs on 2 RTX 2080Ti GPUs. We use the same parameters as suggested in the fine-tuning of BART for the CNN-DM summarization task by fairseq and set MAX-TOKENS to 1024. The training time is 100-140 minutes, depending on the chosen setup (with or without context information).

During inference, we generate candidates using a top-k random sampling scheme (Fan et al., 2018) with the following parameters: length penalty is set to 1.0, n-grams of size 3 can only be repeated once, temperature is set to 0.7, while the minimum and maximum length of the sequence to be generated are 7 and 256 accordingly.

## A.2  Quality Assessment Models

For the automatic assessment of fluency and argument quality, we use the bert-base-cased pre-trained BERT version, as implemented in the huggingface library. The library and pre-trained models have the Apache License 2.0. We finetune the model for two epochs and use the parameters suggested in Skitalinskaya et al. (2021). The accuracy of the trained model for fluency obtained on the train/dev/test split suggested by the authors (Toutanova et al., 2016) is 77.4 and 75.5 for argument quality.

For labeling the missing or unassigned revision types, we use the same bert-base-cased pre-trained BERT model, but in a multi-label setup, where we consider the following 6 classes: claim clarification, typo or grammar correction, correcting or adding links, changing the meaning of the claim, splitting the claim, and merging claims. We fine-tune the model for two epochs using the Adam optimizer with a learning rate of 1e-5 and achieve a weighted F1-score of 0.81.

# B  Alternative Generation Models

For comparison, we provide two additional baseline Seq2Seq model architectures, which help identify the complexity of the model needed for the task:

**LSTM.** Our first baseline is a popular LSTM variant introduced by Wiseman and Rush (2016).

| Model | Strategy | BLEU | RouL | SARI | NoEd↓ | ExM |
|---|---|---|---|---|---|---|
| BART | Top-1 | 64.0 | 0.83 | 39.7 | 0.31 | 7.8% |
| | Random | 62.6 | 0.83 | 38.7 | 0.28 | 6.8% |
| | SVMRank | 55.7 | 0.76 | 38.8 | 0.03 | 4.5% |
| | AutoScore | 59.4 | 0.80 | 43.7 | 0.02 | 8.3% |
| Trans-former | Top-1 | 43.6 | 0.64 | 0.30 | 0.12 | 0.8% |
| | Random | 42.4 | 0.63 | 0.30 | 0.13 | 1.0% |
| | SVMRank | 41.8 | 0.63 | 0.31 | 0.10 | 1.2% |
| | AutoScore | 40.5 | 0.62 | 0.30 | 0.10 | 1.3% |
| LSTM | Top-1 | 36.2 | 0.56 | 0.28 | 0.10 | 0.3% |
| | Random | 36.0 | 0.56 | 0.28 | 0.10 | 0.3% |
| | SVMRank | 36.2 | 0.56 | 0.29 | 0.10 | 1.0% |
| | AutoScore | 34.1 | 0.52 | 0.28 | 0.10 | 1.0% |

Table 7: Automatic evaluation: Results for each combination of generation model and candidate selection strategy on the 600 test samples, in comparison to the human revisions: BLEU (0-100), ROUGE-L (RouL), SARI, ratio of unedited samples (NoEd), % of exact matches to target reference (ExM).

| Model | Strategy | Fluency | Meaning | Argument | Average |
|---|---|---|---|---|---|
| BART | Top-1 | 0.73 | 0.97 | 0.65 | 0.78 |
| | Random | 0.72 | 0.97 | 0.68 | 0.79 |
| | SVMRank | 0.72 | 0.94 | 0.76 | 0.81 |
| | AutoScore | 0.83 | 0.95 | 0.86 | **0.88** |
| Trans-former | Top-1 | 0.44 | 0.76 | 0.40 | 0.53 |
| | Random | 0.41 | 0.76 | 0.38 | 0.52 |
| | SVMRank | 0.50 | 0.76 | 0.45 | 0.57 |
| | AutoScore | 0.68 | 0.75 | 0.61 | **0.68** |
| LSTM | Top-1 | 0.27 | 0.68 | 0.31 | 0.42 |
| | Random | 0.27 | 0.68 | 0.31 | 0.42 |
| | SVMRank | 0.29 | 0.69 | 0.31 | 0.43 |
| | AutoScore | 0.52 | 0.65 | 0.53 | **0.57** |
| Human | | 0.72 | 0.94 | 0.74 | 0.80 |

Table 8: Results for each combination of generation model and candidate selection strategy on the 600 test samples, in comparison to the human revisions based on three quality metrics: fluency, meaning preservation and argument quality.

We use the *lstm_wiseman_iwslt_de_e* architecture, which is a two-layer encoder and decoder LSTM, each with 256 hidden units, and dropout with a rate of 0.1 between LSTM layers.

**Transformer.** The second model is based on the work of Vaswani et al. (2017). We use the *transformer_iwslt_de_en* architecture, a 6-layer encoder and decoder with 512-dimensional embeddings, 1024 for inner-layers, and four self-attention heads.

Tables 7 and 8 compare the automatic evaluation scores of all generation-content selection combinations.

## B.1 Automatic Evaluation

We use the following python packages and scripts to perform automatic evaluations: nltk (BLEU (Papineni et al., 2002)), rouge-score (ROUGE (Lin, 2004)), https://github.com/cocoxu/simplification/SARI.py (SARI (Xu et al., 2016))

## C Claim Optimization Examples

For all eight optimization categories, we provide one or more examples illustrating each action in Table 9.

## D Manual Quality Assessment Guidelines

Figure 3 shows the annotation guidelines for the Amazon Mechanical Turk study.

## E System Outputs

Table 10 provides examples of candidates selected by different content selection strategies along with human references illustrating common patterns found in the results. Table 11 provides examples of candidates generated with and without utilizing context knowledge with insertions and deletions being highlighted in green and red fonts accordingly.

| Type | Examples |
|------|----------|
| Specification | Nipples are the openings of female-only exocrene glands that can have abnormal [secretions] <LINK> during any time of life, get erected by cold stimulation or sexual excitement (much more visibly than in men), get lumps or bumps and change color and size of areola during the menstrual cycle or pregnancy, so their display can break [personal space] <LINK> and privacy (which is stressful), affect public sensibilities and also be a [window] <LINK> for infections, allergies, and irritation. <br><br> The idea behind laws, such as limiting the amount of guns, is to reduce the need to defend yourself from a gun or rapist. <br><br> It is very common for governments to actively make certain forms of healthcare [harder for minority groups to access] <LINK>. They could also, therefore, make cloning technology hard to access. |
| Simplification | Very complex, ~~cognitively meaningful behavior such as~~ behaviours like creating art are evidence of free will, ~~because they exhibit the same lack of predictability as stochastic systems, but are intelligible and articulate clearly via recognizable vehicles~~. |
| Reframing | It reduces the oversight of the BaFin and thus increases ~~the risk of financial crisis~~ market failures. |
| Elaboration | It takes 2-4 weeks for HIV to present any symptom. The incubation period risk ~~can't be ruled out for~~ is higher for a member of high risk group, ~~effectively and timely~~ even though member of a low risk group is not completely safe. The decision is based on the overall risk, not on individual level. |
| Corroboration | [Person-based predictive policing technologies] <LINK> - that focus on predicting who is likely to commit crime rather than where is it likely to occur - violate the [presumption of innocence.] <LINK>. |
| Neutralization | Biden ~~does not~~ lacks the support ~~or agree with several key issues that are important to liberal voters.~~ of many liberal voting groups due to his stance on key issues concerning them. |
| Disambiguation | The USSR had [passed legislation] <LINK> to gradually eliminate religious belief within its borders. However the death penalty was more used in USSR than in Russia. ~~It~~ USSR had 2000 [death penalties] <LINK> per year in the 1980s whereas pre USSR Russia had [banned the death penalty] <LINK> in 1917 and almost never carried it out in the decades before that. <br><br> ~~SRM~~ Solar geoengineering merely serves as a "technological fix" (Weinberg).[harvard.edu] <LINK> |
| Copy Editing | Women are experiencing record ~~level~~ levels of success in primaries. |

Table 9: Illustrative examples of optimization types identified in the paper. The green font denotes additions and the striked out red font denotes the removal of text snippets.

# Instructions

In this task, your goal is to identify whether a claim has been successfully improved, without changing the overall meaning of the text.

Each task contains a set of pairs, where one claim is the "original claim," and the other an optimized candidate. Each of these pairs have the same original text, but different candidate optimizations.

Please rate each candidate along the following three perspectives: argument quality, fluency and semantic similarity. And, finally, please, rank all candidates relative to each other in terms of overall quality.

## Argument Quality

**Scale (1-5)**: 1 (notably worse than original), 2 (slightly worse), 3 (same as original), 4 (slightly improved), 5 (notably improved)

Does the optimized claim improve the argument quality compared to the original claim? Relevant changes include, but are not limited to:

- further specifying or explaining an existing fact or meaning
- removing information or simplifying the sentence structure with the intent to reduce the complexity or breadth of the claim
- rephrasing a claim with the intent to specify or generalize the claim, or to add clarity
- adding (substantive) new content or information to the claim or inserting an additional fact with the intent of making it more self-contained, more sound or stronger
- adding, editing or removing evidence in the form of links that provide supporting information or external resources to the claim
- removal of bias or biased language
- removal uncertainty. e.g. by replacing pronouns referring to concepts that have been mentioned in other claims of the debate, or by replacing acronyms with what they stand for
- improving the grammar, spelling, tone, or punctuation of a claim

## Meaning

**Scale (1-5)**: 1 (entirely different), 2 (substantial differences), 3 (moderate differences), 4 (minor differences), 5 (identical)

Does the transformed claim still have the same overall meaning as the original? It is OK if extra information is added, as long as it doesn't change the underlying people, events, and objects described in the sentence. You should also not strongly penalize for meaning transformations which aim to generalize or specify some aspects of the claim.

## Fluency

**Scale (1-3)**: 1 (major errors, disfluent), 2 (minor errors), 3 (fluent)

Is this sentence fluent English and does it make sense?

Examples with explanations can be found here.

Figure 3: Guidelines provided to the crowdworkers to annotate automatically generated claim optimizations.

| Strategy | Candidate |
|---|---|
| Original | Easy to cheat, and procastinate |
| Top-1 | It is easy to cheat and procastinate . |
| SVMRank | The easy to cheat, and procastinate |
| AutoScore | The UBI is easy to cheat and easy to manipulate. |
| Human | Flexible schedules can easily increase procrastination |
| Original | Many games often attempt to depict the norms and attitudes of the real world. |
| Top-1 | Many games often attempt to depict the norms and attitudes of the real world. |
| SVMRank | Many games often attempt to depict the norms and attitudes of the real world, and thus are likely to be accurate and relevant to what people are thinking and feeling in the real-world. |
| AutoScore | Many games often attempt to depict the norms and attitudes of the real world, in order to give players a greater understanding of the world around them. |
| Human | Many games often attempt to depict the norms and attitudes of the real world; gamers are likely to see social attitudes espoused in these games on issues of gender equality as indicative of progress in the real world. |
| Original | Owners of PWA take home [more profit] <LINK> from their sales than mobile apps. |
| Top-1 | Owners of PWA take home [more profit] <LINK> from their sales than mobile apps. |
| SVMRank | Owners of PWA take home [more profit] <LINK> from their sales than mobile apps. This allows them to reinvest in their business and increase their profits. |
| AutoScore | Owners of PWA take home [more profit] <LINK> from their sales than those of mobile apps. |
| Human | Owners of PWAs take home [more profit] <LINK> from their sales than mobile apps. |
| Original | Scientists found out that Saturn's moon [Enceladus] <LINK> houses a giant, global ocean beneath its icy outer shell. |
| Top-1 | Scientists found out that Saturn's moon [Enceladus] <LINK> houses a giant, global ocean beneath its icy outer shell. |
| SVMRank | Scientists found out that Saturn's moon [Enceladus] <LINK> houses a giant, global ocean beneath its icy outer shell. This suggests that this icy, water-filled place also houses other life forms, such as asteroids and comets. |
| AutoScore | Scientists found out that Saturn's moon [Enceladus] <LINK> houses a giant, global ocean beneath its icy outer shell, suggesting that the moon may hold life beyond the surface. |
| Human | Scientists found ~~out~~ that Saturn's moon [Enceladus] <LINK> houses a giant, global ocean beneath its icy outer shell. |
| Original | Responsibility requires conscious decisions and [free will] <LINK>. Acting according to governmental orders is far from that. |
| Top-1 | Responsibility requires conscious decisions and [free will] <LINK>. Acting according to governmental orders is far from that. |
| SVMRank | Responsibility requires conscious decisions and [free will] <LINK>. Acting according to governmental orders is far from that. Therefore, people are not responsible for their actions, and are therefore not morally responsible. |
| AutoScore | Responsibility requires conscious decisions and [free will] <LINK>. Acting according to governmental orders is far from that, as it is not a conscious decision but rather a result of an involuntary reaction. |
| Human | The coercion to vote eliminates the idea of responsibility, which requires conscious decisions and [free will] <LINK>. Acting according to governmental orders is far from that. |
| Original | There was much progress, advancement and forward looking thinking in the 20th century in Europe across all disciplines. |
| Top-1 | There was much progress, advancement and forward looking thinking in the 20th century in Europe across all disciplines. Therefore, it is logical to conclude that the world has changed in the last 150 years. |
| SVMRank | There was much progress, advancement and forward looking thinking in the 20th century in Europe across all disciplines, which shows the progress that was made during the century. |
| AutoScore | There was much progress, advancement, and forward looking thinking in the 20th century in Europe across all disciplines. |
| Human | There was much progress, advancement and forward ~~looking~~ thinking in the 20th century in Europe across all disciplines. |
| Original | A [catch-22] <LINK> situation currently exists in regards to researching the medicinal applications of some illegal drugs, due to the laws surrounding how they are categorised ([p. 12] <LINK>). |
| SVMRank (Top-1) | A [catch-22] <LINK> situation currently exists in regards to researching the medicinal applications of some illegal drugs, due to the laws surrounding how they are categorised ([p. 12] <LINK>). This prevents researchers from knowing how to fully assess the medicinal potential of drugs in a manner that is best suited for the patients who need it. |
| AutoScore | A [catch-22] <LINK> situation currently exists in regards to researching the medicinal applications of some illegal drugs, due to the laws surrounding how they are categorised ([p. 12] <LINK>). This prevents researchers from researching possible medicinal uses for illegal drugs. |
| Human | A [catch-22] <LINK> situation currently exists in regard to researching the medicinal applications of some illegal drugs, due to the laws surrounding how they are categorised ([p. 12] <LINK>). |

Table 10: Examples of optimizations selected by different candidate selection strategies along with human references. The green font denotes additions and the striked out red font denotes the removal of text snippets.

| Control | Candidate |
|---|---|
| Topic | Should Men Be Involved And Trusted In Childcare? |
| Previous | Perpetrators of child abuse are [way more likely] <LINK> to be the child's parents. |
| Original | Mostly the mothers, too. |
| Claim only | The mothers are the most vulnerable. |
| + Previous | The majority of perpetrators are the mothers, too. |
| + Topic | It is mostly mothers who are responsible for childcare. |
| Topic | Would Trump or Biden be a better President? |
| Previous | The [Healthy Indiana Plan 2.0] <LINK>, initiated by Mike Pence, has provided access to healthcare to many people who could not previously afford it. |
| Original | The Healthy Indiana Plan has high penalties for non-payment of [premiums] <LINK>. |
| Claim only | The Healthy Indiana Plan has high penalties for non-payment of [premiums] <LINK>. This means that people with low income could be penalized for failure to pay. |
| + Previous | The Healthy Indiana Plan has high penalties for non-payment of [premiums] <LINK> in exchange for providing healthcare. |
| + Topic | The Healthy Indiana Plan has high penalties for non-payment of [premiums] <LINK>,which were imposed in 2015. |
| Topic | Artificial Intelligence (AI): Should an artificial general intelligence be created? |
| Previous | There is the chance an AGI might be able to think for itself and turn on its creators. |
| Original | AGI are susceptable. |
| Claim only | Many AGI are susceptable to the danger of death. |
| + Previos | AGI are susceptible. They may be influenced by humans and by the biases of their creators. |
| + Topic | AGI are susceptibleto being hacked. |
| Topic | Who should the Conservative Party choose to be the next UK Prime Minister? |
| Previous | [Jeremy Hunt] <LINK> |
| Original | Hunt is the best placed candidate to unify Parliament in order to get a Brexit deal passed. |
| Claim only | Hunt is the best placed candidate to unify Parliament and the public in order to get a Brexit deal passed. |
| + Previous | Hunt could use his position as the Prime Minister to unify Parliament in order to get a Brexit deal passed. |
| + Topic | Hunt is the best placed candidate to [unify Parliament] <LINK> in order to get a Brexit deal passed |
| Topic | Should high-income countries take in refugees? |
| Previous | Taking in refugees will increase criminality within host countries. |
| Original | Differences in criminal activity are a product, in part, of childhood [socioeconomic] <LINK> conditions. This is exacerbated by the longer [path] <LINK> to employment faced by refugees compared to other communities. |
| Claim only | Differences in criminal activity are a product, in part, of childhood [socioeconomic] <LINK> conditions. This is exacerbated by the longer [path] <LINK> to employment faced by refugees compared to other communities, making them more likely to get involved in crime. |
| + Previous | Differences in criminal activity are a product, in part, of childhood [socioeconomic] <LINK> conditions. This is exacerbated by the longer [path] <LINK> to employment faced by refugees compared to other communities. This will not increase criminality. |
| + Topic | Differences in criminal activity are a product, in part, of childhood [socioeconomic] <LINK> conditions. This is exacerbated by the longer [path] <LINK> to employment faced by refugees compared to other communities, which make it harder to find a job. |
| Topic | Mark Twain used the N-word in The Adventures of Huckleberry Finn. Should it be censored? |
| Previous | Changing the N-word would skip a piece of the *linguistic* past and thus everyday life. As a result, people could start to forget this part of history. |
| Original | In Huckleberry Finn, Twain captured the essence of "[everyday midwest American English] <LINK>". |
| Claim only | In Huckleberry Finn, Twain captured the essence of "[everyday midwest American English] <LINK>".This is a common trait of the American English language. |
| + Previous | In Huckleberry Finn, Twain captured the essence of "[everyday midwest American English] <LINK>"by using the N-word in everyday conversation. |
| + Topic | In Huckleberry Finn, Twain captured the essence of "[everyday midwest American English] <LINK>", which is a language that is often used by people who do not share his values. |

Table 11: Examples of different candidates generated by BART + AutoScore with and without context information. The green font denotes additions of text snippets.

# ChatGPT's Information Seeking Strategy: Insights from the 20-Questions Game

**Leonardo Bertolazzi**
University of Trento
leonardo.bertolazzi@unitn.it

**Davide Mazzaccara**
University of Trento
davide.mazzaccara@unitn.it

**Filippo Merlo**
University of Trento
filippo.merlo@studenti.unitn.it

**Raffaella Bernardi**
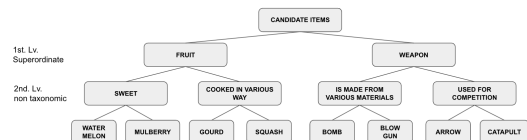University of Trento
raffaella.bernardi@unitn.it

## Abstract

Large Language Models, and ChatGPT in particular, have recently grabbed the attention of the community and the media. Having reached high language proficiency, attention has been shifting toward its reasoning capabilities. In this paper, our main aim is to evaluate ChatGPT's question generation in a task where language production should be driven by an implicit reasoning process. To this end, we employ the 20-Questions game, traditionally used within the Cognitive Science community to inspect the information seeking-strategy's development. This task requires a series of interconnected skills: asking informative questions, stepwise updating the hypothesis space, and stopping asking questions when enough information has been collected. We build hierarchical hypothesis spaces, exploiting feature norms collected from humans vs. ChatGPT itself, and we inspect the efficiency and informativeness of ChatGPT's strategy. Our results show that ChatGPT's performance gets closer to an optimal agent only when prompted to explicitly list the updated space stepwise.

## 1 Introduction

ChatGPT's impressive ability to solve numerous natural language tasks has put it in the spotlight of Academia and media attention (Bang et al., 2023; Laskar et al., 2023). The success on a variety of tasks has brought people to even claim that GPT-4 "could reasonably be viewed as an early (yet still incomplete) version of an artificial general intelligence (AGI)" (Bubeck et al., 2023). Others are more cautious, showing the weakness of the model's reasoning abilities, (e.g., Bang et al. 2023).

A core aspect of human intelligence is the implicit connection between the reasoning process and language production. This connection strongly drives the generation of questions in information-seeking scenarios which, therefore, have been largely studied in Cognitive Science. After the



Figure 1: **Upper part**: an example of a Hierarchical Hypothesis Space built with ChatGPT-feature norms. **Bottom part**: an example of an optimal Questioner which always divide the space (Sp) into half (starting with 8 candidates, going to 4:4, and then 2:2). Questions at turns 1-2 are constraint-seeking (CS), while the 3rd and 4th are hypothesis-scanning (HS). With the halfsplit procedure, the target can be guessed with just 3 turns or at most with 4, as in the example.

pioneering work by Mosher and Hornsby (1966), the 20-Questions game has been employed to observe children's cognitive developmental trajectory: A player thinks of an entity, the second player is given a set of candidates (e.g., *cat*, *dog*, *bird*) and has to identify the target entity among the possible candidates by making Yes-or-No questions. This and following experiments have shown that through the developmental trajectory, children learn to recognize object-general features, cluster similar objects into categories and use such categorization to ask context dependent informative questions: they shift from *Hypothesis Scanning* questions ("*Is it a dog?*") to *Constraint-Seeking* questions ("*Does it has four legs?*"). Such a shift let the elder children be more efficient in their information seeking process. Moreover, pre-scholar children tend to continue asking questions when enough information has been collected (i.e., the space has reduced to one candidate). They do not know when to stop (Ruggeri et al., 2016), a core skill of infor-

mation search and decision-making (Todd et al., 1999). Interestingly, Ruggeri et al. (2021) uses a hierarchical version of the 20-Q game, in which candidates are organized into three category levels based on shared features; by providing children with the object-related features needed to halfsplit the space, children were able to target such higher category levels, reaching the solution more efficiently. Inspired by this literature, we use a hierarchical 20-Q game to evaluate whether ChatGPT is able to generate questions driven by its reasoning over the Hypothesis Space (HypSp).

We leverage on the widely used feature norms elicited from human annotators, McRae-norms (McRae et al., 2005) to build the hierarchical hypothesis spaces. Such norms reflect humans' knowledge representation which could differ from ChatGPT's knowledge. To mitigate this potential difference, we build also hypothesis spaces using norms elicited from ChatGPT itself, (GPT-norms (Hansen and Hebart, 2022)). Figure 1 reports an example of an 8 candidate symmetric nested space based on GPT-norms.

We prompt ChatGPT[1] to play the 20-Q game, both in the role of the Questioner and of the Answerer.[2] We aim to understand whether the (a) Questioner is able to identify the high-level property that clusters the space and hence asks whether the target has that property; (b) it also knows whether all the other candidates have or do not have that property, and is able to use such information to update the HypSp stepwise; and finally (c) it understands when to stop asking questions, i.e., the HypSp is reduced to a singleton. Figure 1 includes a dialogue an optimal agent could ask, if driven by an half-split search. Our results show that ChatGPT's performance is far from an optimal agent when having to update the space internally and it is closer to it when prompted to explicitly list the updated space stepwise.

## 2 Related Work

Our work put together two research lines: the current effort of the AI community to evaluate Chat-GPT language and reasoning skills, and the cognitive science literature focusing on the developmental trajectory of information search strategies in humans.

**ChatGPT evaluation** Bang et al. (2023) run a deep and broad evaluation of ChatGPT on a variety of well recognized benchmarks in the Natural Language Processing community. ChatGPT results to be State-of-the-Art in zero-shot setting for most natural language understanding tasks. Though it is more suitable for open-domain dialogue tasks, it performs well also in task-oriented dialogues, and it is able to keep track of information given in previous turns, when answering follow up questions. Moreover, Bang et al. (2023) evaluate ChatGPT reasoning skills: though it lacks inducting reasoning skills, it performs well on deductive reasoning in clean settings. However, as other LLMs (Ott et al., 2023), ChatGPT as well encounters problems with complex deductive reasoning involving multi-hops, viz. a combination of facts spread in different passages of a corpus. Zhu et al. (2023) challenged ChatGPT on the Visual Dialogue task, originally proposed by Das et al. (2017). The informativeness of the question is measured on the quality of the caption it summarises out of the dialogue. As far as we know, this is the first work to evaluate the information seeking strategy of a LLM using the 20 Questions game. Our research question is whether and to what extent the language generation of a LLM is tied to reasoning.

**Developmental and Cognitive Psychology** Starting with Mosher and Hornsby (1966), the 20-Q search task has been largely used in developmental and cognitive psychology. Among the measures to evaluate the question's informativeness, Expected Information Gain (Lindley, 1956) emerges as one of the most used. It values questions with respect to the uncertainty reduction, and it is usually connected with the prior probability. Subjects have been evaluated with the 20-Q game considering both scenarios simulating *prior* expectations and scenarios with uniform distribution (e.g., Ruggeri and Lombrozo 2015; Meder et al. 2019; Ruggeri et al. 2021; Testoni 2023). Our scenario is the uniform distribution.

It is widely accepted that children's search strategies are less efficient than adults' ones. Rather than identifying high-order properties splitting efficiently the Hypothesis Space, indeed, children tend to scan the space item by item. In complex scenarios, it has been shown that adults do not efficiently plan ahead; they tend to follow a half-splitting strategy: ask the question that more closely approximates a division of the space into half (Meder

---

et al., 2019). Rothe et al. (2018) show that people can accurately evaluate questions quality, but have limited ability to optimize the informativeness of their questions. By leverage of feature norm collections, we work on a simplified scenario where adults would more easily stay close to an optimal agent.

**Feature Norms**  Feature norms refer to minimal semantic descriptions that capture the typical attributes associated with a collection of objects or concepts (e.g., a *dog* can be described by features such as *has fur* and *does bark*). One common method of acquiring semantic features for concepts is to ask individuals to list properties associated with a given concept. A broadly used collection is the McRae-norms (McRae et al., 2005) which comprise 2524 unique features collected from approximately 725 participants, which are in turn categorized according to Wu and Barsalou (2009)'s taxonomy of relations (WB). These norms encompass 541 animate and inanimate concrete concepts, with an average of 30 participants providing feature listings for each of them. The features included in the McRae-norms are of various types, such as physical (perceptual) properties, functional properties, taxonomic properties, and encyclopedic facts. Inspired by this work and to obtain a large-scale collection, (Hansen and Hebart, 2022) instructed GPT-3 (Brown et al., 2020) to generate semantic feature norms for a diverse set of 1,854 concrete concepts which have been annotated with 84561 unique features elicited from the model through 30 runs, pre-processed and filtered. These feature norms were then released by the authors for public use and exploration; the authors expanded their method to other models of the GPT family when they became available. We use the feature norms obtained from GPT-3.5-turbo, and refer to them as GPT-norms.[3]

## 3   Hierarchical version of the 20-Q game

Following Ruggeri et al. (2021), we created a hierarchical version of the 20-Q game. In other words, the hypothesis spaces are built out of two subsets of equal size (N:N), and iteratively divided into further subsets based on some other features. We exploit McRae (McRae et al., 2005) and ChatGPT (Hansen and Hebart, 2022) feature norms to build the nested

structures. We consider hierarchies of two levels (8 candidates, divided into 4:4, and 2:2) and of three levels (16 candidates, divided into 8:8, 4:4, and 2:2). The first level is always based on superordinate properties ($F_{1a}$ vs. $F_{1b}$), which are by definition mutually exclusive (e.g., *bird* vs. *mammal*, *fruit* vs. *weapon*, etc.). The subsets of the other levels instead are obtained from all the other feature norms associated with the candidates (e.g., items that are *fruit* could be divided into those that are *sweet* vs. those that are *cooked in various way*). We make sure that the feature that is shared by half of the candidates is not listed for any of the item in the other half, and viceversa. The leaves of the hierarchy are randomly selected among the concepts of the corresponding groups. See Figure 1 for the schema and an example with 8 candidates organized based on GPT-norms.

**Hierarchical Hypothesis Space creation**  Our starting point are the concepts in McRae et al. (2005) and in Hansen and Hebart (2022), 541 and 1854, respectively. For McRae-norms, we selected the superordinates frequent enough to let us create spaces of 8 and of 16 candidates.[4]  For the second level, we use features of the other 8 most frequent WB relations (51 unique features). For the Hypothesis Spaces built from ChatGPT-norms, we selected the same 6 superordinates for the first level splits, and other 806 most frequent unique features for the second levels. We built the hypothesis space through a recursive process that guarantees variety and randomness of the selection (See the Supplementary Material for details). We will refer to these two types of hypothesis spaces as 8 vs. 16 candidates sets (cds), distinguishing the former into McRae- and GPT-based 8 cds; henceforth, `8-McRae`, `16-McRae`, `8-GPT`.

**Game creation**  A game consists of a set of candidates, assigned to the Questioner player, and a target among them, assigned to the Answerer. We build 90 games for each of the three types of Hypothesis Space as follows. First of all, out of the 6 selected superordinate features, we build all possible pairs, viz. 15 ($F_{1a}$, $F_{1b}$); we then randomly select 6 sets of candidates for each of the 15 pairs, yielding 90 unique sets (total 270 sets). Finally, we build the 90 games by randomly selecting the target 3 times from the candidates that share $F_{1a}$

---

[3] The norms collected with `gpt-3.5-turbo` are available at `https://github.com/ViCCo-Group/semantic_features_gpt_3`.

[4] The 6 superordinates we use to build the first level splits are: mammal, bird, clothing, weapon, fruit, vegetable.

and 3 times from those that share $F_{1b}$. This process guarantees variety of the concepts and targets.

## 4 Agent Roles

Below we describe how we employ ChatGPT as game players to generate the dialogues and as diagnostic agents to evaluate the Questioner's information seeking ability.

### 4.1 Game Players

To generate the dialogues, ChatGPT is instructed to play the role of the Questioner (**ChatGPT-Q**) and of the Answerer (**ChatGPT-A**) with a similar system prompts. The shared part of the prompts explicitly states the only possible answers are 'yes' and 'no'. ChatGPT-Q is told to ask as few questions as possible; the Questioner starts by asking the first question, which is appended to the Answerer's prompt in order to generate the first answer. In this way, the dialogue history is iteratively increased after each turn. ChatGPT-A is told to acknowledge when the Questioner has correctly guessed the item by answering "Yes! That's correct.". Focusing on ChatGPT's capabilities of reasoning about the hypothesis space and asking questions that reflect such reasoning, we retain only successful dialogues. More precisely, the dialogue is kept if the Answerer considers the target reached. Our evaluation is focused on the Questioner role, hence, for it we define theoretically an upper and a lower bound as described below.

We take as upper-bound a model that similarly to adults seeks for a property shared by several items in the space. In particular, we use the **optimal agent** which acts similarly to a binary search algorithm: at each turn, it divides in half the hypothesis space under discussion ($N/2$). When only two items are left, the optimal agent makes a guess that has the 50% chance of being the correct target. This half-split strategy takes on average $\log_2 N + 1/2$ turns to solve the game, where $N$ is the number of items at the beginning of the game.

As lower-bound we consider a model close to the 4-Y child who tends to scan the space item by item. Therefore, our **baseline agent** acts similarly to a linear search algorithm: at each turn, it divides the space into $1\ vs.\ N-1$. Given $N$ items at the beginning of the game, it takes on average $N/2$ turns to solve the game.

## 4.2 Diagnostic Agents

To evaluate the model's ability to stepwise reduce the hypothesis space we exploit ChatGPT in the role of an external Oracle (ChatGPT-Oracle), and of an external Guesser (ChatGPT-Guesser). Moreover, we activate the Guesser internal to the Questioner by prompting the model to update the candidates at each turn (ChatGPT-Q-stepwise).

**ChatGPT-Oracle** is given a question in the dialogue sets described in the previous section and for each item in the hypothesis space of the corresponding game says whether the item has or does not have the required property.[5] This provides us with Y/N-annotation of the hypothesis space that we use to obtain a "ground truth" updated space at each turn. The feasibility of such method relies on the fact that the dialogues are rather simple and no actual linguistic dependencies are in place between the turns (See Supplementary Material for details). **ChatGPT-Guesser** is given chucks of the dialogue histories generated by the game players and is asked to list the candidates till the given turn. Finally, we modify the prompt of ChatGPT-Q by asking it to list the candidates under discussion stepwise before asking a new question (**ChatGPT-Q-stepwise**). The prompts used for each role are reported in the Supplementary Material.

## 5 Experimental Setup

We are interested in understanding whether ChatGPT's language generation is driven by its reasoning process. To answer this question, we propose a number of measures aimed to shed light on the reasoning processes that are implicit in the game: identify the high-level property shared by several items, update the space stepwise, and efficiently arrive to a space with just one possible candidate and realize that it is time to stop asking questions. Not having the possibility to run an ablation study of the model, we simulate it by comparing ChatGPT-Q, simply prompted to play the game, with ChatGPT-Q-stepwise which is explicitly told to update the space turn by turn.

**Information seeking strategy** Following the method used in the Cognitive Science literature to evaluate children's developmental phases, we evaluated the information seeking strategy used

---

[5]We verified the reliability of the annotation by evaluating the model's accuracy on a sample of 180 questions: it correctly answered $83\%$ of the questions.

by ChatGPT-Q by observing the type of questions it asks and their informativeness. First of all, we compute the percentage of questions that are *Hypothesis-Scanning (HS)* and *Constraint-Seeking (CS)*. A question is considered *HS* iff it explicitly mentions one of the candidates in the hypothesis space. All the other questions are considered *CS*. We compute the percentage of HS and CS questions within a game and by the position of the turn within the dialogue.[6]

Following Ruggeri et al. (2016), Meder et al. (2019) and Testoni (2023), we compute the *Expected Information Gain (EIG)* of each question and report the average EIG per turn.[7] As clearly explained in Meder et al. (2019), the information gain of a question is the entropy in the space (given by the number of items and the associated probability) at turn $t_i$ before asking the question minus the expected entropy after asking it ($t_{i+1}$):

$$IG = H_{t_i} - H_{t_{i+1}}$$

As in Meder et al. (2019), in our case, we consider all items in the space to be equally likely to be the target. Hence, what defines entropy is the number of items within the subsets answered with Yes vs. No, based on the external Oracle annotation.

A model that asks fewer HS, especially in the earlier turns is closer to the more efficient strategy used by adults. Its question EIG is expected to be very high in the early turns and to decrease in the later ones.

**Hypothesis space update**  A core skill of the Questioner playing the 20-Q game is the ability to *mentally* keep the space of the hypothesis updated stepwise. We evaluate whether ChatGPT-Q is able to update at each turn the hypothesis space based on the given dialogue history. Again, we consider the Yes/No-annotation obtained from ChatGPT-Oracle as the ground truth and compute the hypothesis space at turn $HypSp_{t_i}$ by filtering out from $HypSp_{t_{i-1}}$ the items which do not have the property required at $t_i$. We compare the ground truth Hypothesis Space with a) the one generated by the external Guesser, ChatGPT-Guesser, and b) the one generated by the Questioner itself when

---

[6]A third type of questions are the *Pseudoconstraint-seeking (PCS)* which ask about a property but actually refer to only one item among the candidates. For the sake of simplicity, we do not consider them in our analysis, but see the Supplementary Material for statistics on them.

[7]We computed the EIG adapting the code by Testoni (2023).

prompted to explicitly update the list of candidates stepwise (ChatGPT-Q-stepwise). To this end, we compute the symmetric difference between the generated sets with the ground-truth ones for every question, and report the average symmetric difference of each game turn. The symmetric difference between two sets $A$ and $B$ is denoted by $A\Delta B$ and is defined as follows:

$$A\Delta B = (A - B) \cup (B - A)$$

For ChatGPT-Guesser, a high difference would mean that the model has difficulty integrating the information collected through the dialogue history. While for ChatGPT-Q-stepwise it would signal a difficulty in integrating the answer with the question turn by turn.

**Search efficiency**  We measure the efficiency of the ChatGPT-Q's game strategy by computing the *average number of questions per game, (AQ)*. In addition to this, and as in (Ruggeri et al., 2016), we consider a question *unnecessary (UQ)*, if the preceding dialogue history already contained the information to identify the target. Again, we use the Y/N-annotation by ChatGPT-Oracle to determine whether this point has been reached by ChatGPT-Q.

The more the search strategy is effective, the shorter is the dialogue. The higher the number of UQ the closer is the model to pre-scholar children, who have been shown not to have learned the stopping rule yet. If the model asks just one UQ as last turn, that would still qualify it adult-like, since adults have been shown to ask a confirmation question before making the final guess (Testoni et al., 2022).

**Experimental Settings**  We expect that bigger candidate sets could challenge the model's capacity to keep track of the information obtained through the dialogue, since they might require longer interactions. Moreover, with the GPT-based 8 cds the model should have all the knowledge to quickly arrive to identify the target. Hence, if ChatGPT's knowledge properly drives its question generation the dialogues of the games based on it should display an almost optimal information seeking strategy. Based on these conjectures, we compare the model when playing games whose hypothesis space a) consists of 8 and 16 candidate sets, b) is built out of McRae- or ChatGPT-feature norms. If the question generation is driven by the reasoning process on the space, we expect the model's
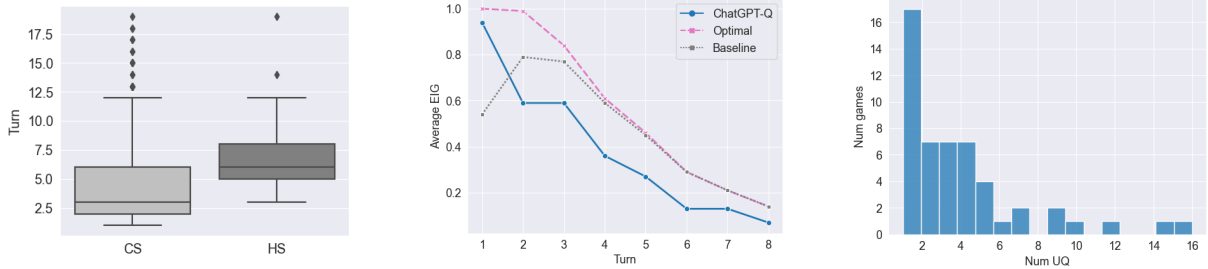
Figure 2: **Left**: ChatGPT-Q asks CS questions mostly in the earlier turns and HS in the later ones. Mann-Whitney U test shows the difference is statistically significant $p < 0.001$; **Middle**: ChatGPT's EIG is almost as high as the optimal model's EIG at the first turn, but it is lower than of the baseline's EIG at later turns. **Right**: Distribution of unnecessary questions.

## 6 Results

In this section, we show the results we obtained following the experimental setup defined in Section 5. We start by evaluating ChatGPT's performance on the games with 8 candidates selected with the McRae-norms (8-McRae), and we then move to compare these results with those obtained by the model when challenged with spaces containing an higher number of candidates (16-McRae), or whose nested structure reflects the model's knowledge representation (8-GPT). Finally, we move to evaluate ChatGPT on the 8-McRae games when asked to play the game (ChatGPT-Q) and when asked to explicitly update the hypothesis space stepwise while playing the game (ChatGPT-Q-stepwise).

### 6.1 ChatGPT-Q on the McRae-8 games

Through the measures introduced above, here we aim to take a picture of how well and efficiently ChatGPT searches for information by considering McRae-8 games.

**Information seeking strategy** The results for the type of questions asked by ChatGPT-Q, the *optimal agent* and the *baseline* can be seen in Table 1. By construction, the optimal agent asks $log_2 N - 1$ CS question per game, followed by 1 or 2 HS questions (hence 1.5 on average), until it guesses the target; in other words, 57.14% of its questions are CS, and 42.86% are HS. Instead, the baseline asks only HS questions and, on average, it guesses the target in 4 questions. ChatGPT-Q asks mostly Constraint-Seeking questions (73.77%), it

| 8 cds based on McRae-norms | | | |
|---|---|---|---|
| | HS | CS | AQ |
| Optimal | 42.86 | 57.14 | 3.50 |
| Baseline | 100 | 0 | 4.00 |
| ChatGPT | 26.33 | 73.77 | 7.24 |

Table 1: Information seeking strategy: Upper and lower bound of the overall percentage of hypothesis scanning (HS) vs. constraint seeking (CS) questions, and the average number of questions per game (AQ) -the difference is statistically significant based on the Wilcoxon signed-rank test (e.g., wrt the optimal agent, $p < 0.001$.

tends to ask CS questions in the early turns and HS questions towards the end of the dialogue – when indeed the latter becomes more efficient to split the space (Figure 2, left).

Moreover, our results show that the EIG of ChatGPT-Q's questions through the dialogue is far from the optimal agent's EIG (that half-split the space at each turn) and even lower than the baseline's (that splits the space into 1 vs. the other candidate at each turn) (see Figure 2, middle). Summing up, on the surface level, the strategy used by ChatGPT-Q reflects what an adult would do. However, the EIG analysis shows that ChatGPT-Q asks more uninformative questions compared both to the optimal agent and the baseline.

**Hypothesis Space update** We evaluate whether ChatGPT-Guesser is able to list the candidates that are still possible candidates based on the Question-Answer exchanges between ChatGPT-Q and ChatGPT-A. We do so by computing the difference, turn by turn, of such list with those considered as "ground truth" based on ChatGPT-Oracle annotations. The pattern we find (see the blue line in Figure 3, right) suggests the model has difficulty in integrating the information collected through the
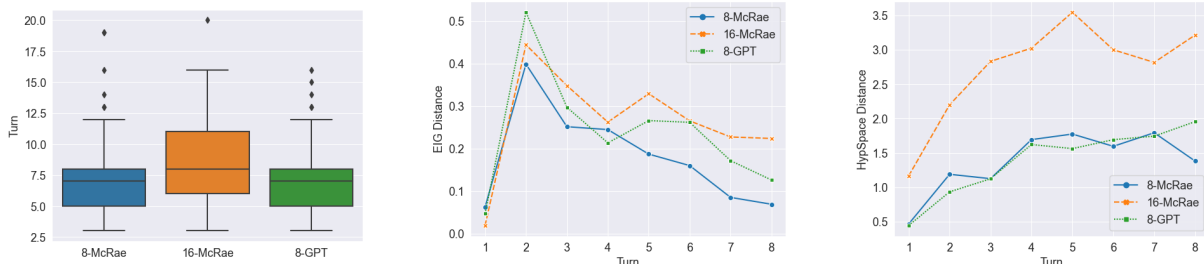
158

Figure 3: Increasing the candidates (**Blue vs. Orange**) causes the generation of longer dialogues (left) – maintaining a similar distance to the optimal model's AQ (8.67 vs 4.5); does not impact much the questions' EIG (middle), while it makes the guessing of the candidate by turn harder (right). Using candidates space structure based on GPT-norms (**Blue vs. Green**) does not impact any measure.

dialogue history; this weakness could impact its ability to stop asking questions when it has reached the singleton set.

**Search efficiency** By construction the optimal agent asks on average 3.5 questions per game, whereas the baseline asks on average 4 questions. ChatGPT-Q asks way more questions per game (7.24) compared to both models (Table 1). Such difference is statistically significant based on a Wilcoxon signed-rank test (e.g., wrt the optimal agent, $p < 0.001$). Moreover, ChatGPT-Q asks unnecessary questions (UQ) (questions asked after the singleton set has been reached) in 56.67% of games – 29.29% of its questions are unnecessary (See Figure 2, right for the distribution of UQ.) Summing up, in terms of search efficiency, ChatGPT's behavior is similar to that pre-scholar children who tend to not stop asking questions once there is only one item left in the hypothesis space.

## 6.2 Changing the games

Figure 3 illustrates how ChatGPT-Q's performance is effected by the Hypothesis space size (Blue vs. Orange) and of the features used to build it (Blue vs. Green). Hence, it compares the results discussed above (McRae-8) with those obtained in the other two settings: McRae-16 and GPT-8. In particular, it shows the comparison based on the number of questions per game (left), the distance between ChatGPT-Q and the optimal agent in terms of EIG (middle), and the average symmetric difference between the ground-truth hypothesis spaces update, based on ChatGPT-Oracle, and the one generated by ChatGPT-Guesser.

**ChatGPT-Q on 16-McRae** Given the difficulty the model has in keeping track of the space update, we expect that by increasing the number of can-

didates ChatGPT-Q' s performance will decrease. The results are not clear-cut: by moving from 8 to 16 candidates, the optimal agent would have an increase of 1 question per game, while Chat-GPT increases of 1.43; the difference in terms of EIG is low, while ChatGPT-Guesser's performance deteriorates.[8]

**ChatGPT-Q on 8-GPT** The games built out of GPT-norms should reflect the model knowledge representation, therefore we expect that on the 8-GPT games ChatGPT-Q performance will increase. Instead, for none of the measures the difference is significant. This suggests that the feature norms used to build the hypothesis spaces do not impact the model's performance.

## 6.3 Changing the prompt

To further understand what causes ChatGPT inefficient strategy, we would need to run an ablation study by isolating the various processes that should be beyond the question generation. To simulate such study, we compare the set of results discussed so far on McRae-8 obtained by ChatGPT-Q simply prompted to play the game, with those obtained by ChatGPT- Q-stepwise, the model that is explicitly told to update the space turn by turn. Our results show that ChatGPT-Q-stepwise gets closer to the optimal model: it asks fewer questions per game (Figure 4, left) compared to ChatGPT-Q (6.4 vs. 7.24), the questions' EIG is higher across all the turns (Figure 4, middle), and it is more precise when updating the hypothesis space (Figure 4, right). This finding confirms the conjecture that ChatGPT main weakness lies in its difficulty in *mentally* updating the hypothesis space.

---

[8]Games generated on 16-McRae are significantly longer than those generated on 8-McRae, based on a Mann-Whitney U test ($p < 0.001$).
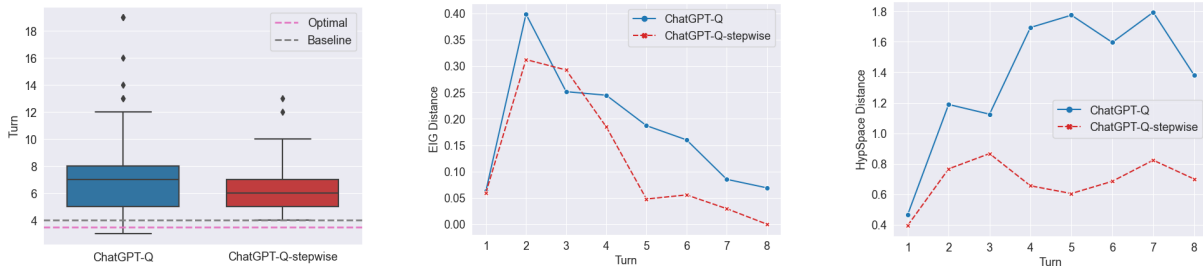
Figure 4: Changing the prompt improves the Questioner performance on all the measures, bringing it closer to the optimal model in terms of number of turns and EIG, an to the ground truth hypothesis space stepwise reduction.

| ChatGPT-Q | | | |
|---|---|---|---|
| | CQ | TQ | SG |
| McRae-8 | 13.50 | 25.15 | 12.22 |
| McRae-16 | 14.74 | 18.59 | 13.33 |
| GPT-8 | 13.22 | 26.90 | 8.89 |
| ChatGPT-Q-stepwise | | | |
| | CQ | TQ | SG |
| 8-McRae | 8.42 | 26.78 | 7.78 |

Table 2: Quality of the dialogue: Contradictory Questions (CQ) decrease when ChatGPT is asked to update the space explicitly. Yet, it is still unaware of Spoiled Games (SG).

## 7  Qualitative Analysis

In this section, we further dive into the quality of the dialogues generated by ChatGPT. First of all, we inspect whether it asks questions that do not reduce the space at all (trivial questions, TQ) or refer to candidates that have already been excluded in previous turns (contradictory questions, CQ). As we can see from the statistics reported in Table 2, ChatGPT is rather coherent through the dialogues; yet, the number of trivial questions is higher than what one would expect from a rational agent. In all the different settings, we observe a low peak in terms of EIG at the second turn. To understand the reason, we look into the 8-candidates sets (McRae and ChatGPT): 39.44% of the second questions are uninformative (EIG=0), with a large majority of trivial questions (92% of the uninformative). Interestingly, neither the size of the space nor the norms used to build it impact the number of trivial and contradictory questions. Instead, the coherence of the dialogues improves when the prompt is changed and ChatGPT is asked to update the space stepwise before asking the next question (ChatGPT-Q-stepwise on 8-McRae).

By inspecting the dialogues, we realized that in

all the settings (8 vs 16 cds, McRae vs. ChatGPT-norms, explicit vs. implicit-update), there are games in which ChatGPT continues asking questions even when the Answerer has accidentally revealed the target (Spoiled games, SG – see an example in the Supplementary Material). This suggests that the model is pretending to play the game without having actually grasped the actual purpose of it. Most probably, a spoiler would not pass unobserved by a 4Y-child.

## 8  Conclusion

Our work shows that ChatGPT is able to identify superordinate features shared by items and ask questions that efficiently reduce the hypothesis space. At the first turn, it is close to an optimal agent using a half-split search. In later turns, however, it has difficulties making questions with respect to the updated space of the hypothesis. This weakness might be behind the high number of games in which it keeps on asking questions even though the dialogue history had led to identifying a possible target. We conjecture this behavior is not due to the lack of knowledge required by the game since it is displayed not only within the games based on McRae norms but also on those built out of GPT feature norms. Our conjecture is reinforced by the increased performance reached by the model when prompted to explicitly update the space before asking the next question. In this setting, the dialogue becomes shorter with fewer contradictory questions. Yet, even in such scenario, it does not notice when the Answerer reveals the target accidentally. Our results call for attention to modeling the human ability to keep a *mental scoreboard*, echoing what stated in Lewis (1979); Madureira and Schlangen (2022); Mazuecos et al. (2021). Finally, our work relates to the Chain-of-Thought (Wei et al., 2022) and similar prompting strategies, which we plan to investigate in the future.

## 9 Limitation

The backbone of the hierarchical space we built are feature norms. For the first level split, we used superordinates which by definition are disjoint. For the second level, we used all other feature norms by making sure that the feature that is shared by a subset is not listed in any of the members of the other subset; this process does not guarantee disjointedness of the two subsets. McRae features norms associated to a concept should be salient for it, while the absence of a feature from the list could be either because the feature does not hold for that concept or because it is not salient. Nevertheless we choose to use McRae-norms because they reflect human representation of the world and gave us the possibility of having a straight comparison with the games built out of ChatGPT-norms – comparison which shows that the knowledge used for building the hypothesis spaces does not impact the model's performance. We evaluated the model also on games built with taxonomic relations extracted from WordNet for both levels of the hierarchy: the patterns are very similar to those obtained with McRae- and GPT-norms (See the Supplementary Material for details.)

A second limitation is due to ChatGPT being a closed-source model, for which the exact training data is not known. We leave for future work the study of a LLM open source. Finally, we have not compared the model results on those that humans playing the games would achieve, instead we rely on the results obtained within the Cognitive Science literature about the 20-Q game.

## References

Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Zi-

wei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. 2023. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *arXiv preprint arXiv:2302.04023*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

S. Bubeck, V. Chandrasekaran, R. Eldan, J. Gehrke, E. Horvitz, E. Kamar, P. Lee, Y.T. Lee, Y. Li, S. Lundberg, H. Nori, H. Planagi, M. T. Ribeiro, and Y. Zhang. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. ArXiv:2303.12712.

Abhishek Das, Satwik Kottur, José M.F. Moura, Stefan Lee, and Dhruv Batra. 2017. Learning Cooperative Visual Dialog Agents with Deep Reinforcement Learning. In *2017 IEEE International Conference on Computer Vision*, pages 2951–2960.

H. Hansen and M. Hebart. 2022. Semantic features of object concepts generated with gpt-3. In *In Proceedings of the Annual Meeting of the Cognitive Science Society*.

Md Tahmid Rahman Laskar, M Saiful Bari, Mizanur Rahman, Md Amran Hossen Bhuiyan, Shafiq Joty, and Jimmy Huang. 2023. A systematic study and comprehensive evaluation of ChatGPT on benchmark datasets. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 431–469, Toronto, Canada. Association for Computational Linguistics.

D. Lewis. 1979. Scorekeeping in a language game. *Semantics from different points of view*, pages 172–187.

D. V. Lindley. 1956. On a Measure of the Information Provided by an Experiment. *The Annals of Mathematical Statistics*, 27(4):986 – 1005.

Brielen Madureira and David Schlangen. 2022. Can visual dialogue models do scorekeeping? Exploring how dialogue representations incrementally encode shared knowledge. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 651–664, Dublin, Ireland. Association for Computational Linguistics.

Mauricio Mazuecos, Franco M. Luque, Jorge Sánchez, Hernán Maina, Thomas Vadora, and Luciana Benotti. 2021. Region under Discussion for visual dialog. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4745–4759, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

K. McRae, G. S. Cree, M. S. Seidenberg, and C. McNorgan. 2005. Semantic feature production norms for a large set of living and nonliving things. *Behavior research methods*, 4(37):547–559.

B. Meder, J. D. Nelson, M. Jones, and A. Ruggeri. 2019. Stepwise versus globally optimal search in children and adults. *Cognition*.

E. Mosher and J. R Hornsby. 1966. On asking questions. *Studies in cognitive growth*.

S. Ott, K. Hebenstreit, V. Liévin, C. E. Hother, M. Moradi, M. Mayrhauser, R. Praas, O. Winther, and M. Samwald. 2023. Thoughtsource: A central hub for large language model reasoning data. Arxiv.2301.11596.

A. Rothe, B.M. Lake, and T.M. Gureckis. 2018. Do people ask good questions? *Comput Brain Behav*, 1:68–89.

A. Ruggeri, T. Lombrozo, T. L. Griffiths, and F Xu. 2016. Sources of developmental change in the efficiency of information search. *Developmental psychology*.

A. Ruggeri, C.M. Walker, T. Lombrozo, and A. Gopnik. 2021. How to help young children ask better questions? *Frontiers in Psychology*, 11. Doi: 10.3389/fpsyg.2020.586819.

Azzurra Ruggeri and Tania Lombrozo. 2015. Children adapt their questions to achieve efficient search. *Cognition*, 143:203–216.

A. Testoni. 2023. *Asking Strategic and Informative Questions in Visual Dialogue Games: Strengths and Weaknesses of Neural Generative Models*. Ph.D. thesis, IECS Doctoral School, University of Trento.

A. Testoni, C. Greco, and R. Bernardi. 2022. Artificial intelligence models do not ground negation, humans do. guesswhat?! dialogues as a case study. *Frontiers in big data*, 4(736709).

Peter Todd, Jean Ortega, Jennifer Davis, Gerd Gigerenzer, Daniel Goldstein, Adam Goodie, Ralph Hertwig, Ulrich Hoffrage, Kathryn Laskey, Laura Martignon, and Geoffrey Miller. 1999. *Simple Heuristics That Make Us Smart*. Oxford University Press.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837. Curran Associates, Inc.

L. L. Wu and L. W. Barsalou. 2009. Perceptual simulation in conceptual combination: evidence from property generation. *Acta psychologica*.

D. Zhu, J. Chen, K. Haydarov, X. Shen, W. Zhang, and M. Elhoseiny. 2023. Chatgpt asks, blip-2 answers: Automatic questioning towards enriched visual descriptions. Arxiv.2303.06594.

# This is not correct! Negation-aware Evaluation of Language Generation Systems

**Miriam Anschütz** and **Diego Miguel Lozano** and **Georg Groh**
School for Computation, Information and Technology
Technical University of Munich, Germany
{miriam.anschuetz, diego.miguel}@tum.de, grohg@in.tum.de

## Abstract

Large language models underestimate the impact of negations on how much they change the meaning of a sentence. Therefore, learned evaluation metrics based on these models are insensitive to negations. In this paper, we propose *NegBLEURT*, a negation-aware version of the BLEURT evaluation metric. For that, we designed a rule-based sentence negation tool and used it to create the *CANNOT* negation evaluation dataset. Based on this dataset, we fine-tuned a sentence transformer and an evaluation metric to improve their negation sensitivity. Evaluating these models on existing benchmarks shows that our fine-tuned models outperform existing metrics on the negated sentences by far while preserving their base models' performances on other perturbations.

## 1 Introduction

Previous work has shown that large language models such as BERT ([Devlin et al., 2019](#)) lack understanding of negated phrases and do not attribute sufficient importance to the word "not" ([Ettinger, 2020](#); [Hosseini et al., 2021](#)). Nevertheless, many widely-used metrics to evaluate natural language generation (NLG) systems, such as BERTScore ([Zhang* et al., 2020](#)) or BLEURT ([Sellam et al., 2020](#)), are based on these models. Automatic evaluation is indispensable when language models are published nearly every day. Moreover, large benchmark datasets make a human evaluation of language models infeasible. Therefore, metric scores are among the most important model selection criteria. However, here we see a severe issue when these metrics fail to distinguish between sentences and their negated versions. Especially when considering fact-checking or entailment prediction, an uninterpreted "not" can invalidate the entire output of the model and, thus, reduce the trustworthiness of such systems.

While there have been approaches to improve the negation-awareness in natural language in-



Figure 1: Existing metrics like BLEURT-20 fail to score negated sentences correctly. We propose *NegBLEURT* that overcomes this problem while preserving detection performance on other perturbations.

ference (NLI) or sentiment classification models ([Moore and Barnes, 2021](#)), the task of negation-sensitive evaluation of such systems is lacking behind ([Karpinska et al., 2022](#)). An example of such failure is shown in [Figure 1](#), where a reference sentence is both paraphrased and negated. The well-established BLEURT-20 metric ([Pu et al., 2021](#)) gives a relatively high score of 0.67 to this negated sentence, suggesting that it does not fully capture the negation in the sentence.

To extend negation research to the topic of evaluation, we present a negation-aware version of BLEURT, named *NegBLEURT* ([Figure 1](#)). In addition, we released a negation-aware sentence transformer ([Reimers and Gurevych, 2019](#)) based on an MPNet model ([Song et al., 2020](#)) that extends the application of negation sensitivity to a broader range of tasks. Both models were fine-tuned on a labeled dataset with about 30 thousand sentence pairs in both their negated and paraphrased versions. We publish this dataset and the sentence negator used to create it together with our models. More specifically, our contributions are:

- We open-sourced a rule-based, sentence-level negation tool and released it as a Python package.

- Based on this negator, we built a Compilation of ANnotated, Negation-Oriented Text-pairs (CANNOT). This negation evaluation dataset can be used to fine-tune evaluation metrics for negation awareness or probe their sensitivity.

- We fine-tuned an MPNet model on our negation dataset. This model returns sentence embeddings that are sensitive to negations.

- We published NegBLEURT, a negation-aware version of the BLEURT evaluation metric.

Our models were evaluated on various benchmark datasets showing that they greatly outperform their base model on negated sentences while delivering similar performance on other tasks.

## 2    Related work

This section highlights existing work investigating the negation awareness of the BERT language model and different NLG evaluation metrics. Furthermore, we present approaches to improve this awareness with negation pre-training.

### 2.1    Studies on negation understanding

Ettinger (2020) showed that BERT is insensitive to negations. She designed a completion task where the hypernym description of a word was masked. In addition, a "not" was added to the sentences, resulting in affirmative and negated versions of each sentence. BERT predicted correct hypernyms for both versions, meaning that the model failed to consider the negation indicator. Similar results were achieved by Kassner and Schütze (2020). However, they obtained correct completions when the model was fine-tuned on a negation classification task.

To investigate if BERT-based metrics inherit this lack of negation awareness, Leung et al. (2022) inspected evaluation metrics such as BERTScore (Zhang* et al., 2020) and Sentence-BERT (Reimers and Gurevych, 2019). They used these metrics to calculate the semantic similarity between 20 equivalent and negated sentence pairs. The BERT-based metrics returned high similarity values, indicating they were robust to negations. More large-scale experiments were performed by Karpinska et al. (2022) and Sai et al. (2021) where different metrics, including BERTScore and BLEURT (Sellam et al.,

2020), were evaluated on a collection of sentence and word-level perturbations, including negations and antonyms. Both studies show that most suggested evaluation metrics struggle to detect negations and deviate strongly from human evaluations. Koch et al. (2022) examined the robustness towards these perturbations on a continuous scale by gradually introducing more perturbations to the sentences and, hence, decreasing their quality step by step. While the metrics' scores lowered for other perturbations, the scores for the negated sentences remained relatively high, indicating insensitivity towards negation.

### 2.2    Improving negation awareness

Negation awareness is crucial for the task of natural language inference (NLI), in which models predict if two sentences entail or contradict each other. Hence, multiple datasets with negated samples and models trained on them have been published (Geiger et al., 2020; Helwe et al., 2022). As such, Hosseini et al. (2021) created BERTNOT by training on the negated LAMA dataset with an unlikelihood training objective. Other negation-aware BERT models are NegBERT (Khandelwal and Sawant, 2020) and CueBERT (Truong et al., 2022a), which were trained for the task of negation cue detection and negation scope resolution. Another task that heavily relies on negation awareness is sentiment classification. Moore and Barnes (2021) proposed multi-task learning with a negation speculation auxiliary task to improve the model's performance on negated samples.

While there has been extensive work on negation understanding in NLI and other tasks, we could not find approaches to improve negation awareness for NLG evaluation metrics. This paper tries to close this gap by pre-training metrics on negated sentences.

## 3    Contrastive negation dataset

To make an evaluation metric aware of negations, we need a dataset containing pairs of reference and candidate sentences and a label of how well the candidate fits the reference. To have a balanced dataset, we not only need negated, but also meaning-preserving paraphrases of the reference sentence. As described in the previous section, there exist multiple datasets focusing on negations. However, most of these datasets are either targeted towards specific tasks such as NLI or only contain

negated sentence pairs. Thus, we processed and aggregated the existent datasets producing a Compilation of ANnotated, Negation-Oriented Text-pairs (CANNOT), which addresses and solves these issues. More specifically, our negation-evaluation dataset is based on the following resources:

- *Not another Negation Benchmark* (Truong et al., 2022b): This dataset was published to improve negation awareness in NLI and includes negated sentence pairs. We filtered for samples with the label "contradiction", resulting in 117 negated pairs.

- *Automated Fact-Checking of Claims from Wikipedia* (Sathe et al., 2020): This dataset contains claim-refutation pairs from texts extracted from Wikipedia. The refutation, i.e., the factually incorrect sentence, is often created by negating the claim or replacing one of its words with an antonym. Including more nuanced negations as antonyms and other semantic artifacts diversify the negations in our dataset, making the models trained on it more robust to different negation forms. Nonetheless, many refuted sentences also included further augmentations, such as hallucinations. To only focus on negations, we discarded sentence pairs that had a Jaccard similarity coefficient of less than $0.55$ or differed in length by four or more words. The word splits were obtained with simple white-space tokenization. After the processing, $14,970$ samples were kept.

- *GLUE Diagnostic Dataset* (Wang et al., 2018): Again, this dataset is targeted to NLI and contains changes beyond pure negation. As with the other datasets, we selected only samples labeled as contradiction and dropped pairs with low Jaccard similarity coefficients or large differences in their lengths. This selection resulted in 154 samples.

- *Sentiment-annotated reviews* (Kotzias et al., 2015): This dataset contains online reviews with a strong positive or negative sentiment and, thus, broadens the domains covered by our data. We selected sentences with an auxiliary verb and at most 33 words. Then, we used our rule-based negation tool (see following subsection 3.1) to create negated versions

of the sentences. In total, $2,110$ further samples were collected.

This resulted in a dataset with negated sentence pairs. To extend it with meaning-preserving paraphrases, we used a PEGASUS model pre-trained for this task [1] and created paraphrased versions of each of the references. Finally, the dataset was augmented by adding a swapped version of each pair. This results in a dataset of $68,780$ sentence pairs with equal distribution of negated and equivalent samples. The pre-processed versions of the underlying datasets and our resulting dataset are publicly available on GitHub [2].

### 3.1 Rule-based sentence negator

While previous work used rule-based negation to create negation datasets before, their negators are often not open-source or lost in large repositories with code for the overall goal of the project. Therefore, we publish a lightweight and open-source sentence negation tool as Python module[3] that can be used beyond the scope of this paper.

Our negation tool focuses on verbal negations and supports the addition and deletion of negation cues on a sentence level. The flowchart for the negator is shown in Figure 2, accompanied by example sentences in Table 1. To determine whether a sentence is negated and to distinguish between auxiliary verb forms and common verbs, we first apply the POS tagger provided by the spaCy package (Honnibal et al., 2020). A negated sentence is a sentence where a token in the dependency tree is labeled as "neg" (branches (1) and (2) in Figure 2). We differentiate between the auxiliary "do" and other auxiliary verbs to remove this negation particle. We either entirely remove the negated "do" (e.g., *don't like* → *like*) or remove the negation particle from the auxiliary (e.g., *isn't* / *is not* → *is*). Afterward, the remaining verb is conjugated to match the form of the auxiliary[4].

To negate an affirmative sentence (branches (3)-(5) in Figure 2), we extract the root verb of the dependency tree. If this verb is a full verb and not an auxiliary, we add a negated "do" matching the

---

Figure 2: Flow chart for rule-based sentence negation. The negator can delete negation cues from already negated sentences as well as add them to negate a sentence.

| Branch | Input sentence | Negated sentence |
|--------|----------------|------------------|
| (1) | I *didn't know* what to do. | I *knew* what to do. |
| (2) | I *have never been* to Paris. | I *have been* to Paris. |
| (3) | I *enjoyed* it so much. | I *did not enjoy* it so much. |
| (4) | I *will be* there. | I *won't be* there. |
| (5) | I*'m* very hungry. | I*'m not* very hungry. |

Table 1: Example sentences for different branches in our rule-based negator. Examples (1) and (2) remove a negation from the sentence while examples (3)-(5) add one. The user can decide whether the system should prefer contracted versions like "won't" instead of "will not".

conjugation of the respective verb. If the root is an auxiliary verb, we either negate its first auxiliary child, if any, or otherwise negate the auxiliary itself. The user can decide if the negator should prefer the contracted version, e.g., "don't", or write all words separately, e.g., "do not".

## 3.2 NLG evaluation dataset

We aimed to make evaluation metrics more sensitive towards negations while preserving their ability to detect other errors. Therefore, we added data from the WMT Metrics Shared Task (Bojar et al., 2017) to our dataset. This human-annotated data

focuses on common errors in machine translation outputs and was used to train multiple evaluation metrics before (Sellam et al., 2020). We limited ourselves to the datasets from the years 2015 to 2017 since, upon manual review, the more recent datasets were noisier and contained misannotations. We filtered for samples with a score above $-1$, resulting in $9,264$ labeled samples. Most of the scores range between $0$ and $1$, where a score of $1$ or higher indicates a perfect paraphrase. Our negation data was labeled with a score $0$ for the negated pairs and a score $1$ for the meaning-preserving pairs to match the annotation schema of the WMT

data. Finally, we split both the WMT data and our negation dataset into training, evaluation, and test subsets with a ratio of 80:10:10 and combined the respective subsets. This results in training data with $62,435$ samples and test data with $7,804$ samples, with a 12%-88% distribution of WMT and negation data both.

## 4 Negation-aware models

We publish two different models fine-tuned on our CANNOT data. On the one hand, we fine-tuned a sentence transformer to return negation-aware sentence embeddings. While the cosine similarity of two embeddings can be applied as an evaluation metric, the embedding representations have broader use cases, e.g., for topic modeling (Reimers and Gurevych, 2019). On the other hand, we fine-tuned BLEURT (Sellam et al., 2020), explicitly targeted towards the evaluation task. Both models are published in our GitHub repository[5] and on the Hugging Face Model Hub[6]. Our NegBLEURT can also be utilized within the Hugging Face Evaluate library (Wolf et al., 2020)[7].

### 4.1 Sentence Transformer fine-tuning

Our negation sentence encoder is based on an `all-mpnet-base-v2`[8] model and fine-tuned with the Sentence Transformer library (Reimers and Gurevych, 2019). We trained on our negation training data for one epoch with a batch size of $64$, a learning rate of $2e^{-5}$, and an AdamW optimizer. We utilized a multiple negatives ranking loss to increase the latent distance between correctly paraphrased and negated samples. To create an evaluation metric based on this sentence transformer, embeddings for both the reference and candidate sentence are computed and then scored by their cosine similarity. We call this model the negated MPNet, *NegMPNet*. The cosine-similarity metric based on this model achieves a Spearman correlation of $0.72$ with the ratings in the CANNOT-WMT test set.

---

### 4.2 Negation aware evaluation metric

BLEURT (Sellam et al., 2020) is a reference-based NLG evaluation metric that encodes the references and candidates with a BERT model (Devlin et al., 2019) and predicts a quality score between $0$ and $1$ with a linear regression layer on top of the BERT model. After pre-training with augmented Wikipedia data, BLEURT was trained on WMT data (Bojar et al., 2017). We chose a BLEURT (Sellam et al., 2020) metric as the base for our evaluation metric and selected the `bleurt_tiny` checkpoint, published as test checkpoint on the official GitHub page[9]. This checkpoint is very lightweight with a hidden size of only $128$, instead of $768$ as in standard BERT models. We used the fine-tuning script provided by the authors, and thus, their original hyperparameters and regression (L2) loss. We fine-tuned on the CANNOT training data for 500 steps, resulting in our final *NegBLEURT* checkpoint. This model has a Spearman correlation of $0.65$ with the scores of our test set.

## 5 Evaluation

A common issue with fine-tuning is catastrophic forgetting (Goodfellow et al., 2013), i.e., the models forget their initial knowledge and overfit the new task. Our fine-tuning approach is successful if it improves negation awareness while retaining performance on other tasks, e.g., not corrupting the detection quality of other errors in candidate sentences. We test our models on common embedding and evaluation benchmarks to test our approach and compare their performances against their respective base models.

### 5.1 Massive Text Embedding Benchmark (MTEB)

MTEB (Muennighoff et al., 2023) is an embedding benchmark that evaluates embeddings on multiple tasks such as classification, clustering, and semantic textual similarity (STS). It is one of the most extensive collections of tasks, and thus, we evaluated our NegMPNet on this benchmark. As our work targets English, we only evaluate on the English version of the benchmark. The results averaged per task, and the overall macro average are presented in Table 2, while the performances on the single datasets are provided in Appendix A. Unfortunately, some of the datasets (one reranking

---

Figure 3: Metrics sensitivities to different degrees of perturbation impairment as introduced by Koch et al. (2022). Both our proposed models match the performance of their base models on the word swap, word drop and repetition perturbations but clearly outperform them on the negation detection task.

and seven retrieval) returned errors, and hence, we excluded them from our evaluation. We copied the scores for all-mpnet-base-v2 from the official leaderboard[10]. NegMPNet outperforms its base model in the classification and summarization tasks but shows a decreased performance for clustering, pair classification, and retrieval. When averaging the performances among all tasks, both models perform equally. This benchmark is not targeted towards negation, and, therefore, the results indicate that fine-tuning on our negation data does not harm NegMPNet's general embedding quality.

| Benchmark (num datasets) | Model | |
|---|---|---|
| | all-mpnet-base-v2 | NegMPNet |
| **Average** (60) | 58.78 | 57.16 |
| **Classification** (12) | 65.07 | 70.83 |
| **Clustering** (11) | 43.69 | 38.45 |
| **Pair Classification** (3) | 83.04 | 79.05 |
| **Retrieval** (20) | 43.10 | 36.12 |
| **Reranking** (3) | 68.83 | 68.24 |
| **STS** (10) | 80.28 | 77.58 |
| **Summarization** (1) | 27.49 | 29.84 |

Table 2: Comparison of NegMPNet and its base model on the Massive Text Embedding Benchmark (MTEB). We evaluate on different task categories and macro average the scores.

## 5.2 Improved negation awareness

In this section, we analyze the improved negation awareness beyond the performance on the CANNOT-WMT test set. We selected two NLG evaluation benchmarks that probe negation sensitivity in different metrics and tested NegMPNet (with cosine similarity) and NegBLEURT on them. The results are presented in the following sections.

### 5.2.1 Metrics Comparison benchmark

Koch et al. (2022) probed current, learned metrics for their sensitivity to word swap, repetition, certain word drops, and negations. They gradually increased the level of impairment for all perturbations except negation to measure if the metrics could reflect upon this gradual deterioration. We evaluated our two metrics on their codebase and report the results in Figure 3. NegBLEURT matches the performance of its BLEURT base model (Sellam et al., 2020) and is sensitive to word drops and repetitions but unaware of word swaps. The same holds for NegMPNet, which shows performances similar to its base model for word swap, word drop, and repetition perturbations. In contrast, NegBLEURT and NegMPNet clearly outperform their base models with a mean difference score of up to 0.5 for the negation perturbation. These results demonstrate that our models are aware of negation but do not overfit on them and, thus, preserve their performance on tasks aside pure negation detection.

### 5.2.2 DEMETR benchmark

DEMETR (Karpinska et al., 2022) is a diagnosing benchmark dataset for machine-translation output. It contains reference-candidate pairs with different perturbed versions of the candidates, spanning semantic, syntactic, and morphological errors. These

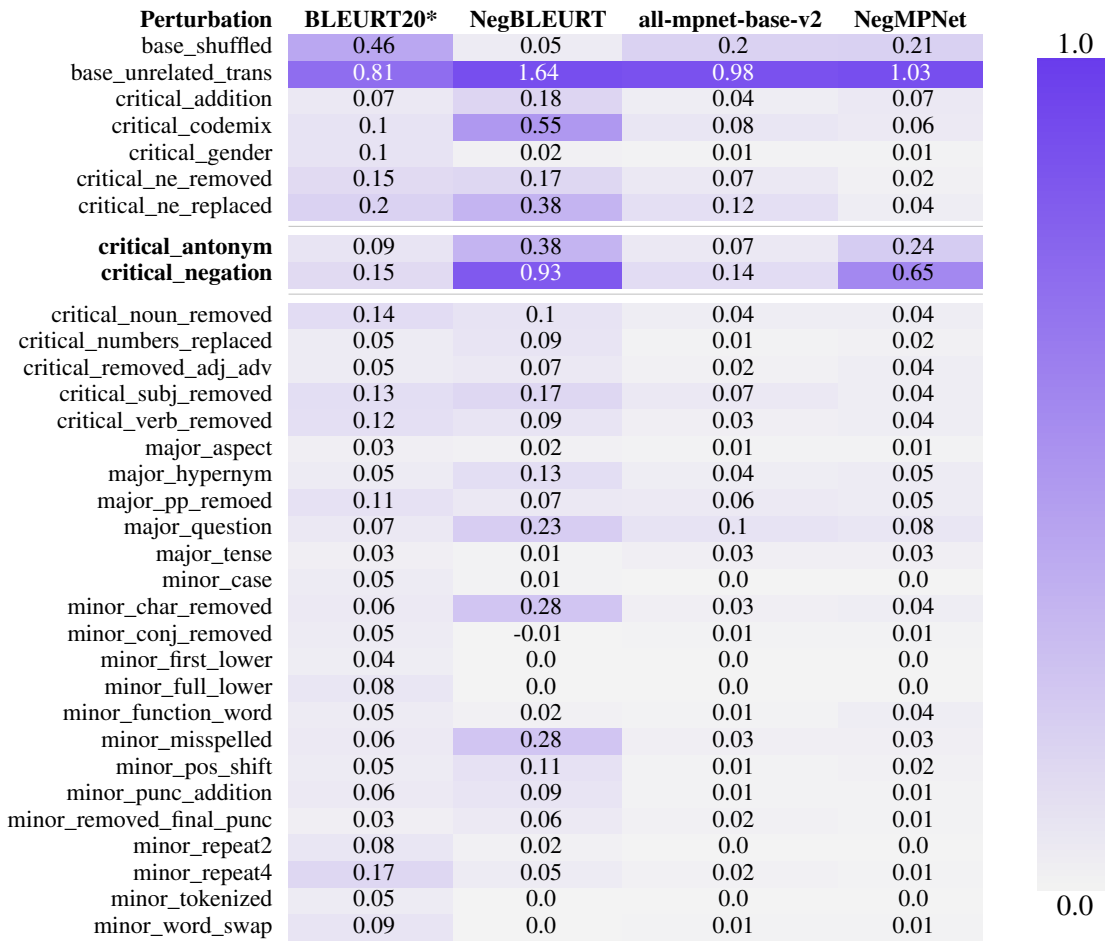| Perturbation | BLEURT20* | NegBLEURT | all-mpnet-base-v2 | NegMPNet |
|---|---|---|---|---|
| base_shuffled | 0.46 | 0.05 | 0.2 | 0.21 |
| base_unrelated_trans | 0.81 | 1.64 | 0.98 | 1.03 |
| critical_addition | 0.07 | 0.18 | 0.04 | 0.07 |
| critical_codemix | 0.1 | 0.55 | 0.08 | 0.06 |
| critical_gender | 0.1 | 0.02 | 0.01 | 0.01 |
| critical_ne_removed | 0.15 | 0.17 | 0.07 | 0.02 |
| critical_ne_replaced | 0.2 | 0.38 | 0.12 | 0.04 |
| **critical_antonym** | 0.09 | 0.38 | 0.07 | 0.24 |
| **critical_negation** | 0.15 | 0.93 | 0.14 | 0.65 |
| critical_noun_removed | 0.14 | 0.1 | 0.04 | 0.04 |
| critical_numbers_replaced | 0.05 | 0.09 | 0.01 | 0.02 |
| critical_removed_adj_adv | 0.05 | 0.07 | 0.02 | 0.04 |
| critical_subj_removed | 0.13 | 0.17 | 0.07 | 0.04 |
| critical_verb_removed | 0.12 | 0.09 | 0.03 | 0.04 |
| major_aspect | 0.03 | 0.02 | 0.01 | 0.01 |
| major_hypernym | 0.05 | 0.13 | 0.04 | 0.05 |
| major_pp_remoed | 0.11 | 0.07 | 0.06 | 0.05 |
| major_question | 0.07 | 0.23 | 0.1 | 0.08 |
| major_tense | 0.03 | 0.01 | 0.03 | 0.03 |
| minor_case | 0.05 | 0.01 | 0.0 | 0.0 |
| minor_char_removed | 0.06 | 0.28 | 0.03 | 0.04 |
| minor_conj_removed | 0.05 | -0.01 | 0.01 | 0.01 |
| minor_first_lower | 0.04 | 0.0 | 0.0 | 0.0 |
| minor_full_lower | 0.08 | 0.0 | 0.0 | 0.0 |
| minor_function_word | 0.05 | 0.02 | 0.01 | 0.04 |
| minor_misspelled | 0.06 | 0.28 | 0.03 | 0.03 |
| minor_pos_shift | 0.05 | 0.11 | 0.01 | 0.02 |
| minor_punc_addition | 0.06 | 0.09 | 0.01 | 0.01 |
| minor_removed_final_punc | 0.03 | 0.06 | 0.02 | 0.01 |
| minor_repeat2 | 0.08 | 0.02 | 0.0 | 0.0 |
| minor_repeat4 | 0.17 | 0.05 | 0.02 | 0.01 |
| minor_tokenized | 0.05 | 0.0 | 0.0 | 0.0 |
| minor_word_swap | 0.09 | 0.0 | 0.01 | 0.01 |

1.0

0.0

Figure 4: Sensitivity scores of different NLG metrics on the DEMETR benchmark. The values represent ratios as introduced in Karpinska et al. (2022). A higher value denotes a higher sensitivity and is marked in a darker color. Both our models clearly have a higher sensitivity towards negations than their base versions.
* copied from the original paper

errors are categorized by their severity, with categories being critical, major, and minor. The authors measured the sensitivity of the metric to a specific perturbation by predicting the metrics' scores for the reference-candidate and reference-perturbed candidate pairs and calculating the weighted difference between the scores. In their original work, BERTScore achieved the best negation (0.21) and antonym (0.15) detection scores. However, these values still need to catch up to detection scores of other perturbations.

Figure 4 shows the sensitivity scores of our metrics compared to their base versions. Sensitivity scores for other metrics such as BERTScore or COMET are presented on the original DEMETR paper by Karpinska et al. (2022). NegBLEURT clearly outperforms all metrics on the critical negation and antonym perturbations while preserving or even improving the detection rate on other perturbations compared to BLEURT20 (Pu et al., 2021).

Both all-mpnet-base-v2 and NegMPNet show no sensitivity for most of the perturbations, indicating that these sentence transformers were not trained for the task of NLG evaluation. Nevertheless, NegMPNet shows a competitive detection rate on negations and even antonyms. An increased sensitivity towards antonyms indicates that our fine-tuning approach yields embeddings that distinguish between negated and affirmative sentences beyond the presence of the word "not". Although our dataset mainly focuses on simple verbal negations, fine-tuning on it teaches the models to capture antonym-related nuances better. This suggests that our models do not simply learn the artefacts in our CANNOT dataset by hard but can distinguish between different types of contradictions.

## 6  Ablation study

Our CANNOT dataset consists of a diverse collection of datasets, as presented in (section 3).

| DEMETR Perturbation | Neg-BLEURT | w/o Not another Negation Benchmark (Truong et al., 2022b) | w/o Automated Fact-Checking of Claims from Wikipedia (Sathe et al., 2020) | w/o GLUE Diagnostic Dataset (Wang et al., 2018) | w/o Sentiment-annotated reviews (Kotzias et al., 2015) with rule-based negations |
|---|---|---|---|---|---|
| numbers_replaced | 0.09 | 0.06 | 0.04 | 0.11 | 0.05 |
| gender | 0.02 | 0.03 | 0.02 | 0.02 | 0.03 |
| shuffled | 0.05 | 0.04 | 0.02 | 0.06 | 0.04 |
| adj_adv_removed | 0.07 | 0.06 | 0.04 | 0.09 | 0.09 |
| verb_removed | 0.09 | 0.1 | 0.06 | 0.01 | 0.11 |
| noun_removed | 0.1 | 0.3 | 0.07 | 0.21 | 0.15 |
| subj_removed | 0.17 | 0.09 | 0.1 | 0.2 | 0.09 |
| ne_removed | 0.17 | 0.14 | 0.1 | 0.17 | 0.15 |
| codemix | 0.55 | 0.33 | 0.31 | 0.52 | 0.5 |
| addition | 0.18 | 0.2 | 0.05 | 0.18 | 0.17 |
| antonym | 0.38 | 0.29 | 0.1 | 0.41 | 0.39 |
| negation | 0.93 | 0.74 | 0.35 | 1.02 | 0.82 |
| ne_replaced | 0.38 | 0.26 | 0.25 | 0.4 | 0.4 |
| **Number of removed samples** | 0 | 281 | 53.747 | 399 | 7.475 |

Table 3: Ablation study of CANNOT subsets. Each data source in the CANNOT-WMT data was removed individually, and the resulting BLEURT checkpoint evaluated on the DEMETR perturbations marked as critical (Karpinska et al., 2022). Cells with an orange color indicate a decreased performance compared to NegBLEURT, while green indicates an improvement. The darker the color, the larger the difference.

Here, we perform an ablation study to measure the impact of the individual subsets on the model's overall performance. We individually removed each subset from the CANNOT-WMT training and evaluation data, and fine-tuned a new BLEURT checkpoint based on the remaining subsets and the WMT data. The setup for training was the same as for NegBLEURT, which means we fine-tuned the BLEURT-tiny checkpoint for 500 steps using the fine-tuning script from the BLEURT GitHub page. We did not modify the CANNOT-WMT test split, and thus, the number of samples in the subsets deviate from the numbers in section 3.

To compare the impacts of the different datasets, we evaluated the fine-tuned models on all critical perturbations in the DEMETR evaluation benchmark (Karpinska et al., 2022). The model's sensitivities towards specific perturbations are presented in Table 3. The Wiki-Factcheck dataset (Sathe et al., 2020) is by far the biggest subset. As expected, removing it from the training data results in sensitivity drops for nearly all perturbations, espe-

cially for antonyms and negations. Removing the Truong et al. (2022b) dataset from the training data results in a substantial sensitivity loss towards the negation perturbation. This is especially remarkable as the dataset is very small, with less than 300 samples. The BLEURT checkpoint fine-tuned without the GLUE subset shows an increased sensitivity towards negations, indicating that the dataset contains some noise introduced by the selection of contradicting samples. These samples may cover contradictions beyond pure negations that decrease NegBLEURT's performance.

## 7 Conclusion

In this work, we created a sentence negation tool that we made available to the research community as a Python package. In addition, we released CANNOT, a data collection for negation detection that can be used to improve negation awareness of language models. We leveraged this dataset to fine-tune a sentence transformer and an NLG evaluation metric. Both models show a strong negation

detection ability while preserving task-specific performances compared to their base models. Considering the many papers that pointed out the negation weaknesses of model language models, our work is an important step towards negation awareness.

In the future, we will extend our negator to support more advanced negations beyond the verb level and make our negation dataset multilingual so that multilingual NLG evaluation metrics can be improved as well.

## 8 Limitations

Our rule-based negation system works on a verb level and fails for cases that do not match our defined sentence structure. In addition, there are special cases like the sentence "She's determined" that could use both the verbs "is" and "has". When removing contractions, the negator has to select one of the verbs and may, hence, change the meaning of the input sentence beyond the pure negation. Moreover, in sentences like "I have not yet been there.", the adverb "yet" must be removed or replaced by "already", both of which still need to be added to our tool.

We evaluated our negation-aware models on two evaluation metric benchmark datasets. We would have wished to extend this evaluation to further benchmarks, but unfortunately, the lack of published code or datasets, and insufficient code documentation prevented us from doing so.

While NegMPNet achieves remarkable negation detection scores on the considered metric evaluation benchmarks, we must admit that it fails with most of the other perturbations. The all-mpnet-base-v2 model was initially trained as a sentence transformer and not as an evaluation metric, and the small percentage of WMT data in our dataset is insufficient to train it to be such. Therefore, NegMPNet can produce negation-sensitive sentence embeddings but needs further work to be applied as an evaluation metric directly.

## Ethical Statement

As stated in section 1, trustworthy automatic evaluation metrics are indispensable for selecting and deploying large language models. Metrics that capture negations and reduce the overall score for models that mix up negated and original sentences are, therefore, an important step to increase trust in the metrics themselves, but also in the evaluated models. Moreover, models can be trained to improve

negation sensitivity with metrics that detect negation insensitivity, as well as the CANNOT dataset. Therefore, we do not see ethical concerns with our negation-aware metrics or datasets.

However, our negation tool can add or remove a negation to any input sentence. If applied to sentences from the Internet, such as news articles or Twitter posts, it can easily alter the information provided. The negated and original versions still look very similar, and thus, people might oversee the missing or added negation cues when comparing the provided information with other sources. Consequently, we are aware that our negator may be used in a malicious way to spread misinformation. Nevertheless, negation-aware sentence embeddings and evaluation metrics could again detect such modifications. We believe that the benefits of an open-source tool for researchers, as well as the simplified dataset creation it enables, outweigh the drawbacks of potential misuse.

**Supplementary Materials Availability Statement:** All material used in this paper is available to the research community. The sentence negation tool is published as a Python package and in a GitHub repository. The dataset and source code for fine-tuning on this data is also open-sourced on GitHub and Hugging Face. The checkpoints of our models are available on the Hugging Face Model Hub. The links to the individual resources are referenced in their respective paper sections.

## References

Ondřej Bojar, Yvette Graham, and Amir Kamran. 2017. Results of the WMT17 metrics shared task. In *Proceedings of the Second Conference on Machine Translation*, pages 489–513, Copenhagen, Denmark. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Allyson Ettinger. 2020. What BERT is not: Lessons from a new suite of psycholinguistic diagnostics for language models. *Transactions of the Association for Computational Linguistics*, 8:34–48.

Atticus Geiger, Kyle Richardson, and Christopher Potts. 2020. Neural natural language inference models

partially embed theories of lexical entailment and negation. In *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 163–173, Online. Association for Computational Linguistics.

Ian J. Goodfellow, Mehdi Mirza, Da Xiao, Aaron Courville, and Yoshua Bengio. 2013. An empirical investigation of catastrophic forgetting in gradient-based neural networks.

Chadi Helwe, Simon Coumes, Chloé Clavel, and Fabian Suchanek. 2022. TINA: Textual inference with negation augmentation. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 4086–4099, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. spaCy: Industrial-strength Natural Language Processing in Python.

Arian Hosseini, Siva Reddy, Dzmitry Bahdanau, R Devon Hjelm, Alessandro Sordoni, and Aaron Courville. 2021. Understanding by understanding not: Modeling negation in language models. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1301–1312, Online. Association for Computational Linguistics.

Marzena Karpinska, Nishant Raj, Katherine Thai, Yixiao Song, Ankita Gupta, and Mohit Iyyer. 2022. DEMETR: Diagnosing evaluation metrics for translation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9540–9561, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Nora Kassner and Hinrich Schütze. 2020. Negated and misprimed probes for pretrained language models: Birds can talk, but cannot fly. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7811–7818, Online. Association for Computational Linguistics.

Aditya Khandelwal and Suraj Sawant. 2020. NegBERT: A transfer learning approach for negation detection and scope resolution. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 5739–5748, Marseille, France. European Language Resources Association.

Philipp Koch, Matthias Aßenmacher, and Christian Heumann. 2022. Pre-trained language models evaluating themselves - a comparative study. In *Proceedings of the Third Workshop on Insights from Negative Results in NLP*, pages 180–187, Dublin, Ireland. Association for Computational Linguistics.

Dimitrios Kotzias, Misha Denil, Nando de Freitas, and Padhraic Smyth. 2015. From group to individual labels using deep features. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '15, page

597–606, New York, NY, USA. Association for Computing Machinery.

Wai Ching Leung, Shira Wein, and Nathan Schneider. 2022. Semantic similarity as a window into vector- and graph-based metrics. In *Proceedings of the 2nd Workshop on Natural Language Generation, Evaluation, and Metrics (GEM)*, pages 106–115, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Andrew Moore and Jeremy Barnes. 2021. Multi-task learning of negation and speculation for targeted sentiment classification. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2838–2869, Online. Association for Computational Linguistics.

Niklas Muennighoff, Nouamane Tazi, Loic Magne, and Nils Reimers. 2023. MTEB: Massive text embedding benchmark. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2014–2037, Dubrovnik, Croatia. Association for Computational Linguistics.

Amy Pu, Hyung Won Chung, Ankur Parikh, Sebastian Gehrmann, and Thibault Sellam. 2021. Learning compact metrics for MT. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 751–762, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERT-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.

Ananya B. Sai, Tanay Dixit, Dev Yashpal Sheth, Sreyas Mohan, and Mitesh M. Khapra. 2021. Perturbation CheckLists for evaluating NLG evaluation metrics. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7219–7234, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Aalok Sathe, Salar Ather, Tuan Manh Le, Nathan Perry, and Joonsuk Park. 2020. Automated fact-checking of claims from Wikipedia. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 6874–6882, Marseille, France. European Language Resources Association.

Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online. Association for Computational Linguistics.

Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2020. Mpnet: Masked and permuted pre-training for language understanding. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS'20, Red Hook, NY, USA. Curran Associates Inc.

Thinh Truong, Timothy Baldwin, Trevor Cohn, and Karin Verspoor. 2022a. Improving negation detection with negation-focused pre-training. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4188–4193, Seattle, United States. Association for Computational Linguistics.

Thinh Hung Truong, Yulia Otmakhova, Timothy Baldwin, Trevor Cohn, Jey Han Lau, and Karin Verspoor. 2022b. Not another negation benchmark: The NaN-NLI test suite for sub-clausal negation. In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 883–894, Online only. Association for Computational Linguistics.

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels, Belgium. Association for Computational Linguistics.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

## A    MTEB full results

| Task | Benchmark | Model | |
|------|-----------|-------|---|
| | | all-mpnet-base-v2* | NegMPNet |
| **Classification** (12 datasets) | AmazonCounterfactualClassification | 65.27 | 73.96 |
| | AmazonPolarityClassification | 67.13 | 86.1 |
| | AmazonReviewsClassification | 31.92 | 41.85 |
| | Banking77Classification | 81.86 | 84.23 |
| | EmotionClassification | 39.72 | 45.98 |
| | ImdbClassification | 70.72 | 68.4 |
| | MTOPDomainClassification | 92.08 | 93.38 |
| | MTOPIntentClassification | 70.21 | 78.45 |
| | MassiveIntentClassification | 69.57 | 74.38 |
| | MassiveScenarioClassification | 76.01 | 78.12 |
| | ToxicConversationsClassification | 60.86 | 66.15 |
| | TweetSentimentExtractionClassification | 55.46 | 58.99 |
| **Retrieval** (20 datasets) | ArguAna | 46.52 | 19.51 |
| | CQADupstackAndroidRetrieval | 56.49 | 53.46 |
| | CQADupstackEnglishRetrieval | 52.29 | 49.36 |
| | CQADupstackGamingRetrieval | 60.03 | 52.16 |
| | CQADupstackGisRetrieval | 44.27 | 40.42 |
| | CQADupstackMathematicaRetrieval | 34.21 | 31.8 |
| | CQADupstackPhysicsRetrieval | 50.97 | 44.38 |
| | CQADupstackProgrammersRetrieval | 44.17 | 41.24 |
| | CQADupstackStatsRetrieval | 38.15 | 36.38 |
| | CQADupstackTexRetrieval | 33.35 | 30.7 |
| | CQADupstackUnixRetrieval | 45.41 | 42.18 |
| | CQADupstackWebmastersRetrieval | 44.24 | 42.92 |
| | CQADupstackWordpressRetrieval | 35.94 | 33.39 |
| | DBPedia | 32.09 | 23.08 |
| | FiQA2018 | 49.96 | 26.89 |
| | NFCorpus | 33.29 | 27.97 |
| | SCIDOCS | 23.76 | 20.1 |
| | SciFact | 65.57 | 30.83 |
| | TRECCOVID | 51.33 | 58.9 |
| | Touche2020 | 19.93 | 16.72 |
| **Clustering** (11 datasets) | ArxivClusteringP2P | 48.38 | 42.53 |
| | ArxivClusteringS2S | 39.72 | 37.92 |
| | BiorxivClusteringP2P | 39.62 | 33.7 |
| | BiorxivClusteringS2S | 35.02 | 33.45 |
| | MedrxivClusteringP2P | 35.58 | 29.97 |
| | MedrxivClusteringS2S | 32.87 | 31.48 |
| | RedditClustering | 54.82 | 44.31 |
| | RedditClusteringP2P | 56.77 | 45.43 |
| | StackExchangeClustering | 53.8 | 49.4 |
| | StackExchangeClusteringP2P | 34.28 | 30.14 |
| | TwentyNewsgroupsClustering | 49.74 | 44.7 |
| **Reranking** (3 datasets) | AskUbuntuDupQuestions | 65.85 | 65.11 |
| | SciDocsRR | 88.65 | 87.75 |

| Task | Benchmark | Model | |
| --- | --- | --- | --- |
| | | all-mpnet-base-v2* | NegMPNet |
| | StackOverflowDupQuestions | 51.98 | 51.87 |
| **STS** (10 datasets) | BIOSSES | 80.43 | 64.45 |
| | SICK-R | 80.59 | 76.71 |
| | STS12 | 72.63 | 71.23 |
| | STS13 | 83.48 | 84.62 |
| | STS14 | 78 | 79.39 |
| | STS15 | 85.66 | 84.7 |
| | STS16 | 80.03 | 82.17 |
| | STS17 | 90.6 | 90.77 |
| | STS22 | 67.95 | 57.63 |
| | STSBenchmark | 83.42 | 84.18 |
| **PairClassification** (3 datasets) | SprintDuplicateQuestions | 90.15 | 77.69 |
| | TwitterSemEval2015 | 73.85 | 75.98 |
| | TwitterURLCorpus | 85.11 | 83.48 |
| **Summarization** (1 datasets) | SummEval | 27.49 | 29.84 |

Table 4: Detailed performance on MTEB by task and Benchmark dataset.
* copied form the official leaderboard at `https://huggingface.co/spaces/mteb/leaderboard`, as of 17.05.2023.

# Guidance in Radiology Report Summarization:
## An Empirical Evaluation and Error Analysis

**Jan Trienes**[†]    **Paul Youssef**[‡]    **Jörg Schlötterer**[†¶‡]    **Christin Seifert**[†‡]

[†]University of Duisburg-Essen, University Hospital Essen
[¶]University of Mannheim    [‡]University of Marburg

jan.trienes@uni-due.de
{paul.youssef, joerg.schloetterer, christin.seifert}@uni-marburg.de

## Abstract

Automatically summarizing radiology reports into a concise impression can reduce the manual burden of clinicians and improve the consistency of reporting. Previous work aimed to enhance content selection and factuality through guided abstractive summarization. However, two key issues persist. First, current methods heavily rely on domain-specific resources to extract the guidance signal, limiting their transferability to domains and languages where those resources are unavailable. Second, while automatic metrics like ROUGE show progress, we lack a good understanding of the errors and failure modes in this task. To bridge these gaps, we first propose a domain-agnostic guidance signal in form of variable-length extractive summaries. Our empirical results on two English benchmarks demonstrate that this guidance signal improves upon unguided summarization while being competitive with domain-specific methods. Additionally, we run an expert evaluation of four systems according to a taxonomy of 11 fine-grained errors. We find that the most pressing differences between automatic summaries and those of radiologists relate to content selection including omissions (up to 52%) and additions (up to 57%). We hypothesize that latent reporting factors and corpus-level inconsistencies may limit models to reliably learn content selection from the available data, presenting promising directions for future work.

## 1 Introduction

The radiology report is an important tool for radiologists to communicate examination results with other clinicians. Typically, these reports contain three sections: the background section describing the exam and patient context, the findings section providing a detailed description of observations, and the impression section, which concisely summarizes the key findings (Kahn et al., 2009). In the clinical process, the impression is of high importance as it informs further treatments. However,



Figure 1: Example radiology report. We guide abstractive summarization with extractive summaries. We propose to adapt the length of the guidance signal to each report rather than using a fixed setting across all reports which helps to accommodate varying target lengths.

writing the impression can be time-consuming and error-prone, which is why automatic text summarization systems can substantially improve the quality of clinical reporting (Gershanik et al., 2011).

From a summarization perspective, this task involves both an *extractive* component, where important findings are copied verbatim into the summary, and an *abstractive* component, forming those findings into a concise conclusion taking into account the full report (example in Figure 1). Although abstractive methods generate fluent and relevant summaries, they are prone to hallucinations and their output is difficult to control (Maynez et al., 2020; Kryściński et al., 2020; Huang et al., 2020). Therefore, current methods for radiology report summarization employ *guided text summarization* to control the summary content through carefully selected guidance signals such as salient ontology terms (Sotudeh et al., 2020), facts (Zhang et al., 2020b), and clinical entities (Hu et al., 2022).

While summary quality has improved steadily,

**A** Guided summarization framework

Candidate Impression (ŷ)

Decoder

Document Encoder | Guidance Encoder

Findings (**x**) | Guidance (**g**)

**B** Manual error analysis of four methods

BertAbs
GSum
WGSum
WGSum+CL

Finding/interpretation | Comparison | Ref. to prior report
Communication/followup | Contradicting finding

**Reference Impression** ①
Right lower lobe opacity, possibly atelectasis, with associated moderate sized effusion.

**Candidate Impression** ②
Persistent right lower lobe opacity with associated effusion, mildly progressed from the preceding radiograph.

☐ Incorrect *location* of a finding? ③    Any other error? ⑤
☐ Incorrect *severity* of a finding? ④    Please describe...
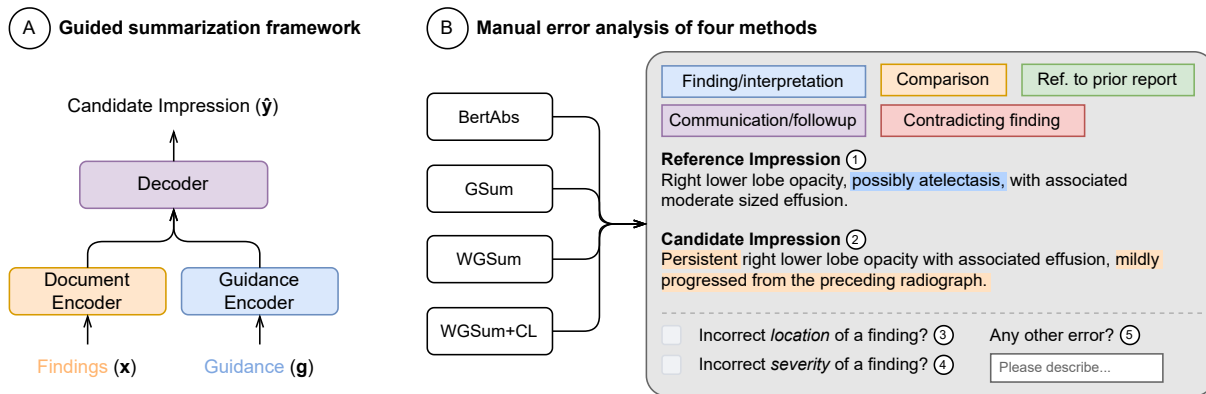
Figure 2: Paper overview. <u>Left:</u> We evaluate guided methods, where the decoder attends to both the input and a pre-selected guidance signal. <u>Right:</u> We task domain experts to identify errors in summaries of unguided and guided methods. Annotation is done on a span-level for omissions from the reference (①) and additions to the candidate (②), and on an instance-level when both texts report the same finding but with incorrect semantics (③,④). Annotators can flag any other error in free-form (⑤).

two key issues have received little attention. First, the success of current methods heavily relies on the availability and quality of the domain-specific guidance extractors (i.e., ontologies, fact extractors, clinical entity taggers). As these resources are expensive to develop, and as they are only available for a few languages and clinical domains, it is of great interest to investigate to what extent we can use domain-agnostic guidance signals to make guided radiology report summarization methods more easily adoptable. Second, while we do see improvements in automatic metrics and human assessments of coarse quality criteria such as fluency, correctness and completeness, we lack a good understanding and quantification of the errors and failure modes of current methods. We address the two issues as follows (overview in Figure 2).

**Contribution 1: a variable-length extractive guidance signal (Sections 3 and 4).** Motivated by the observation that summaries have a large extractive component, we investigate *extractive summaries* as guidance (Dou et al., 2021). We identify that the *length of the extractive summaries* is critical for the effectiveness of this guidance signal. Intuitively, generating longer summaries requires more guidance than shorter ones. Therefore, we outline two approaches to adapt the guidance length to each report: (i) a classifier that predicts a suitable length, and (ii) a threshold-based method. This variable-length guidance signal improves upon unguided summarization, is competitive with recent domain-guided methods, and cheap to adopt as it does not require any domain-specific resources.

**Contribution 2: an error analysis (Section 5).** We conduct an error analysis of unguided, guided and domain-guided methods to identify avenues for improvements of all methods. We ask domain experts to identify errors in the outputs of four systems and to characterize them along 11 fine-grained categories. This analysis uncovers three aspects: (1) despite guidance, a significant portion of candidates shows problems in content selection, (2) some content selection decisions are likely only explained by latent factors, (3) there are some dataset-level issues in MIMIC-CXR, including duplicate findings with different impressions, and impression segments without any grounding in the reports.

We make experiment code, full pre-processing pipeline for the datasets, and 1,200 expert assessments of model outputs publicly available.[1]

## 2 Related Work

**Guided text summarization.** To address shortcomings in content selection and factuality of neural abstractive summarization methods (Rush et al., 2015; Nallapati et al., 2016), guided methods aim to control the content of summaries through carefully selected guidance signals such as keywords (Li et al., 2018), sentences (Chen and Bansal, 2018), entities (Fan et al., 2018; Narayan et al., 2021), templates (Cao et al., 2018) and prompts (He et al., 2022).

In the radiology domain, Zhang et al. (2018) proposed to guide generation with the background section of reports using a dual-encoder model. With a

---

[1]github.com/jantrienes/inlg2023-radsum

177

similar architecture, subsequent work explored the use of salient ontology terms (MacAvaney et al., 2019; Sotudeh et al., 2020). Besides text-based guidance, Hu et al. (2021, 2022) propose a graph-guided decoder which attends both to the report text and to a word-graph of clinical entities. In contrast, we explore extractive summaries as guidance signal (Liu and Lapata, 2019; Dou et al., 2021). Crucially, this guidance signal can be extracted without any domain-specific resources such as ontologies and clinical entity taggers. To this end, our method is similar to the approach by Zhu et al. (2023), which guides summarization with reference summaries from the training set, selected at random or by an oracle. Similar to our approach, this guidance signal can vary in length.

**Alternatives to guided summarization.** Several parallel research lines aim to enhance radiology report summarization with a different methodological focus. First, several studies optimize factual consistency through reinforcement learning (Zhang et al., 2020b; Delbrouck et al., 2022) or post-hoc reranking (Xie et al., 2023). Second, Karn et al. (2022) devise an extract-then-abstract pipeline with multi-agent reinforcement learning. Last, recent work explores domain-adaptation techniques for pre-trained language models to better accommodate radiology reports (Cai et al., 2021; Van Veen et al., 2023). Our work is orthogonal to these efforts and future work could investigate how to combine them with guided summarization.

**Evaluation of radiology report summarization.** Evaluating text summarization systems is a long standing issue. Since automatic metrics have a limited correlation with human judgment (Fabbri et al., 2021), manual evaluation is still regarded as the gold standard. For the task of radiology report summarization, most manual evaluations focus on coarse criteria such as accuracy, completeness, conciseness, and readability (Zhang et al., 2018; Hu et al., 2022; Cai et al., 2021). Yet, these evaluations only provide limited insights into directions for improvement. To support the interpretation of automatic and manual evaluations, and to understand the pitfalls of current methods, we conduct an error analysis (van Miltenburg et al., 2021). In this line of work, Yu et al. (2022) evaluated the ability of automatic metrics to capture six fine-grained errors of radiograph-to-impression models. We extend this taxonomy in our error analysis.

## 3 Method

We formulate the task of summarizing radiology reports as follows. Given the findings section of a report, represented as a sequence of tokens $x = (x_1, x_2, \ldots, x_N)$, the goal is to generate an impression $y = (y_1, y_2, \ldots, y_M)$ that accurately summarizes the most significant findings. The guided summarization framework extends this problem setting with an additional input signal $g = (g_1, g_2, \ldots, g_L)$ which aims to improve the quality of generated summaries by indicating salient information in $x$.

### 3.1 Model and Extractive Guidance

As a concrete implementation of the guided summarization framework, we adopt GSum (Dou et al., 2021). This sequence-to-sequence model extends a transformer-based architecture for abstractive text summarization (Liu and Lapata, 2019) with an additional encoder for the guidance signal $g$. To create a guidance-aware representation of the input, the decoder first attends to the encoded representation of $g$, and afterwards to the whole input document $x$ using cross-attention (Vaswani et al., 2017). The authors demonstrate that GSum is effective at controlling the content of summaries, leading to good results on several non-medical datasets.

**Extractive guidance.** While $g$ can take any form, Dou et al. (2021) found the output of an extractive summarization to be highly effective. Intuitively, this guidance signal informs the model about which input sentences should be highlighted in a summary. An important implementation detail of GSum is the mechanism to obtain the extractive sentences. Dou et al. (2021) distinguish between the **oracle** setting and the **automatic** setting. In the oracle setting, the guidance sentences are greedily picked from $x$ such that they maximize ROUGE with respect to $y$ (Nallapati et al., 2017). In the automatic setting, this oracle is approximated by an extractive summarization method (BertExt, Liu and Lapata, 2019). The training labels for BertExt are derived using the same greedy matching, thus BertExt can be considered an approximation of the oracle guidance. Selecting the guidance signal from BertExt follows a top-$k$ strategy: scoring all sentences for relevance and selecting the highest scoring sentences until a predefined length threshold is reached (Nallapati et al., 2017; Liu and Lapata, 2019). Following Dou et al. (2021), the oracle signal is used during training of GSum. During inference, we state explicitly whether we use the oracle or automatic variant.

## 3.2 Variable-length Extractive Guidance

We empirically find that extracting a fixed-length summary with the top-$k$ approach has a negative impact on the effectiveness of GSum (Section 4.2). To address this problem, we propose two methods to select a variable-length extractive guidance signal from BertExt. Formally, for a given document $\boldsymbol{x}$ and its sequence of sentences $(s_1, \dots, s_N)$, with $s_i$ being the $i$-th sentence in $\boldsymbol{x}$, these methods have to select $L < N$ sentences as guidance $\boldsymbol{g}$.

**Method 1: predicting oracle length.** As described in Section 3.1, BertExt is trained to assign a label $y \in \{0, 1\}$ to each sentence $s_i$. The predicted probability $p(y = 1|s_i)$ indicates if $s_i$ should be included in the summary. The ground-truth labels are derived from an extractive oracle $f_{\text{oracle}}(\boldsymbol{x}, \boldsymbol{y})$ which greedily selects a subset of sentences of length $[0, 3]$ that maximizes ROUGE against the gold summary $\boldsymbol{y}$ (Liu and Lapata, 2019). Instead of taking a fixed number of sentences with highest probability (top-$k$), we train a sequence-classification model to predict the length of the extractive oracle $f_{\text{approx}}(\boldsymbol{x}) = L \in [0, 3]$, and select the top-$L$ sentences as guidance signal.

**Method 2: threshold calibration.** Instead of considering the full ranked list of sentences, this method constrains selection with a threshold-based approach inspired by Jia et al. (2021). Recall that $p(y = 1|s_i)$ denotes the probability that BertExt assigns to the positive class. We then select the set of sentences that exceed a probability threshold $T$ as guidance signal:

$$\boldsymbol{g} = \{s_i \in \boldsymbol{x} | p(y = 1|s_i) \geq T\}.$$

We optimize $T \in [0, 1]$ on a validation set to maximize ROUGE-1.

## 4 Technical Evaluation

> **RQ1.** To what extent are extractive summaries an effective guidance signal for radiology report summarization?
>
> **RQ2.** How does adapting the extractive guidance length to each report impact the overall quality of summaries?

### 4.1 Experimental Setup

**Datasets.** We use two public datasets of English chest x-ray reports: **MIMIC-CXR** (Johnson et al., 2019a) and **OpenI** (Demner-Fushman et al., 2015). Consistent with prior work (Zhang et al., 2018; Sotudeh et al., 2020; Hu et al., 2022), we retain reports

| Aspect | MIMIC-CXR | OpenI |
|---|---|---|
| Reports | 122,500 / 963 / 1,598 | 2,342 / 334 / 670 |
| Avg. $|\boldsymbol{x}|_t$ | $56 \pm 25.2$ | $37 \pm 16.4$ |
| Avg. $|\boldsymbol{x}|_s$ | $5.5 \pm 1.9$ | $4.6 \pm 1.6$ |
| Avg. $|\boldsymbol{y}|_t$ | $15 \pm 13.5$ | $8 \pm 8.1$ |
| Avg. $|\boldsymbol{y}|_s$ | $1.6 \pm 0.9$ | $1.4 \pm 0.8$ |
| Novelty | 73.4% | 86.8% |
| Compression | 73.8% | 76.1% |

Table 1: Statistics of the benchmark datasets, including the number of reports (train/valid/test), length and standard deviation in tokens/sentences ($|\cdot|_t$ and $|\cdot|_s$), novelty as average percentage of bigrams in impression $\boldsymbol{y}$, but not in findings $\boldsymbol{x}$, and average compression ($\frac{|\boldsymbol{y}|_t}{|\boldsymbol{x}|_t}$).

with exactly one findings and one impression section, where both have an acceptable length ($\geq 10$ tokens in findings, $\geq 2$ tokens in impression), and we discard the background section.[2] Following Hu et al. (2022), we use the official training, validation and test splits of MIMIC-CXR and a random split with a 70/10/20 ratio for OpenI. We use SPACY for tokenization and NLTK for sentence segmentation.[3] Table 1 reports the dataset statistics.

**Baselines.** We compare with three categories of baselines: (1) unguided methods, (2) vanilla GSum with fixed-length extractive guidance (Dou et al., 2021), and (3) domain-specific guided methods. Regarding unguided methods, we use **Oracle-Ext** (Nallapati et al., 2017) which greedily selects sentences from the findings that maximize ROUGE with respect to the impression. Furthermore, we use **BertExt** and **BertAbs** (Liu and Lapata, 2019) which are extractive/abstractive transformer-based models initialized with pre-trained BERT (Devlin et al., 2019). Regarding domain-specific methods, we compare with **WGSum** (Hu et al., 2021) which employs a graph-guided decoder to attend to a graph of clinical entities extracted with Stanza (Zhang et al., 2021), and with **WGSum+CL** (Hu et al., 2022) which refines this guidance signal through contrastive learning.

**Automatic evaluation metrics.** We evaluate the quality of generated impressions with ROUGE $F_1$ (Lin, 2004) to measure unigram and bigram overlap as a proxy for relevance (R-1, R-2) and the longest common subsequence as a proxy for fluency (R-L). In addition, we report BERTScore as

---

[2]To compare the relative utility of guidance signals, including the background section is not necessary. For completeness, we report results with background section in Section 5.3.

[3]`spacy.io` and `nltk.org`

a measure of soft-alignment (Zhang et al., 2020a). As factual correctness is critical, we also calculate a factuality $F_1$ (Zhang et al., 2020b; Hu et al., 2022, Fact.). This metric is based on a rule-based fact-extraction method, CheXpert (Irvin et al., 2019), which labels the status (present, absent, uncertain) of 14 radiological observations. By applying this procedure to both the reference and candidate summary, we can calculate a precision/recall of facts.

**Implementation and hyperparameters.** For all summarization models, we use the hyperparameters and code of the original papers. Below, we focus on deviations from those settings and report all hyperparameters in Appendix B.

For BertExt, BertAbs and GSum, we make three adaptations: (i) the summary length of BertExt is set to the average number of sentences selected by OracleExt, rounded to the nearest integer,[4] (ii) we reduce the training steps to 20,000 to account for the smaller datasets, and (iii) to address an exploding gradient problem, we reduce the initial learning rate by a factor of 10. For final testing, we take the checkpoint with lowest validation loss on MIMIC-CXR. On OpenI, we found the loss to be unstable, so opted to select models by validation R-1.

Regarding the guidance-length prediction models (Method 1 in Section 3.2), we experiment with two classifiers. First, a multinomial logistic regression classifier with unigram bag-of-words features (LR-APPROX). Second, as this model may be too simplistic to accurately predict the guidance length, we implement a transformer-based classifier (BERT-APPROX) on top of DistilBERT (Sanh et al., 2019).

## 4.2 Fixed-length Guidance (RQ1)

We first aim to understand if extractive summaries can be a useful guidance signal for radiology report summarization. To this end, we compare BertAbs (i.e., unguided) with GSum in its default configuration (Part 1 in Table 2).

We find that **GSum with fixed-length extractive guidance (Dou et al., 2021) does not generalize to the radiology domain.** Compared with BertAbs, effectiveness decreases by 4.5% and 3.2% in R-1 for MIMIC-CXR and OpenI, respectively. This is surprising as GSum is highly effective on multiple non-medical summarization benchmarks under the same experimental conditions (Dou et al., 2021). Our hypothesis is that highly varying summary lengths make the standard fixed-length guidance

in GSum ineffective on this data.[5] We empirically verify this hypothesis in the following experiments.

**Comparing oracle and automatic guidance.** To get an upper-bound estimate for extractive guidance signals, we analyze GSum in an unrealistic oracle setting. Recall from Section 3.1 that during training of GSum, the guidance signal is extracted by OracleExt, whereas during inference guidance is extracted by BertExt with a summary length fixed to $k = 1$ across all reports. If we instead also use OracleExt as guidance extractor during inference, we see a substantial increase in all metrics (R-1 46.3→58.8 on MIMIC, and R-1 60.1→68.8 on OpenI, all metrics in Appendix Table 4). **This oracle experiment demonstrates (i) that GSum learned to rely on guidance, and (ii) that extractive summaries can be a highly effective guidance signal if selected in the right way.**

Given that GSum is effective when we use the oracle guidance (OracleExt), it is important to understand how this guidance signal differs from the automatically extracted guidance (BertExt). We find that a characterizing difference between the two guidance signals is the length of the resulting summaries. OracleExt produces summaries with 0/1/2/3 sentences for 2/52/32/14% of the MIMIC-CXR reports, and for 15/67/14/3% of the OpenI reports. This implies that a guidance signal with a length of $k = 1$ is too short for 46% of the MIMIC-CXR reports, whereas on OpenI it is too short for 17% and too long for 15%.

## 4.3 Variable-length Guidance (RQ2)

We next evaluate the utility of our proposed variable-length extractive guidance signal (Part 2 of Table 2). We make several observations.

First, we find that variable-length extractive guidance substantially improves the effectiveness of GSum. On MIMIC-CXR, our adaptation is also better than unguided summarization (BertAbs). In particular, we observe a large increase in factuality, which is critical in the clinical domain. While we see a similar improvement of GSum on OpenI, this guided summarization model does not improve over BertAbs. One potential reason is that OpenI is more abstractive than MIMIC-CXR, as indicated by the high degree of novelty (Table 1) and the relatively low scores of the extractive methods (BertExt, OracleExt in Table 2). This corroborates the findings by Dou et al. (2021), where GSum

---

[4]On MIMIC-CXR and OpenI $|\text{OracleExt}(\boldsymbol{x}, \boldsymbol{y})| \approx 1$.

[5]Appendix A.1 gives the length distribution of targets.

| Method | MIMIC-CXR | | | | | OpenI | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | R-1 | R-2 | R-L | BS | Fact. | R-1 | R-2 | R-L | BS | Fact. |
| *Part 1: Baselines and reproduction of GSum* | | | | | | | | | | |
| OracleExt | 44.0 | 25.4 | 40.6 | 50.1 | 55.1 | 30.5 | 11.9 | 29.2 | 33.7 | 53.5 |
| BertExt (Liu and Lapata, 2019) | 32.7 | 18.1 | 30.0 | 41.9 | 44.5 | 23.6 | 7.4 | 22.6 | 32.2 | **42.8** |
| BertAbs (Liu and Lapata, 2019) | 48.4 | 34.1 | 46.6 | 58.8 | 47.3 | 62.0 | 52.7 | 61.7 | 69.2 | 39.3 |
| GSum (Dou et al., 2021) | 46.3 | 32.7 | 44.7 | 57.4 | 46.6 | 60.1 | 49.6 | 59.8 | 67.0 | 40.0 |
| *Part 2: GSum adapted with a variable-length guidance signal (ours)* | | | | | | | | | | |
| GSum w/ LR-Approx | 48.9 | 34.2 | 47.0 | 59.1 | 48.2 | 62.0 | 51.2 | 61.6 | 67.9 | 41.7 |
| GSum w/ BERT-Approx | 49.4 | 34.5 | 47.4 | 59.5 | 50.6 | 62.5 | 51.6 | 62.2 | 68.4 | 39.6 |
| GSum w/ Thresholding | **49.9** | 34.3 | **47.8** | **59.8** | 49.0 | 62.2 | 50.8 | 61.8 | 68.6 | 40.4 |
| *Part 3: Comparison with domain-specific methods* | | | | | | | | | | |
| WGSum (Hu et al., 2021) | 48.4 | 32.8 | 46.5 | 58.6 | 49.8 | 61.1 | 50.0 | 60.8 | 67.9 | 38.4 |
| WGSum+CL (Hu et al., 2022) | 49.5 | **35.3** | 47.8 | 59.5 | **51.1** | **64.7** | **57.1** | **64.5** | **70.0** | 37.2 |
| WGSum (Hu et al., 2021)† | 48.3 | 33.3 | 46.6 | — | — | 61.6 | 50.9 | 61.7 | — | — |
| WGSum+CL (Hu et al., 2022)† | 49.1 | 33.7 | 47.1 | — | — | 64.9 | 55.5 | 64.4 | — | — |

Table 2: Technical evaluation of unguided, guided and domain-guided methods on two datasets. Metrics are ROUGE-1/2/L, BERTScore (BS) and CheXpert factuality F1 (Fact). All results were obtained by re-implementing the models with the official code of respective papers, results directly cited are indicated with †.

| BertExt length ($k = \cdot$) | MIMIC-CXR | | OpenI | |
|---|---|---|---|---|
| | R-1 (Prec./Rec.) | $|\hat{y}|$ | R-1 (Prec./Rec.) | $|\hat{y}|$ |
| Fixed ($k = 1$) | 32.7 (**38.5**/34.2) | 1.0 | **23.6** (**24.6**/26.9) | 1.0 |
| LR-APPROX | 34.5 (35.7/40.0) | 1.4 | 23.5 (23.9/27.2) | 1.1 |
| BERT-APPROX | 35.2 (34.6/42.0) | 1.5 | 23.5 (23.7/27.5) | 1.1 |
| Thresholding | **36.1** (34.1/**46.3**) | 1.7 | 23.2 (22.9/**29.0**) | 1.2 |
| $k = |\text{OracleExt}|$ | 36.9 (35.3/44.2) | 1.6 | 24.3 (23.2/29.2) | 1.2 |

Table 3: Comparing strategies for extracting variable-length summaries with BertExt by measuring ROUGE against the gold summary. Average summary length $|\hat{y}|$ given in sentences. All methods are tested as guidance signal for GSum in Table 2.



Figure 3: Evaluating summaries by target length on MIMIC-CXR (equal number of samples per bucket).

was less effective on more abstractive datasets. For future work, it would be interesting to study the interplay between the degree of abstraction, and the utility of extractive guidance signals.

Second, regarding the different strategies to obtain variable-length extractive summaries, we cannot conclude that one is superior over another. The classifier-based approaches (LR-Approx, BERT-Approx), and the thresholding-based approach (Thresholding) lead to similar results when the extracted guidance is used downstream in GSum. For each guidance extraction strategy, we calculate the ROUGE scores of the guidance signal with respect to the gold summaries. From Table 3, we see that all strategies have the desired effect of increasing content recall, with a smaller sacrifice in precision.

Third, to better understand how guidance influences the quality of summaries, we plot the R-1 scores across different target summary lengths (Figure 3). We find that variable-length guidance im-
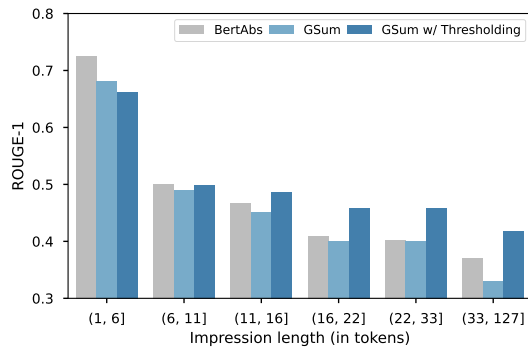
proves the quality of longer summaries, while for shorter targets, extractive guidance is not beneficial. By manual inspection, we find that short targets are standard phrasings of negative results (e.g., *"No evidence of acute findings"*), whereas longer targets have a higher extractive component by reiterating particular findings. In practice, it could be interesting to combine unguided and guided methods by letting the radiologist decide whether a long or short summary should be generated.

**Comparison with domain-specific guided summarization (WGSum, WGSum+CL).** Lastly, compared with the domain-specific guided methods (Part 3 of Table 2), we find on MIMIC-CXR that GSum with variable-length extractive guidance is just as effective as WGSum and WGSum+CL which use a graph of clinical entities. On OpenI, our approach improves over WGSum, but is slightly worse than WGSum+CL.

**Summary of RQ1/RQ2.** Overall our results show that extractive summaries are a promising guidance signal for clinical reports without requiring any domain-specific resources. We envision that this makes it easier to adopt guided summarization in other clinical domains and languages, for which domain-specific resources like ontologies and clinical NER models are not widely available.

## 5 Error Analysis

> **RQ3.** What are the errors and failure modes of unguided and guided methods for radiology report summarization?

### 5.1 Evaluation Setup

Inspired by the Multidimensional Quality Metrics framework for evaluation of machine translation systems (Lommel et al., 2014), we conduct a span-based error annotation. We task annotators to highlight erroneous text spans and to classify them according to an error taxonomy. As a starting point, we use the taxonomy proposed by Yu et al. (2022). Based on two pilot runs, we extended this taxonomy from initially 6 to 11 fine-grained error categories (see Figure 4) and developed a definition and examples for each. Following Yu et al. (2022), we opt for a reference-based evaluation. We want to understand how the system generated summary differs from the clinician summary both in content and correctness of the presented facts. Therefore, our errors can be grouped into additions (spans in the candidate), omissions (spans in the reference), and binary choices for the correctness of presented facts. Further, we ask annotators to flag any additional errors they encounter as a free-form answer. We provide full annotation guidelines in Appendix D.

**Materials.** We randomly select 100 reports from the official test set of MIMIC-CXR which is stratified to cover both frequent and less frequent inputs/clinical observations (Johnson et al., 2019b). For each input, we generate four candidate summaries using BertAbs (representative of unguided systems), GSum w/ Thresholding (representative of systems with domain-agnostic guidance), and WG-Sum/WGSum+CL (representative of systems with domain-specific guidance). We present the reference summary and all candidates (in random order) at once to annotators to reduce effort and ensure consistent annotation of similar summaries. Each set of summaries is completed by three annotators resulting in 1,200 error annotations (100 reports × 4 candidates × 3 annotators). We form a "gold standard" from the triple annotation by majority voting (example aggregation in Appendix C).

**Annotators.** To account for the domain knowledge necessary for the annotation task, we hired 6 senior medical students in their fifth year of training. All annotators are fluent in English. We compensated annotators with 10.5€ per hour (standard rate for student assistants in Germany). The annotation took 23.1 hours (avg. 4.6 min/sample), plus additional time for pilot rounds and discussions.

### 5.2 Results (RQ3)

We report aggregated error counts and example annotations in Figure 4.[6]

Overall, we find that the prevalence of errors is comparable across the investigated methods, and that only 14–22% of generated summaries are error-free. The most common errors are omissions and additions of findings, which indicates that the models struggle to select relevant content (1a. 43–52%; 2a. 44–57%). Compared with unguided summarization, there is a slight trend that guided methods reduce the risk of omissions, while only WG-Sum+CL succeeds at doing this without sacrificing precision. Even though additions are common, they rarely contradict the reference (2e. 0–3%). Similarly, when both the reference and candidate present the same findings, errors related to their clinical correctness are rare (3. 5–8%; 4. 6–9%).

A surprising finding is the common omission and addition of clinicians' communications (1d. 16–20%; 2d. 3–8%). By manual inspection (examples in Figure 4), we find that these are specific actions that a clinician performed after the examination such as informing colleagues about the findings, or recommending additional analysis. Additions of this kind have likely no grounding in the underlying report. To successfully generate such statements, models would require additional context information or guidance from a user.

### 5.3 Discussion

Overall, our error analysis reveals that the key differences between model-generated impressions and radiologists' impressions relate to content selection (i.e., a tension between completeness/recall and

---

[6]To measure inter-annotator agreement (IAA), we calculate $F_1$ for span-annotations (Deleger et al., 2012) and Krippendorffs' Alpha for binary judgments (Krippendorff, 1970). Aggregated IAA: 1. Omissions: 0.61, 2. Additions: 0.60, 3. Incorrect Location: 0.25, and 4. Incorrect Severity: 0.41. IAA by error category for span-level annotations in Appendix C.

| # | Error Category | M1 (%) | M2 (%) | M3 (%) | M4 (%) |
|---|---|---|---|---|---|
| 0 | No error | 20 (20) | 18 (18) | 14 (14) | **22** (22) |
| | *Omissions from reference* | | | | |
| 1a | Finding/interpretation | 70 (52) | **58** (43) | 62 (48) | 64 (47) |
| 1b | Comparison | 23 (19) | **16** (15) | 19 (16) | 23 (19) |
| 1c | Ref. to prior report | **1** (1) | 3 (3) | 2 (2) | 2 (2) |
| 1d | Communication/followup | 20 (19) | **18** (16) | 19 (17) | 19 (17) |
| | Total | 114 (66) | **95** (58) | 102 (63) | 108 (61) |
| | *Additions to candidate* | | | | |
| 2a | Finding/interpretation | **51** (44) | 72 (57) | 61 (50) | 54 (46) |
| 2b | Comparison | 11 (8) | 10 (9) | 9 (9) | **7** (6) |
| 2c | Ref. to prior report | **0** (0) | 1 (1) | **0** (0) | **0** (0) |
| 2d | Communication/followup | 5 (5) | 8 (6) | 8 (8) | **4** (3) |
| 2e | Contradicting finding | **0** (0) | 1 (1) | 3 (3) | 1 (1) |
| | Total | 67 (49) | 92 (63) | 81 (58) | **66** (48) |
| | *Semantics of intersecting findings* | | | | |
| 3 | Incorrect location | **5** (5) | 8 (8) | 8 (8) | 7 (7) |
| 4 | Incorrect severity | **6** (6) | 7 (7) | 7 (7) | 9 (9) |
| 5 | Other error | 31 (23) | **30** (23) | 33 (29) | **30** (21) |

**Reference:** Interval increase in vascular engorgement. No frank interstitial edema. No focal consolidations identified.
**Candidate (M3):** interval increase in pulmonary vascular congestion without evidence of interstitial edema. small right-sided pleural effusion.

**Reference:** Right lower lobe opacity, possibly atelectasis, with associated moderate sized effusion.
**Candidate (M4):** persistent right lower lobe opacity with associated effusion, mildly progressed from the preceding radiograph.

**Reference:** Multiloculated right pleural effusion unchanged since _. New linear and nodular opacities in the left upper lobe may represent carcinomatosis. Findings were relayed to Dr. _ by Dr. _ _ following review on _ at approximiately 11:00 via telephone.
**Candidate (M1):** stable appearance of multiple loculated right pleural effusion.

**Reference:** Unchanged size and position of right-sided hydropneumothorax over the last _-hour examination interval.
**Candidate (M3):** development of new right-sided hydropneumothorax in this patient with history of newly placed pigtail catheter. referring physician, _. _ was paged at 4:45 p.m.

**Reference:** Little change in the severe bronchiectasis and emphysema.
**Candidate (M3):** unchanged bibasilar bronchiectasis and bibasilar bronchiectasis.

Figure 4: Results of manual error analysis of 100 MIMIC-CXR reports. Left: number of times each error occurred per method (percent of reports in gray, least errors per row in bold). Right: example error annotations. Models: BertAbs (**M1**), GSum w/ Thresholding (**M2**), WGSum (**M3**), and WGSum+CL (**M4**) [best viewed in color].

relevance/precision). We offer two hypotheses to explain the models' difficulties in this area.

First, there may be latent factors that explain which findings are included in the impression. Among those factors could be patient demographics, the radiograph, prior exams and the clinical question. Typically, this information is available to radiologists through the electronic health records, and is partly documented in the background section of radiology reports. Early work explored using the background section as guidance (Zhang et al., 2018), but more recent work commonly excluded it in pre-processing (Sotudeh et al., 2020; Hu et al., 2021, 2022). We present evaluation results when including the background and observe an overall improvement in almost all metrics for abstractive methods (Appendix A.5). This improvement indicates that (i) additional context supports content selection, and (ii) it could be useful to explicitly model the background in guided summarization.

Second, we anecdotally observed a substantial degree of duplication in the MIMIC-CXR corpus, where reports with identical findings have different impressions (examples in Appendix A.6).[7] This may lead to corpus-level inconsistencies preventing models to reliably learn the selection of findings. We note that there can be numerous reasons for these duplication induced inconsistencies, including the presence of latent factors (see above) and

remaining subjectivity/uncertainty in radiologists' assessments. We leave the investigation of this aspect of data quality and potential effects of training data deduplication for future work.

## 5.4 Limitations

We note two limitations of this error analysis.

First, the analysis is based on comparing candidate impressions with reference impressions. In the absence of the full clinical context, we argue that this is the most reliable benchmark for completeness and relevance of summaries. However, we recognize that we cannot draw any conclusions about the factuality of additions with respect to the full report. To give a first factuality estimate, we conducted a post-hoc analysis with RadNLI (Miura et al., 2021). Let $x_i$ be a sentence in report $\boldsymbol{x} = (x_1, \ldots, x_n)$, and $s$ be an addition span. If RadNLI predicts a contradiction for any $(x_i, s)$ pair, we label this span as contradicting and neutral/entailed otherwise. We find that between 23.4% (BertAbs) and 29.3% (GSum w/ Thresholding) of additions are contradicting, indicating that factuality is another challenge for current models (details in Appendix A.7).

Second, the sample size was driven by time and resource constraints ($N = 100$). To estimate representativeness of this sample, we compare descriptive statistics of the sample with those of the whole test set (length, novelty, compression), and observe that these largely agree (see Appendix C). While we believe that this sample is sufficient to support

---

[7]11.9% of the 122,500 MIMIC-CXR training reports have a findings section occurring more than once. Among those reports are only 1036 distinct impressions.

the qualitative conclusions about the failure modes of current methods, a larger study is warranted when the goal is to quantitatively compare the efficacy of different methods.

# 6 Conclusion

In this work, we revisited guided abstractive summarization of radiology reports. We demonstrated that extractive summaries can be an effective guidance signal for the task, if we allow the length of this guidance signal to vary across reports, and thereby make the gap between domain-agnostic and domain-specific guidance smaller. Furthermore, through a fine-grained error analysis of unguided and guided models we found that guidance successfully steers the content of summaries but that significant deficits in content selection persist.

We hope that this paper motivates future efforts on content selection mechanisms for radiology report summarization, their evaluation in other domains and languages, and on more comprehensive evaluation suites. We release our error annotations which can serve as a starting point for evaluating the efficacy of metrics in capturing these errors.

## Ethical Considerations

**Privacy sensitive datasets.** Both the MIMIC-CXR dataset (Johnson et al., 2019a), and the OpenI dataset (Demner-Fushman et al., 2015) were fully de-identified by the dataset authors in compliance with applicable privacy laws (HIPAA). This includes the removal of any protected health information that may directly or indirectly identify a patient. Nevertheless, the data is still privacy sensitive, and special care was taken to only process it within secured computing infrastructure.

**Intended use.** We believe that the proposed methods can improve the workflow of clinicians both by reducing the documentation effort and encouraging higher-quality reporting, and thereby improving patient care. However, as our results and discussion show, state-of-the-art summarization methods may not have the desired level of quality that is needed in high-stakes domains such as the clinical context. Therefore, our work is not to be understood in the context of a system that can be deployed, but rather as a step toward a better understanding of the shortcomings of current text summarization methods and providing insight into how these can solved.

## Supplementary Materials Availability Statement

- Detailed analysis, hyperparameters, and annotation guidelines are available in Appendices A to D.
- Source code to reproduce all experiments is available from `github.com/jantrienes/inlg2023-radsum/`
- The expert annotations of summarization errors are available from `github.com/jantrienes/inlg2023-radsum/` under the PhysioNet Credentialed Health Data License 1.5.0.
- The MIMIC-CXR (v2.0.0) dataset is available from `physionet.org/content/mimic-cxr/` under the PhysioNet Credentialed Health Data License 1.5.0.
- The OpenI dataset is available from `openi.nlm.nih.gov` (no license terms stated)

## Acknowledgements

## References

Xiaoyan Cai, Sen Liu, Junwei Han, Libin Yang, Zhenguo Liu, and Tianming Liu. 2021. ChestXRayBERT: A pretrained language model for chest radiology report summarization. *IEEE Transactions on Multimedia*, 25:845–855.

Ziqiang Cao, Wenjie Li, Sujian Li, and Furu Wei. 2018. Retrieve, rerank and rewrite: Soft template based neural summarization. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 152–161.

Yen-Chun Chen and Mohit Bansal. 2018. Fast abstractive summarization with reinforce-selected sentence rewriting. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 675–686.

Jean-Benoit Delbrouck, Maya Varma, and Curtis P. Langlotz. 2022. Toward expanding the scope of radiology report summarization to multiple anatomies and modalities. In *Proceedings of the 2nd Machine Learning for Health symposium (ML4H), Extended Abstract Collection*.

Louise Deleger, Qi Li, Todd Lingren, Megan Kaiser, Katalin Molnar, Laura Stoutenborough, Michal

Kouril, Keith Marsolo, and Imre Solti. 2012. Building gold standard corpora for medical natural language processing tasks. In *AMIA Annual Symposium Proceedings*, pages 144–153.

Dina Demner-Fushman, Marc D. Kohli, Marc B. Rosenman, Sonya E. Shooshan, Laritza Rodriguez, Sameer Antani, George R. Thoma, and Clement J. McDonald. 2015. Preparing a collection of radiology examinations for distribution and retrieval. *Journal of the American Medical Informatics Association*, 23(2):304–310.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL)*, pages 4171–4186.

Zi-Yi Dou, Pengfei Liu, Hiroaki Hayashi, Zhengbao Jiang, and Graham Neubig. 2021. GSum: A general framework for guided neural abstractive summarization. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL)*, pages 4830–4842.

Alexander R. Fabbri, Wojciech Kryściński, Bryan McCann, Caiming Xiong, Richard Socher, and Dragomir Radev. 2021. SummEval: Re-evaluating summarization evaluation. *Transactions of the Association for Computational Linguistics*, 9:391–409.

Angela Fan, David Grangier, and Michael Auli. 2018. Controllable abstractive summarization. In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 45–54.

Esteban F. Gershanik, Ronilda Lacson, and Ramin Khorasani. 2011. Critical finding capture in the impression section of radiology reports. In *AMIA Annual Symposium Proceedings*, pages 465–469.

Junxian He, Wojciech Kryściński, Bryan McCann, Nazneen Rajani, and Caiming Xiong. 2022. CTRLsum: Towards generic controllable text summarization. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5879–5915.

Jinpeng Hu, Jianling Li, Zhihong Chen, Yaling Shen, Yan Song, Xiang Wan, and Tsung-Hui Chang. 2021. Word graph guided summarization for radiology findings. In *Findings of the Association for Computational Linguistics (ACL-IJCNLP)*, pages 4980–4990.

Jinpeng Hu, Zhuo Li, Zhihong Chen, Zhen Li, Xiang Wan, and Tsung-Hui Chang. 2022. Graph enhanced contrastive learning for radiology findings summarization. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 4677–4688.

Dandan Huang, Leyang Cui, Sen Yang, Guangsheng Bao, Kun Wang, Jun Xie, and Yue Zhang. 2020. What have we achieved on text summarization? In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 446–469.

Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silviana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, Jayne Seekins, David A. Mong, Safwan S. Halabi, Jesse K. Sandberg, Ricky Jones, David B. Larson, Curtis P. Langlotz, Bhavik N. Patel, Matthew P. Lungren, and Andrew Y. Ng. 2019. CheXpert: A large chest radiograph dataset with uncertainty labels and expert comparison. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 590–597.

Ruipeng Jia, Yanan Cao, Haichao Shi, Fang Fang, Pengfei Yin, and Shi Wang. 2021. Flexible non-autoregressive extractive summarization with threshold: How to extract a non-fixed number of summary sentences. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 13134–13142.

Alistair E. W. Johnson, Tom J. Pollard, Seth J. Berkowitz, Nathaniel R. Greenbaum, Matthew P. Lungren, Chih-ying Deng, Roger G. Mark, and Steven Horng. 2019a. MIMIC-CXR, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific Data*, 6(317).

Alistair E. W. Johnson, Tom J. Pollard, Nathaniel R. Greenbaum, Matthew P. Lungren, Chih ying Deng, Yifan Peng, Zhiyong Lu, Roger G. Mark, Seth J. Berkowitz, and Steven Horng. 2019b. MIMIC-CXR-JPG, a large publicly available database of labeled chest radiographs. *CoRR*, abs/1901.07042.

Charles E. Kahn, Curtis P. Langlotz, Elizabeth S. Burnside, John A. Carrino, David S. Channin, David M. Hovsepian, and Daniel L. Rubin. 2009. Toward best practices in radiology reporting. *Radiology*, 252(3):852–856.

Sanjeev Kumar Karn, Ning Liu, Hinrich Schütze, and Oladimeji Farri. 2022. Differentiable multi-agent actor-critic for multi-step radiology report summarization. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 1542–1553.

Klaus Krippendorff. 1970. Bivariate agreement coefficients for reliability of data. *Sociological methodology*, 2:139–150.

Wojciech Kryściński, Bryan McCann, Caiming Xiong, and Richard Socher. 2020. Evaluating the factual consistency of abstractive text summarization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9332–9346.

Chenliang Li, Weiran Xu, Si Li, and Sheng Gao. 2018. Guiding generation for abstractive text summarization based on key information guide network. In

*Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL)*, pages 55–60.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81.

Yang Liu and Mirella Lapata. 2019. Text summarization with pretrained encoders. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3730–3740.

Arle Lommel, Hans Uszkoreit, and Aljoscha Burchardt. 2014. Multidimensional quality metrics (MQM): A framework for declaring and describing translation quality metrics. *Tradumàtica*, 2014(12):455–463.

Sean MacAvaney, Sajad Sotudeh, Arman Cohan, Nazli Goharian, Ish Talati, and Ross W. Filice. 2019. Ontology-aware clinical abstractive summarization. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, pages 1013–1016.

Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 1906–1919.

Yasuhide Miura, Yuhao Zhang, Emily Tsai, Curtis Langlotz, and Dan Jurafsky. 2021. Improving factual completeness and consistency of image-to-text radiology report generation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 5288–5304.

Ramesh Nallapati, Feifei Zhai, and Bowen Zhou. 2017. SummaRuNNer: A recurrent neural network based sequence model for extractive summarization of documents. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 3075–3081.

Ramesh Nallapati, Bowen Zhou, Cicero dos Santos, Çağlar Gülçehre, and Bing Xiang. 2016. Abstractive text summarization using sequence-to-sequence RNNs and beyond. In *Proceedings of the 20th SIGNLL Conference on Computational Natural Language Learning (CoNLL)*, pages 280–290.

Shashi Narayan, Yao Zhao, Joshua Maynez, Gonçalo Simões, Vitaly Nikolaev, and Ryan McDonald. 2021. Planning with learned entity prompts for abstractive summarization. *Transactions of the Association for Computational Linguistics*, 9:1475–1492.

Alexander M. Rush, Sumit Chopra, and Jason Weston. 2015. A neural attention model for abstractive sentence summarization. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 379–389.

Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. In *Proceedings of the 5th Workshop on Energy Efficient Machine Learning and Cognitive Computing (EMC2-NIPS)*.

Sajad Sotudeh, Nazli Goharian, and Ross W. Filice. 2020. Attend to medical ontologies: Content selection for clinical abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 1899–1905.

Emiel van Miltenburg, Miruna Clinciu, Ondřej Dušek, Dimitra Gkatzia, Stephanie Inglis, Leo Leppänen, Saad Mahamood, Emma Manning, Stephanie Schoch, Craig Thomson, and Luou Wen. 2021. Underreporting of errors in NLG output, and what to do about it. In *Proceedings of the 14th International Conference on Natural Language Generation (INLG)*, pages 140–153.

Dave Van Veen, Cara Van Uden, Maayane Attias, Anuj Pareek, Christian Bluethgen, Malgorzata Polacin, Wah Chiu, Jean-Benoit Delbrouck, Juan Zambrano Chaves, Curtis Langlotz, Akshay Chaudhari, and John Pauly. 2023. RadAdapt: Radiology report summarization via lightweight domain adaptation of large language models. In *The 22nd Workshop on Biomedical Natural Language Processing and BioNLP Shared Tasks*, pages 449–460.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 5998–6008.

Qianqian Xie, Jiayu Zhou, Yifan Peng, and Fei Wang. 2023. FactReranker: Fact-guided reranker for faithful radiology report summarization. *CoRR*, abs/2303.08335.

Feiyang Yu, Mark Endo, Rayan Krishnan, Ian Pan, Andy Tsai, Eduardo Pontes Reis, Eduardo Kaiser Ururahy Nunes Fonseca, Henrique Min Ho Lee, Zahra Shakeri Hossein Abad, Andrew Y. Ng, Curtis P. Langlotz, Vasantha Kumar Venugopal, and Pranav Rajpurkar. 2022. Evaluating progress in automatic chest x-ray radiology report generation. *medRxiv*.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020a. BERTScore: Evaluating text generation with BERT. In *Proceedings of the 8th International Conference on Learning Representations (ICLR)*.

Yuhao Zhang, Daisy Yi Ding, Tianpei Qian, Christopher D. Manning, and Curtis P. Langlotz. 2018. Learning to summarize radiology findings. In *Proceedings of the 9th International Workshop on Health Text Mining and Information Analysis (LOUHI)*, pages 204–213.

Yuhao Zhang, Derek Merck, Emily Bao Tsai, Christopher D. Manning, and Curtis P. Langlotz. 2020b. Optimizing the factual correctness of a summary: A study of summarizing radiology reports. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 5108–5120.

Yuhao Zhang, Yuhui Zhang, Peng Qi, Christopher D. Manning, and Curtis P. Langlotz. 2021. Biomedical and clinical english model packages for the stanza python NLP library. *Journal of the American Medical Informatics Association*, 28(9):1892–1899.

Yunqi Zhu, Xuebing Yang, Yuanyuan Wu, and Wensheng Zhang. 2023. Leveraging summary guidance on medical report summarization. *CoRR*, abs/2302.04001.

## A  Detailed Analysis

To support replication, this section provides supplementary analysis on the results of the main part.

### A.1  Target Impression Length Distribution and Evaluation by Target Length

We demonstrated in the main part that variable-length guidance helps to adapt to varying target lengths. To better interpret this result, we plot the length distribution of target summaries and the ROUGE-1 score by target-length interval in Figure 5. It can be observed that the length distribution has a long tail with a peak around 4-5 tokens. Impressions of this length are standard phrasings to indicate that no abnormalities were found (e.g., *"No evidence of acute findings"*).

### A.2  Evaluating GSum in an Oracle Setting

As a supplement to the oracle experiment in Section 4.2, we provide all metrics for the three inference settings of GSum in Table 4: (i) automatic fixed-length guidance (i.e., extracted from BertExt with $k = 1$), (ii) automatic variable-length guidance but with an oracle length (i.e., BertExt with $k = |\text{OracleExt}(\boldsymbol{x}, \boldsymbol{y})|$), and (iii) oracle guidance (i.e., $\boldsymbol{g} = \text{OracleExt}(\boldsymbol{x}, \boldsymbol{y})$).

### A.3  BertExt: Evaluating Fixed-length Settings

To evaluate if larger values for $k$ in the fixed-summary length setting would improve the effectiveness of BertExt, we generate summaries for all settings of $k = \{1, ..., 5\}$. Analogously, we provide these summaries as guidance signal to GSum. Table 5 reports the results of this experiment. While we find that larger settings of $k$ lead to an increase in recall, we see an equally strong drop in precision,

| MIMIC-CXR | R-1 | R-2 | R-L | BS | Fact. |
|---|---|---|---|---|---|
| *Guidance signal for GSum* | | | | | |
| Fixed (Dou et al., 2021) | 46.3 | 32.7 | 44.7 | 57.4 | 46.6 |
| Oracle Length | 51.7 | 36.3 | 49.6 | 61.2 | 52.4 |
| Oracle Length + Content | 58.5 | 42.0 | 56.2 | 66.0 | 60.0 |
| **OpenI** | | | | | |
| *Guidance signal for GSum* | | | | | |
| Fixed (Dou et al., 2021) | 60.1 | 49.6 | 59.8 | 67.0 | 40.0 |
| Oracle Length | 63.9 | 53.0 | 63.5 | 69.4 | 42.3 |
| Oracle Length + Content | 68.8 | 56.7 | 68.3 | 72.7 | 45.1 |

Table 4: Evaluating GSum in an oracle setting. *Fixed* is reproduced from Table 2.

| | MIMIC-CXR | | | | | OpenI | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | R-1 | R-2 | R-L | $B_P$ | $B_R$ | R-1 | R-2 | R-L | $B_P$ | $B_R$ |
| BertExt with fixed-length summaries | | | | | | | | | | |
| $k = 1$ | 32.7 | 18.1 | 30.0 | **45.2** | 40.1 | **23.6** | **7.4** | **22.6** | **33.6** | 32.3 |
| $k = 2$ | **34.1** | **18.6** | **31.3** | 40.9 | 50.1 | 19.7 | 6.7 | 18.9 | 28.3 | 39.9 |
| $k = 3$ | 31.7 | 17.0 | 29.2 | 37.0 | 53.5 | 17.4 | 6.1 | 16.6 | 25.9 | 42.8 |
| $k = 4$ | 29.1 | 15.4 | 26.8 | 34.0 | 54.6 | 15.8 | 5.5 | 15.1 | 24.0 | 43.7 |
| $k = 5$ | 27.2 | 14.3 | 25.2 | 32.2 | **54.9** | 15.1 | 5.2 | 14.4 | 23.3 | **44.1** |
| GSum with fixed-length guidance extracted from BertExt | | | | | | | | | | |
| $k = 1$ | 46.3 | 32.7 | 44.7 | 64.6 | 52.8 | 60.1 | 49.6 | 59.8 | 67.0 | 68.5 |
| $k = 2$ | 46.3 | 30.3 | 44.2 | 58.1 | 58.5 | 54.3 | 43.2 | 53.9 | 61.2 | 66.2 |
| $k = 3$ | 44.1 | 27.7 | 41.9 | 53.6 | 59.9 | 54.6 | 43.2 | 54.1 | 61.6 | 67.3 |
| $k = 4$ | 42.2 | 26.0 | 40.2 | 50.4 | 60.2 | 53.5 | 42.1 | 53.1 | 60.1 | 67.5 |
| $k = 5$ | 40.8 | 24.6 | 38.8 | 48.3 | 60.1 | 52.7 | 41.3 | 52.2 | 59.5 | 67.5 |

Table 5: Testing fixed-length summaries ($k \in [1, 5]$) for BertExt (first block) and as GSum guidance (second block). Metrics are ROUGE-1/2/L and BERTScore precision ($B_P$) and recall ($B_R$)

both on BertExt and GSum which demonstrates the necessity of variable-length extractive guidance.

### A.4  Evaluating Guidance Length Prediction

To predict the length of OracleExt in the variable-length guidance setting, we employ a logistic regression classifier and a BERT-based classifier (cf. Section 4.1). Detailed evaluation results for both classification models are given in Table 6.

### A.5  Including the Background Section

To understand to what extent the background section carries important information for summarizing findings to impression, we prepend it to the findings section and retrain all models. It can be observed that this change improves most abstractive methods on both datasets (Figure 6). For extractive methods results stay largely on par or get worse, indicating that these models do not effectively integrate the background information.
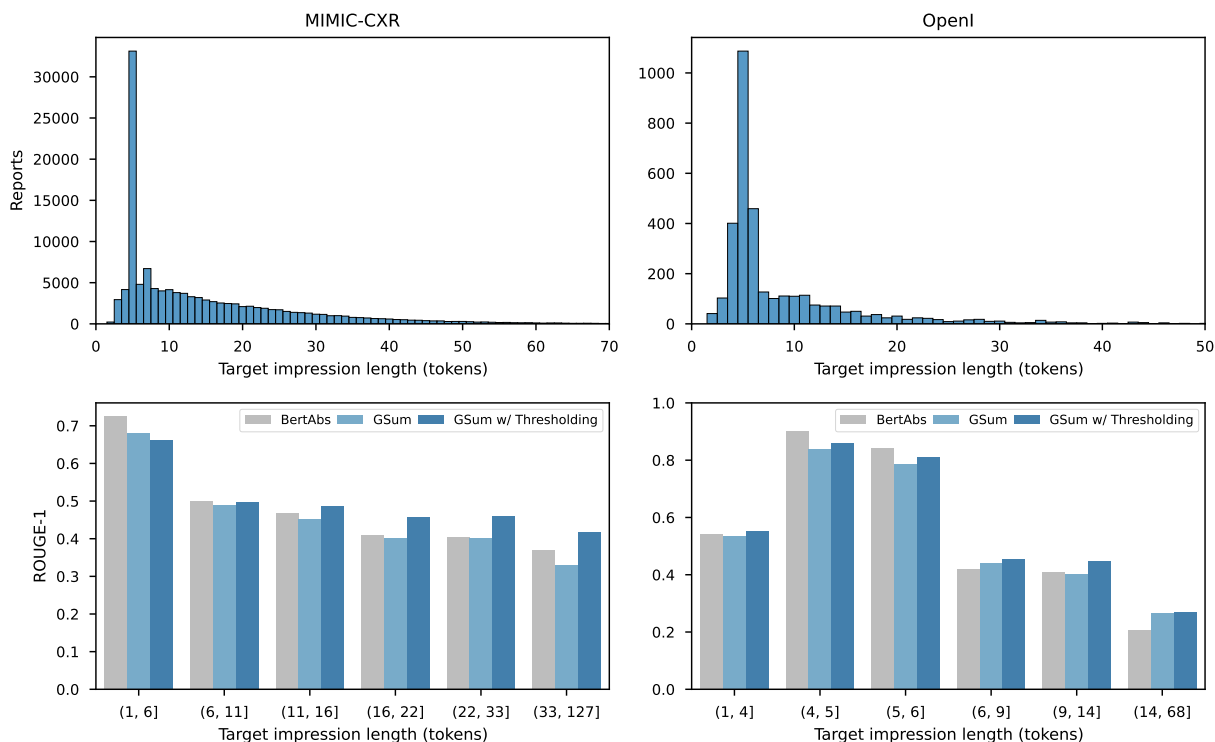
Figure 5: Top row: distribution of target impression lengths. Bottom row: ROUGE-1 by target length for BertAbs (unguided summarization), GSum (fixed-length guidance) and GSum w/ Thresholding (variable-length guidance).



Figure 6: Training models to summarize both background and findings improves most abstractive methods. Scores as absolute delta to the same models without the background section (cf. Table 2).

## A.6 Examples of Duplicated Findings and Impressions

We anecdotally observed a large degree of duplication within MIMIC-CXR which may cause corpus-level inconsistencies (see discussion in Section 5.3). This section further quantifies the degree of duplication and provides several examples. Throughout, we only consider instances of *exact* duplication. Of the 122,500 training reports in MIMIC-CXR, we

find that 11.9% have a findings section occurring more than once. We present examples of duplicate findings with *different* impressions in Table 12. In addition, we calculate a *label entropy* over the probabilities that each impression occurs for a given finding. We posit that duplicate finding-impression pairs may negatively impact model training in two ways. First, for findings with a high label entropy, the training loss cannot not stabilize (i.e., it is not clear which impression the model should favor). Second, for findings with a low label entropy, the model may learn a kind of "majority vote," which in turn may render models not sensitive enough to generate useful summaries for slightly different findings. We leave further investigation of report duplication to future work.

## A.7 Factuality of Additions

As discussed in Section 5.3, we use RadNLI (Miura et al., 2021) to get a first estimate for the factuality of additions marked by annotators in the error analysis. RadNLI obtained an accuracy of 77.8% on a test set of 480 manually labeled sentence pairs in MIMIC-CXR (Miura et al., 2021), which we consider sufficient for an initial exploration of the factuality of additions. Table 7 presents a breakdown of the RadNLI predictions for all addition

|        | **LR-Approx** | **BERT-Approx** |        |
|--------|---------------|-----------------|--------|
| **Target** | F-1 (Prec./Rec.) | F-1 (Prec./Rec.) | Support |
| $k = 0$ | 46.2 (80.0/32.4) | 53.7 (60.0/48.6) | 37 |
| $k = 1$ | 71.1 (63.4/80.9) | 71.7 (68.9/74.9) | 824 |
| $k = 2$ | 39.7 (43.1/36.9) | 46.7 (45.1/48.4) | 512 |
| $k = 3$ | 30.9 (53.3/21.8) | 43.2 (61.5/33.3) | 225 |
| Macro Avg. | 47.0 (59.9/43.0) | 53.9 (58.9/51.3) | 1,598 |
| On training set | 52.3 (64.1/47.6) | 62.5 (69.7/58.5) | 122,500 |

(a) Dataset: MIMIC-CXR

|        | **LR-Approx** | **BERT-Approx** |        |
|--------|---------------|-----------------|--------|
| **Target** | F-1 (Prec./Rec.) | F-1 (Prec./Rec.) | Support |
| $k = 0$ | 77.7 (85.9/70.9) | 84.0 (86.6/81.6) | 103 |
| $k = 1$ | 84.6 (77.2/93.6) | 85.4 (79.8/92.0) | 450 |
| $k = 2$ | 19.8 (36.1/13.7) | 28.4 (39.6/22.1) | 95 |
| $k = 3$ | 15.4 (50.0/9.1) | 8.7 (100.0/4.5) | 22 |
| Macro Avg. | 49.4 (62.3/46.8) | 51.6 (76.5/50.1) | 670 |
| On training set | 58.5 (83.3/54.3) | 51.1 (53.0/51.0) | 2,342 |

(b) Dataset: OpenI

Table 6: Precision, recall and F1 for length prediction of OracleExt. Scores are provided per class on the test set, and as macro-average for both the training and test set. Support indicates the number of samples in each class.

spans and models. It can be seen that the majority of additions is either neutral to the findings section, or entailed by it. Yet, between 23.4% and 29.3% of additions contradict at least one findings sentence, indicating that factuality of radiology report summarization methods can also further be improved.

## A.8 Error Analysis: Responses to *Other* Category

We analyze the annotators' comments from the *other* error category, and categorize these errors into two-level hierarchy using a bottom-up approach. Our categorization alongside definitions, examples and counts is shown in Table 13.

## B Replication Details for Modeling

We report hyperparameters of the summarization models in Table 9, and for models that predict the length of OracleExt (LR-APPROX/BERT-APPROX) in Table 10. All models were trained on NVIDIA RTX A6000 GPUs with 48GB of memory.

## C Replication Details for Error Analysis

**Sample statistics.** For inclusion in the error analysis, samples were drawn uniformly at random from the official test set of MIMIC-CXR. We compare statistics of the sample with those of the full test set in Table 8.

| Model | Entail | Neutral | Contradict |
|-------|--------|---------|------------|
| BertAbs | 31.9% | 44.7% | 23.4% |
| GSum w/ Thresholding | 34.5% | 36.2% | 29.3% |
| WGSum | 32.0% | 44.0% | 24.0% |
| WGSum+CL | 33.3% | 41.2% | 25.5% |

Table 7: Factuality of additions in candidates (i.e., spans categorized as "2a Finding/interpretation"), as per RadNLI (Miura et al., 2021).

| Aspect | **Full Test Set** | **Sample** |
|--------|-------------------|------------|
| Reports | 1,598 | 100 |
| Avg. $|\boldsymbol{x}|_t$ | $70 \pm 27.4$ | $63 \pm 20.4$ |
| Avg. $|\boldsymbol{x}|_s$ | $6.2 \pm 1.9$ | $5.7 \pm 1.6$ |
| Avg. $|\boldsymbol{y}|_t$ | $19 \pm 15.2$ | $18 \pm 12.4$ |
| Avg. $|\boldsymbol{y}|_s$ | $1.8 \pm 1.0$ | $1.7 \pm 0.9$ |
| Novelty | 69.8% | 69.7% |
| CMP | 71.9% | 70.3% |

Table 8: Statistics of the MIMIC-CXR test set and the sample used in the error analysis.

**Aggregating span-based annotations.** From the three annotations we form a "gold standard" as follows: for binary questions we take a majority vote. For span-based annotations, we first group (partially) overlapping spans, and then take a majority vote within each group. We provide an example for the majority voting of span-based annotations below. A1, A2, A3, denote annotators, and [-eX-] denotes an error of category X.

```
Tokens:   a   b   c   d   e   f   g   h
A1    : [-e1-]   [-----e2----]
A2    : [-e1-] [-e1-] [-e2-]
A3    : [-e1-]               [--e1--]
------------------------------------------
Group :    1         2          3
------------------------------------------
Vote  : [-e1-]     [-e2-]
```

**Inter-annotator agreement (IAA).** We calculate $F_1$ for span-annotations (Deleger et al. (2012), categories 1 and 2), and Krippendorffs' Alpha (Krippendorff, 1970) for binary judgments (categories 3 and 4) and report the IAA by category in Table 11.

| Parameter | BertExt | BertAbs | GSum | WGSum | WGSum+CL |
|---|---|---|---|---|---|
| Training Steps (MIMIC) | 20,000 | 20,000 | 20,000 | 50,000 | 100,000 |
| Training Steps (OpenI) | 20,000 | 20,000 | 20,000 | 20,000 | 20,000 |
| LR (Encoder) | 2e−3 | 2e−4 | 2e−4 | 5e−2 | 2e−4 |
| LR (Decoder) | n/a | 2e−2 | 2e−2 | 5e−2 | 5e−2 |
| Warmup (Encoder) | 10,000 | 20,000 | 20,000 | 8000 | 10,000 |
| Warmup (Decoder) | n/a | 10,000 | 10,000 | 8000 | 7000 |
| Dropout | 0.1 | 0.2 | 0.2 | 0.1 | 0.2 |
| Checkpoint freq. (MIMIC) | 1000 | 2000 | 2000 | 2000 | 2000 |
| Checkpoint freq. (OpenI) | 1000 | 2000 | 2000 | 200 | 200 |
| Decoding | n/a | Beam search | Beam search | Beam search | Beam search |
| Prediction length | n/a | $\geq 5$ tokens | $\geq 5$ tokens | $\geq 5$ tokens | $\geq 5$ tokens |
| Training GPUs | 3 | 5 | 5 | 4 | 3 |
| Inference GPUs | 1 | 1 | 1 | 1 | 1 |
| Base model | bert-base-uncased | bert-base-uncased | bert-base-uncased | None | dmis-lab/biobert-base-cased-v1.1 |
| Parameters | 120,512,513 | 180,222,522 | 205,433,914 | 82,260,794 | 221,600,069 |

Table 9: Hyperparameters of BertExt/BertAbs (Liu and Lapata, 2019), GSum (Dou et al., 2021), WGSum (Hu et al., 2021) and WGSum+CL (Hu et al., 2022). Training steps, warmup and learning rates were adapted as described in Section 4.1. Remaining parameters kept as in the original publications.

| Parameter | Setting |
|---|---|
| *LR-Approx* | |
| Features | Bag-of-words, unigrams with minimum document-frequency of 5, tf-idf |
| Parameters | 3718 (MIMIC-CXR), 592 (OpenI) |
| Regularization | L2 regularization with strength $C = 1$ |
| Solver | SAGA |
| Max. Iterations | 1000 |
| *BERT-Approx* | |
| Checkpoint | `distilbert-base-cased` |
| Parameters | 65,784,580 |
| Optimizer | Adam |
| Learning rate | 2e−5 |
| Epochs | 3 |
| Dropout | 0.2 |
| Batch size | 16 |
| Checkpoint freq. | 250 |
| Hardware | 6 GPUs |

Table 10: Hyperparameters for guidance length prediction models.

| # | Category | IAA | Count |
|---|---|---|---|
| *Omissions from reference* | | | |
| 1a | Finding/interpretation | 0.64 | 774 |
| 1b | Comparison | 0.34 | 236 |
| 1c | Ref. to prior report | 0.23 | 43 |
| 1d | Communication/followup | 0.83 | 216 |
| | Total | 0.61 | 1269 |
| *Additions to candidate* | | | |
| 2a | Finding/interpretation | 0.66 | 718 |
| 2b | Comparison | 0.44 | 155 |
| 2c | Ref. to prior report | 0.08 | 17 |
| 2d | Communication/followup | 0.65 | 72 |
| 2e | Contradicting finding | 0.26 | 34 |
| | Total | 0.60 | 996 |
| 3 | Incorrect location | 0.26 | 111 |
| 4 | Incorrect severity | 0.41 | 121 |

Table 11: Inter-annotator agreement (IAA) by category and total number of annotations before majority voting.

| # | Finding | Dups. | % | $|y^*|$ | H | Count | Top-5 Impressions |
|---|---------|-------|---|---------|---|-------|-------------------|
| | *Most frequent duplicates* | | | | | | |
| 1 | PA and lateral views of the chest provided. There is no focal consolidation, effusion, or pneumothorax. The cardiomediastinal silhouette is normal. Imaged osseous structures are intact. No free air below the right hemidiaphragm is seen. | 1141 | 0.93 | 26 | 0.12 | 1061 | No acute intrathoracic process. |
| | | | | | | 45 | No acute intrathoracic process |
| | | | | | | 3 | No acute intrathoracic process. _, MD |
| | | | | | | 3 | No acute intrathoracic process. Specifically, no pneumothorax. |
| | | | | | | 3 | No evidence of pneumonia. |
| 2 | Heart size is normal. The mediastinal and hilar contours are normal. The pulmonary vasculature is normal. Lungs are clear. No pleural effusion or pneumothorax is seen. There are no acute osseous abnormalities. | 1033 | 0.84 | 34 | 0.11 | 974 | No acute cardiopulmonary abnormality. |
| | | | | | | 24 | No evidence of pneumonia. |
| | | | | | | 3 | No radiographic evidence of pneumonia. |
| | | | | | | 2 | No acute cardiopulmonary abnormality. No displaced fracture identified. If there is continued concern for a rib fracture, consider a dedicated rib series. |
| | | | | | | 1 | Improving bibasilar atelectasis and decreasing bilateral effusions. |
| 3 | The lungs are clear without focal consolidation. No pleural effusion or pneumothorax is seen. The cardiac and mediastinal silhouettes are unremarkable. | 753 | 0.61 | 47 | 0.20 | 665 | No acute cardiopulmonary process. |
| | | | | | | 15 | No acute cardiopulmonary process. No focal consolidation to suggest pneumonia. |
| | | | | | | 8 | No pneumonia. |
| | | | | | | 7 | No evidence of pneumonia. No acute cardiopulmonary process. |
| | | | | | | 4 | No acute cardiopulmonary process. No significant interval change. |
| | *Duplicates with highest impression entropy* | | | | | | |
| 4 | The heart is normal in size. The mediastinal and hilar contours appear within normal limits. There is no pleural effusion or pneumothorax. The lungs appear clear. Bony structures appear within normal limits. | 25 | 0.02 | 2 | 0.99 | 14 | No evidence of acute cardiopulmonary disease. |
| | | | | | | 11 | No evidence of acute disease. |
| 5 | The lungs are clear. There is no pneumothorax. The heart and mediastinum are within normal limits. Regional bones and soft tissues are unremarkable. | 25 | 0.02 | 2 | 0.94 | 16 | Clear lungs with no evidence of pneumonia. |
| | | | | | | 9 | Clear lungs. |
| 6 | The lungs are well expanded and clear. Hila and cardiomediastinal contours and pleural surfaces are normal. | 23 | 0.02 | 15 | 0.92 | 6 | Normal. No evidence of pneumonia. |
| | | | | | | 2 | No evidence of pneumonia. |
| | | | | | | 2 | Normal chest radiograph. |
| | | | | | | 2 | No pneumonia. |
| | | | | | | 1 | Normal. No evidence of mass. |

Table 12: Examples of exact duplicates in the training set of MIMIC-CXR. In total, there are 14,596 reports with duplicated findings (11.9% of the training data). The table shows the number of reports with a given finding (**Dups.**), the relative frequency in the training set (**%**), the number of distinct impressions with this finding ($|y^*|$), the entropy over the impression frequencies (**H**), and the top-5 impressions with their respective **Count**.

| (Sub-)Category | Description | Example | Explanation | Count |
|---|---|---|---|---|
| *1. Incorrect findings: the finding in the reference is replaced with a different and incorrect finding.* | | | | 29 |
| Finding | incorrectness affects the main finding. | no acute intrathoracic process. | The reference uses "cardiopulmonary process" instead of "intrathoracic process". | 21 |
| Past state | incorrectness affects a past state of the patient. | increased opacity in the right lung.. | The reference mentions that the opacity is new and did not exist before. | 7 |
| Other | incorrectness affects other aspects. | bilateral pleural effusions,..., slightly improved... | The improvement is used to describe a second finding in the reference. | 1 |
| *2. Imprecise findings: the description of the finding or some of its aspects is imprecise compared to the reference.* | | | | 73 |
| Finding | the description of the finding itself is imprecise compared to the reference. | ...no displaced fractures are seen. | The reference uses "acute fractures" instead of "displaced fractures" (the reference is more general). | 21 |
| Location | the location of the finding is imprecise. | retrocardiac opacity compatible with pneumonia... | The references specifies the exact location: "Left lower lobe pneumonia". | 21 |
| Certainty | the summary is presented with a different degree of certainty. | bilateral middle lobe opacities could represent atelectasis or pneumonia. | The reference is certain about the finding. | 9 |
| Repetition | some findings are repeated. | unchanged bibasilar bronchiectasis and bibasilar bronchiectasis. | bibasilar bronchiectasis is mentioned twice. | 6 |
| Count | the count in the finding is imprise. | right pleural effusion. | The reference adds "Multiloculated", i.e., "Multiloculated right pleural effusion" | 2 |
| Size | the size of the finding is added/omitted/different. | multiple bilateral pulmonary nodules measuring up to 2. 5 cm. | The reference omits the size. | 1 |
| Other | other aspects about the finding are imprecise. | interval resolution of large right pleural effusion... | The reference includes other clinical information. | 13 |
| *3. Minor/secondary: errors that do not affect the finding.* | | | | 21 |
| Limitation | some limitations of the examination are (not) mentioned. | no definite acute cardiopulmonary process. | The reference adds "based on this limited, portable examination". | 15 |
| Phone calls | The time of a telephone call is different. | ...these findings were discussed with dr. _ by _ via telephone on _ at 4 : 45 pm. | The reference mentions a different time for the phone call. | 4 |
| Recommendation | errors related to recommendations. | short radiographic follow up is recommended within _ weeks to document resolution. | The reference omits "within _ weeks". | 2 |

Table 13: Bottom-up categorization of errors from the *Other* category with descriptions, examples and counts.

## D   Annotation Guidelines

**Introduction.**   We consider automatic impression generation for English radiology reports of chest imaging examinations. These reports conventionally have three sections (example in Figure 7).

1. **Background.** A description of the exam, patient information, and relevant prior exams.
2. **Findings.** A description or itemization of the radiologists' observations based on the radiographs.
3. **Impression.** A concise summary of the most important findings, including inferences and any recommendations.

```
                 FINAL REPORT
EXAMINATION: CHEST (PA AND LAT).

INDICATION: ____  year old woman with ?pleural effusion //
?pleural effusion

TECHNIQUE: Chest PA and lateral

COMPARISON: ____

FINDINGS:
Cardiac size cannot be evaluated.  Large left pleural effusion is
new. Small right effusion is new. The upper lungs are clear. Right
lower lobe opacities are better seen in prior CT. There is no pneu-
mothorax.  There are mild degenerative changes in the thoracic spine

IMPRESSION: Large left pleural effusion
```

Figure 7: Example radiology report.

**Study setup.**   We are looking to determine typical failures that automatic impression generation systems make. You will be shown a reference impression and four candidate impressions. The reference was written by a radiologist, and the candidates were generated by automatic systems. Your goal is to judge the accuracy of candidate impressions *based on a reference impression*. For each candidate, you will be asked to identify any errors that it may have.

**Definition of "error."**   We define an error as a difference between the candidate and reference. An error can be one of the following:

1. Omissions
   (a) Omission of finding/interpretation
   (b) Omission of comparison describing a change from a previous examination
   (c) Omission of reference to prior report while making a comparison
   (d) Omission of next steps, recommendation, communications

2. Additions
   (a) Additional finding/interpretation
   (b) Mention a comparison that is not present in reference
   (c) Additional reference to prior report while making a comparison
   (d) Additional mention of next steps, recommendation, communications
   (e) Additional finding/interpretation which contradicts reference
3. Incorrect location/position of finding
4. Incorrect severity of finding
5. Other difference between candidate and reference (please describe...)

Table 14 shows an example for each error category.

**Annotating errors as inline annotations.**   You are asked to annotate errors as *inline annotations*. For each error that you identify, first select the error category and then highlight the relevant text snippet with your mouse. This applies the category. If you have to remove an annotation, press on the highlight and use your backspace/delete key (⌦ or ←Backspace ). If one of the above categories occurs multiple times, please annotate all of them *separately* (see Figure 8). Some general guidelines:

- A candidate may have multiple errors, so please add all that apply.
- Some candidates will be the same, so please assign the same errors to all candidates.
- For additional findings that are plausible, pick 2a. Additional finding/interpretation. In the context of the full report, these additions may be correct. What this category aims to capture is that the system included information which the radiologist chose not to include. If a finding contradicts the reference, select 2e. Additional finding/interpretation which contradicts reference.
- Use 5. Other freely, especially if you find it difficult to assign any of the above categories. These remarks help us to better understand and characterize potential errors.
- You can ignore differences in word choice if they are synonymous. Example: *"may reflect developing consolidation"* is equal to *"could represent early consolidation."*

Finally, always use your best judgment when assessing the reports. If you are in doubt, you can add any questions/comments about the report or the error categories in the given box.

Figure 8: A candidate with two additional findings. Even though they are placed next to each other in the text, apply the category 2a. Additional finding/interpretation twice.

**Corner Cases**

**How to annotate "3. Incorrect location/position of finding" and "4. Incorrect severity of finding"?** Only apply if both reference and candidate mention a finding, *and* when there is a mismatch in severity/location. In the example below, both mention effusion, but the reference does not specify the size of effusion, whereas the candidate states that there are "small" effusions. Therefore, apply 4. Incorrect severity of finding.

> **Reference:** interval worsening of now moderate interstitial pulmonary edema. bilateral pleural effusions.
> **Candidate:** interval worsening of now moderate interstitial pulmonary edema. small bilateral pleural effusions.

**Opacities vs. consolidation.** Often, opacities are used in place of consolidation and vice versa. In those cases, apply 5. Other with a comment similar to "opacities not equal consolidation, but otherwise correct".

> **Reference:** Improved right lower medial lung peribronchial consolidation.
> **Candidate:** right lower medial lung peribronchial opacities have improved.

**No acute abnormality vs. COPD.** Does "no acute abnormality" contradict "COPD"? No, for the purposes of our evaluation, COPD is not an *acute* disease, so this is not contradicting. In the example below, following categories apply: (1) "COPD" is missing → 1a. Omission of finding/interpretation, (2) "opacity is resolved" → 1b. Omission of comparison describing a change from a previous examination, (3) "no acute cardiopulmonary abnormality" → 2a. Additional finding/interpretation.

> **Reference:** Left basilar opacity is resolved. COPD.
> **Candidate:** no acute cardiopulmonary abnormality.

**Misleading grammar or sentence structure.** In general, disregard grammatical errors. However, please pay attention to any *logical flaws* that arise because of grammar errors or a misleading sentence structure. In the example below, the *"and"* in the candidate implies that both "bronchiectasis" and "peribronchial consolidation" have improved, whereas the reference only states that the consolidation has improved. In those cases, apply 5. Other and add a comment similar to "logical error because of grammar."

> **Reference:** Bilateral lower lung bronchiectasis with improved peribronchial consolidation
> **Candidate:** bilateral lower lung bronchiectasis and peribronchial consolidation have improved since _.

194

| Error | Reference | Candidate | Explanation |
|---|---|---|---|
| *Omissions (apply to reference)* | | | |
| 1a. Omission of finding/interpretation | New left lower lobe infiltrate and effusion. | New left lower lobe infiltrate. | Effusion is missing. |
| 1b. Omission of comparison describing a change from a previous examination | In comparison to _ exam, there is interval near-complete resolution of bilateral pleural effusion. | No evidence of acute cardiopulmonary process. | Resolution of effusion is not described, therefore the comparison is missing. |
| 1c. Omission of reference to prior report while making a comparison | Increased pulmonary edema compared to _. | increased pulmonary edema. | While the candidate correctly states that the edema has increased, it lacks the reference to the prior report (or the date of it). |
| 1d. Omission of next steps / recommendation / communications | No pneumothorax or pneumomediastinum. Recommend repeat PA and lateral imaging later today to verify these findings. Otherwise unremarkable chest radiograph. These findings were communicated to Dr. _ at 11:55 a.m. by telephone by Dr. _. | No pneumothorax or pneumomediastinum. | The candidate does not include the followup (*recommend repeat PA*) and the remark about a communication with another doctor (*These findings were communicated [...]*). |
| *Additions (apply to candidate)* | | | |
| 2a. Additional finding / interpretation | Slight increased hazy opacities at the right lung base which may reflect developing consolidation. | slightly increased hazy opacities at the right lung base which may represent atelectasis or developing consolidation. | Atelectasis is not mentioned in the reference. This finding is not contradicting the reference. It may be correct in the context of the full report. *Same as 1a, but in the other direction.* |
| 2b. Mention a comparison that is not present in reference | Mild to moderate pulmonary edema, increased from _. | Mild to moderate pulmonary edema, increased from _. Stable cardiomegaly. | "Stable" suggests that the state of a finding was compared to a previous examination. This comparison is not made in the reference. *Same as 1b, but in the other direction.* |
| 2c. Additional reference to prior report while making a comparison | | | *Same as 1c, but in the other direction.* |
| 2d. Additional mention of next steps / recommendation / communications | | | *Same as 1d, but in the other direction.* |
| 2e. Additional finding / interpretation which contradicts reference | Unchanged size and position of right-sided hydropneumothorax. | Development of new right-sided hydropneumothorax | Unchanged vs. development of new |
| *Incorrect location, Incorrect Severity, Other* | | | |
| 3. Incorrect location/position of finding | New left lower lobe infiltrate | New right lower lobe infiltrate | Left vs. right |
| 4. Incorrect severity of finding | In comparison prior exam, there is near-complete resolution of bilateral pleural effusion | In comparison to _ exam, there is resolution of bilateral pleural effusion | Near complete vs. resolved |
| 5. Other | Slight increased hazy opacities at the right lung base which may reflect developing consolidation | Slight increased hazy opacity at the right lung base which may reflect developing consolidation | Difference in multiplicity |
| 5. Other | left picc terminates within the upper svc. | left picc terminates within the proximal svc. | Ambiguous location |
| 5. Other | No acute abnormalities identified to explain patient's cough and asthma flare. | no acute abnormalities identified to explain patient's cough. | Asthma flare is a symptom, which was not mentioned in the candidate. |

Table 14: Examples for all error categories.

# A Zero-Shot Approach for Multi-User Task-Oriented Dialog Generation

**Shiv Surya**      **Yohan Jo**      **Arijit Biswas**      **Alexandros Potamianos**

Amazon

`{shisurya,jyoha,barijit,potamian}@amazon.com`

## Abstract

Prior art investigating task-oriented dialog and automatic generation of such dialogs have focused on single-user dialogs between a single user and an agent. However, there is limited study on adapting such AI agents to multi-user conversations (involving multiple users and an agent). Multi-user conversations are richer than single-user conversations containing social banter and collaborative decision making. The most significant challenge impeding such studies is the lack of suitable multi-user task-oriented dialogs with annotations of user belief states and system actions. One potential solution is multi-user dialog generation from single-user data. Many single-user dialogs datasets already contain dialog state information (intents, slots), thus making them suitable candidates. In this work, we propose a novel approach for expanding single-user task-oriented dialogs (e.g. MultiWOZ) to multi-user dialogs in a zero-shot setting.

## 1 Introduction

As AI assistants are increasingly becoming a staple in our social environments like homes, offices and classrooms, it is imperative that these assistants are capable of comprehending conversations and understand requests in social settings. Consider the conversations in Figure 1 between human users and agents. The first conversation is an example of a single-user conversation, involving a single human user interacting with an AI agent. The second conversation is a conversation between two users that convey the same intent as the first conversation but in a social setting. In contrast with the single-user conversation, the multi-user conversation shows social banter and collaborative decision making arising from the deliberation of a suggestion between the human users, acceptance or rebuttal of suggestions and social banter. These social and conversational cues distract AI agents from tracking the task-relevant information most relevant to them.

---

**SINGLE-USER DIALOG**

**U:** I am looking for a place to dine. The restaurant should serve Mediterranean food and should be in the south.

**A:** I'm afraid there aren't any. Would you like to try another area or cuisine?

**U:** What restaurants are cheap?

**GENERATED MULTI-USER DIALOG**

**U:** I am looking for a place to dine. The restaurant should serve Mediterranean food and should be in the south.

**A:** I'm afraid there aren't any. Would you like to try another area or cuisine?

**U0:** Another area? I don't want to go to a place that isn't in the south. It's important to me that the restaurant be in the centre of the city. If you don't have a place in the south, I'll go to a restaurant in the centre. It'll be a new experience for me.

**U1:** Maybe you are right. Let me see what I can find. There are some places in the centre that I would like to try. If you can't find anything in the centre, then try the south. There are many good restaurants in the south.

**U0:** Yes, I would like a restaurant in the South. Cheap is fine.

Figure 1: An example of a multi-user dialog generated by our approach is shown here along with the reference single-user dialog. The single-user utterance shown in red is expanded as a multi-user dialog (also in red).

---

Therefore, it is important to develop AI agents that are robust to multi-user conversations.

The development of dialog systems that support multi-user conversations is hindered by the lack of suitable multi-user task-oriented dialog datasets with annotations of user belief states and system actions. Collecting such datasets is expensive and it is impractical to support all domains of interest (a typical digital assistant can support thousands of domains). Training people to conduct multi-user dialogs is nontrivial, and annotating the dialogs with dialog states is time-consuming.

To solve this problem, we propose a zero-shot approach for expanding single-user task-oriented dialogs to automatically generate multi-user task-oriented dialogs (Section 3). We do this with two main goals in mind: (1) leverage existing annotations of dialog states in the source single-user

dialogs, and (2) ensure a balance between the task-oriented tone and a social tone in the generated multi-user dialogs. To that end, our model expands each user utterance in a single-user dialog to a multi-user dialog that leads to the same user belief state as the source single-user utterance. Specifically, our model consists of two dialog generators and a turn planner. A task-oriented generator is trained on single-user task-oriented dialogs (e.g., MultiWOZ) responsible for generating utterances that contain task-relevant information (e.g. informed or requested slots), and a social generator trained on everyday dialogs (e.g., DailyDialog) responsible for generating utterance that express social dynamics of the speakers. We fuse these dialog generators to generate utterances in a multi-user dialog that is conversational, social, and relevant to the task at hand. The turn planner decides on an appropriate mix of dialog generators for generating a particular utterance in multi-user dialog.

Our automatic and human evaluation (Section 4) shows that our approach generates multi-user dialogs that contain social chatter consistent with dialog history and reflect user belief states consistent with source utterances. The main contributions of our work are:

- This is, to our knowledge, the first approach for generating multi-user dialogs from single-user task-oriented dialogs in a zero-shot setting.

- We show that the fusion of dialog generators is effective in maintaining topic relevance of social turns and improving quality of multi-user dialogs.

- We propose rule and model-based turn planners that select an appropriate combination of dialog generators and generate a multi-user dialog.

- Automated metrics and human quality evaluations show that dialog generator fusion and turn planning improves conversational fluidity in multi-user dialogs. It also better reflects the user belief.

## 2 Related Work

While there are public datasets of task-oriented dialogs annotated with dialog states (Andreas et al., 2020; Byrne et al., 2019; Rastogi et al., 2020; Zhu et al., 2020), most of them (if not all) are focused on single-user transactions, i.e., an agent converses with one user at a time. As a result, dialogs in these datasets do not reflect important dynamics of users making decisions together while interacting with an agent. Some dialog datasets cover such dynamics, like social banter (Li et al., 2017) and deliberation (Karadzhov et al., 2021; He et al., 2018), but they are not task-oriented or annotated with dialog states important for training dialog systems. While the FusedChat data (Young et al., 2022) contains task-oriented dialogs with social chatter, the dialogs are still single-user. By contrast, our goal is to build multi-user task-oriented dialogs reflecting social dynamics simultaneously.

The most similar line of work to ours is dialog generation from summaries, as a means for data augmentation for dialog summarization. In one framework (Gunasekara et al., 2021), a conversation generator (a seq2seq model) generates a conversation, and a summary generator (a seq2seq model) summarizes the generated conversation. The similarity between the generated summary and the original summary is used as a reward to inform the conversation generator in a reinforcement learning fashion. Another approach is to take a seed pair of a dialog and its summary, and gradually modify them iteratively (Liu et al., 2022). Specifically, an utterance of the dialog is replaced with a new utterance that is generated based on its context by a seq2seq model. Next, the summary is updated based on the updated dialog. Data augmentation using these methods improves dialog summarization accuracy in few-shot settings. The main difference between these approaches and ours is that they require seed dialogs. By contrast, we assume the more challenging scenario of zero-shot dialog generation, where we do not have enough dialogs to start with. This setting is more realistic because it is impractical to collect seed conversations of more than 100K domains in case of popular voice assistants.

Some studies have addressed generation of single-user task-oriented dialogs. One approach is to iteratively generate a user utterance and predict the resulting user belief state using two models (Kim et al., 2021). Another approach is to use a variational hierarchical dialog autoencoder that generates dialogs and their underlying dialog states simultaneously (Yoo et al., 2020). Importantly, these studies still tackle single-user dialogs,

**Prompt**

```
<belief-state>
attraction area: centre, attraction type: entertainment
<dialog-history>
u_1: Why don't you do something fun?
...
u_(n-1): I don't know what to do. It's so boring
```

**3.1 Dialog Generators**

Task-Oriented Dialog Generator

Social Dialog Generator

**3.3 Turn Planning**

Turn Planner

**(Coarse Plan)**
Task-Oriented?
Social?

**(Fine-Grained Plan)**
Task-Oriented?
Social-Inform?
Social-Commissive?
Social-Directive?
Social-Question?

**3.2 Late Fusion**

Late Fusion

**Output**

```
<belief-state>
attraction area: centre, attraction type: entertainment
<dialog-history>
u_1: Why don't you do something fun?
...
u_(n-1): I don't know what to do. It's so boring
u_n: How about we try to find something entertaining in the centre of town.
```
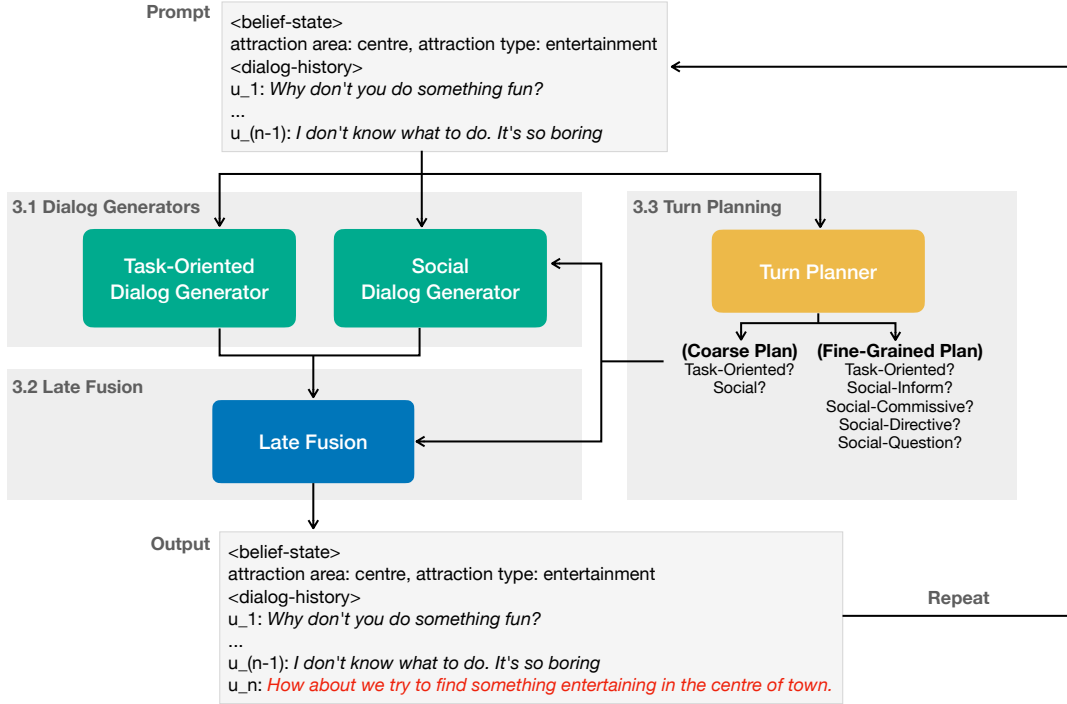
**Repeat**

Figure 2: Architecture of the proposed multi-user dialog generator is shown here.

whereas our task of multi-user dialog generation is new and arguably more challenging as it requires resulting dialogs to reflect the social dynamics of decision-making between multiple users.

## 3 Our Approach: Multi-user Dialog Generation

Our proposed approach automatically generates multi-user task-oriented dialog. The utterances in a multi-user task-oriented dialog can be task-oriented, social or a mix of both. Training a single dialog generator capable of handling different kind of dialog turns e.g social or task-oriented requires an aligned dataset. In the absence of such a dataset, we use a bottom up approach of combining dialog generators capable of generating different kind of utterances (e.g social/task-oriented) to generate a multi-user dialog turn-by-turn. The architecture of this approach is shown in Figure 2. We train a task-oriented generator on single-user task-oriented dialogs responsible for generating task-oriented utterance, and a social generator trained on dialogs from everyday scenarios for generating social utterances (Section 3.1). We control turn level dynamics by training a turn planner. We use a turn planner (Section 3.3) that switches between the two dialog generators utterance-by-utterance to generate a multi-user rewrite. Additionally, we propose a

mechanism to blend generations from the two dialog generators (Section 3.2).

### 3.1 Dialog Generators

Our model consists of two dialog generators. The first of these dialog generators is a task-oriented generator trained on single-user task-oriented dialogs (e.g., MultiWOZ). This generator is responsible for generating task-oriented utterances. We also train a social generator trained on dialogs from everyday scenarios (e.g., DailyDialog) capable of generating social utterances. We use a combination of these dialog generators to generate the utterances in a multi-user dialog. Training and inference details of dialog generators are provided in Appendix A.

#### 3.1.1 Task-Oriented Dialog Generator

The task-oriented dialog generator is responsible for generating task-oriented utterances. This dialog generator is trained to generate an utterance that reflects a user belief state when prompted with the relevant dialog history between user and agent along with the user belief states.

Let $U = \{u_1, \ldots u_n\}$ be a sequence of dialog turns consisting of both user and agent utterances. For a given user utterance $r = u_j, j \in \{1, ..., n\}$ we define $X^h = \{x_1^h, ..., x_{|h|}^h\}$ as the tokens that make up the dialog history of utterances $h = u_{<j}$.

Let $X^b = \{x_1^b, ..., x_{|b|}^b\}$ be the tokens that make up the user belief state $b$ corresponding to user utterance $r$ and $X^r = \{x_1^r, ..., x_{|r|}^r\}$ be the tokens that make up target utterance $r$.

We train the task-oriented dialog generator to produce the tokens $X^r$ in target utterance $r$ while being prompted with the concatenation of tokens in user belief state $b$ and dialog history $h$, $[X^b, X^h]$. Let $X = \{x_1, x_2, ..., x_t, ..., x_{|h|+|b|+|r|}\}$ be the tokens in concatenation of user belief state $b$, dialog history $h$ and target utterance $r$. The task-oriented dialog generator is trained to maximize the conditional log-likelihood of predicting the target tokens $x_t$ for all $t > |h| + |b|$ to achieve this goal.

$$\mathcal{L} = -\sum_{t=|h|+|b|+1}^{|h|+|b|+|r|} \log P_{task}(x_t|x_{<t}) \quad (1)$$

where $|h|, |b|$ and $|r|$ are the number of tokens in dialog history, belief state and target respectively. $P_{task}$ is the likelihood of the task-oriented dialog generator predicting a token given prompt and generated tokens as context.

### 3.1.2 Social Dialog Generator

The social dialog generator is responsible for generating utterances in social chit-chat. This dialog generator is trained to generate an utterance that best continues a dialog. It is prompted with the relevant dialog history between two users and the corresponding social act while generating this utterance. Given a sequence of dialog turns $U = \{u_1, \ldots u_n\}$ consisting of utterances in a dialog from everyday conversations. For a given utterance $r = u_j, j \in \{1, ..., n\}$, we define $X^h = \{x_1^h, ..., x_{|h|}^h\}$ as the tokens that make up the dialog history of utterances $h = u_{<j}$. $X^h$ can optionally include the social act $\alpha$ for $r$. $X^r = \{x_1^r, ..., x_{|r|}^r\}$ are the tokens that make up target utterance $r$.

The social dialog generator is trained to produce the tokens $X^r$ in target utterance $r$ given the tokens in $h$, $[X^h]$ as input. Let $X = \{x_1, x_2, ..., x_t, ..., x_{|h|+|r|}\}$ be the tokens in concatenation of dialog history $h$ and target utterance $r$. To achieve this goal, the social dialog generator is trained to to maximize the conditional log-likelihood of predicting the target tokens $x_t$ for all

$t > |h|$ in $X$.

$$\mathcal{L} = -\sum_{t=|h|+1}^{|h|+|r|} \log P_{social}(x_t|x_{<t}) \quad (2)$$

where $|h|$ and $|r|$ are the number of tokens in dialog history, and target respectively. $P_{social}$ is the likelihood of the social dialog generator predicting a token at time step $t$ given the prompt and generated token history as context.

### 3.2 Late Fusion of Dialog Generators

We propose a mechanism to blend generations between task-oriented and social dialog generators. Utterances generated by the task-oriented dialog generator can be command like. Similarly, utterances generated by the social dialog generator can be unrelated to the task in the utterance. We hypothesize that fusing the generators would allow us to achieve more conversational task-oriented dialog and increase the topic relevance of social dialog turns. Formally, at each step in generating a token from our dialog generator, we combine the token probabilities $P_{task}$ and $P_{social}$ from the task-oriented generator and social generator respectively. Formally, this is defined as:

$$p(x_t|x_{<t}) = \beta \cdot P_{task}(x_t|x_{<t}) + \\ (1-\beta) \cdot P_{social}(x_t|x_{<t}) \quad (3)$$

where $\beta \in [0, 1]$ is a scalar blending factor.

The blending factor allows us to bias the decoding of the utterance towards the social dialog generator or the task-oriented dialog generator. This allows us to generate dialog that is a mix of both social and task-oriented dialog. Thus, we can achieve more natural transitions between turns and lessen abrupt topic changes at an utterance level in the multi-user dialog generation. Examples of how late fusion of dialog generators affects dialog generated are provided in Appendix A.

### 3.3 Turn Planning

We train a turn planner to predict the type of utterance we expect in the multi-user dialog and select the appropriate combination of dialog generators for use in generating a particular utterance. This model based planner predicts the type of utterance (social/task) we use in the next utterance of multi-user dialog, based on history of utterances in the dialog.

We design two variants of the turn planner. The first version, a coarse turn planner, is capable of predicting turn type as social or task. The fine-grained variant of the turn planner is trained to predict social acts of the utterance along with the type of turn. Training and inference details of the turn planners are provided in Appendix A.

### 3.3.1 Coarse Turn Planner

Given previous dialog turns as context, the coarse planner is trained to predict the turn type of the next turn in the dialog as social or task. The coarse turn planner is a pre-trained BERT model (Devlin et al., 2019) tuned in a binary classification setting. Following BERT, we use a fully-connected layer over the final hidden representation corresponding to the special classification token ([CLS]). The coarse turn planner is trained on as subset of the FusedChat dataset (Young et al., 2022).

### 3.3.2 Fine-Grained Turn Planner

The fine-grained turn planner is a classifier trained to predict turn type of the next utterance as social or task. In addition, it also predicts the social act of this utterance. DailyDialog dataset (Li et al., 2017) has annotations for the social act $\alpha \in \{$inform, commissive, directive, question$\}$ for each utterance. This social act $\alpha$ is given as an additional input prompt during finetuning of the social dialog generator. Specifically, the input prompt has the corresponding social act prepended at the beginning. This allows the social dialog generator to learn to maximize $P_{social}(x_t|x_{<t}, \alpha)$ enabling finer control of the social dialog generator via these social dialog acts.

Similar to the setting of the coarse turn planner, the fine-grained turn planner is a BERT classifier. However, this classifier is trained in a multi-label setting. The multi-label setting allows turn type prediction(social/task) along with prediction of the social act. However, annotations for these social acts are not present in the FusedChat dataset used to train our model turn planner. To circumvent this lack of annotations,we label each utterance in the FusedChat data with social acts using distance supervision as follows: (1) First, we train a multi-class BERT classifier on DailyDialog dataset with labels for each type of utterance. (2) We then use this trained multi-class classifier to label every utterance in FusedChat. This propagates distant labels of the utterance type to social and task utterances in FusedChat. The distant labels allow us to train the fine-grained turn planner on a subset of FusedChat to jointly predict the nature of the next turn in the dialog as social or task along with the type of utterance in the next turn.

### 3.3.3 Rule Planner

As a baseline, we design a rule based turn planner. The rule turn planner predicts a random number of social turns followed by a single task turn. In such a plan, each social and task turn is generated by the social or task generator respectively. While rule planners could append social turns, we only consider the setting in which social turns are prepended to obtain a setting similar to the user dynamics modeled by our model based turn planners. This enables a fairer comparison to the settings of the trained model planners in Section 3.3.1–3.3.2.

## 4 Experiments

We compare four variants of our proposed approach for multi-user dialog generations in our evaluations. These variants ablate the effect of (1) late fusion (2) type of turn planner. The variants are described below:

- **Rule Planner (RTP)**: This approach makes a hard choice between the social and task-oriented generator while generating a multi-user dialog. The type of utterance at each turn is determined by a rule planner (see Section 3.3.3)

- **Rule Planner w/ Late Fusion (RTP+LF)**: This approach blends the social and task-oriented generator using late fusion(see Section 3.2) while generating a multi-user dialog. The type of utterance at each turn is determined by a rule planner (see Section 3.3.3).

- **Coarse Planner w/ Late Fusion (CTP)**: This approach blends the social and task-oriented generator using late fusion(see Section 3.2) while generating a multi-user dialog. The type of utterance at each turn is determined by a coarse planner (see Section 3.3.1).

- **Fine-grained Planner w/ Late Fusion (FTP)**: This approach blends the social and task-oriented generator using late fusion(see Section 3.2) while generating a multi-user dialog. The type of utterance and social act at each turn is determined by a fine-grained model planner(see Section 3.3.2).

### 4.1 Datasets

#### 4.1.1 MultiWOZ

We use this dataset as our primary task-oriented dialog dataset for generating and evaluating multi-user dialog generations. It is also used for training the task-oriented dialog generator. This dataset (Zang et al., 2020) contains multi-turn dialog between a single user and agent spanning multiple task domains. The dataset consists of 8,438 training dialogs and 1,000 validation and test dialogs each. In our experiments, we exclude the police, bus and taxi domains from the training set as they have very few dialogs in training.

#### 4.1.2 DailyDialog

This is a corpus (Li et al., 2017) containing dialogs centered around daily life communications written by humans. We use the DailyDialog twice in our proposed approach. Firstly, it is used to train the social dialog generator. It is also used to train a classifier that labels utterances in FusedChat with social acts using distance supervision. The dataset contains 13,118 multi-turn dialogs spanning 10 daily life topics. This dataset is of particular interest to us as it has annotations for social act at an utterance level. We infer speaker turns based on utterance turns as the corpus does not explicitly indicate speaker information.

#### 4.1.3 FusedChat

We use FusedChat to train the coarse and fine-grained turn planners. This is a dataset (Young et al., 2022) based on MultiWOZ (Zang et al., 2020). This work adds expanded social turns between a single user and agent. The dataset consists of MultiWOZ dialogs with prepended and appended social turns. We only consider the subset of FusedChat with social turns prepended to task-oriented turns. In this setting, the intent of the task-oriented turn following the social turns is strictly dependent on the topic of conversation in the social turns. This results in 3670 training dialogs and 500 validation and testing dialogs each.

### 4.2 Evaluation Measures

#### 4.2.1 Automatic evaluation

Our multi-user rewrite of a single user task-oriented dialog should reflect the same user intent as the single user utterance, reliably cover the user belief states expressed in user utterance being rewritten while exhibiting high lexical diversity. We define the following automated metrics targeted at measuring semantic similarity between source utterance and utterances in multi-user rewrite and lexical diversity.

- **Semantic Similarity (SS)**: For a multi-turn rewrite of a single utterance, the cosine similarity of sentence embedding from transformers (Reimers and Gurevych, 2019) is computed between the average embedding of all generated utterances and the single-user utterance being rewritten. This measures the topical alignment of the multi-user dialog with the single user utterance.

- **Slot Coverage (SC)**: We measure the fraction of user belief states expressed in single-user utterance being mentioned in the multi-turn rewrite via a sub-string match. This metric captures the coverage of user belief states from user utterance in the generated multi-user rewrite.

- Additionally, we measure **Distinct-n/D-n (n=1,2,3)** (Li et al., 2015) and **Repetition-n/R-n (n=3)** which measures the percentage of unique n-grams and the fraction of copied 3-gram from prompt in a multi-user rewrite. These metrics measure lexical diversity in the generated multi-user dialog.

#### 4.2.2 Human Evaluation

We conduct a human evaluation on Amazon Mechanical Turk. For this, we randomly sample 150 single-user utterances and corresponding generated multi-user dialogs from our proposed approach. We compare 4 variants of our proposed approach in a staggered manner which results in 450 multi-user rewrites. We conduct pairwise comparisons of generated multi-user rewrites, and each pair was evaluated by three human labelers. We ask humans to pick the better dialog in a pairwise evaluation by assessing:

- **Content Quality**: Consider the fluency, grammatically correctness, logical coherence and utterance ordering to determine the better multi-user dialog.

- **Belief State faithfulness**: Choose the multi-user dialog generation system that better reflects user beliefs expressed in the single user utterance.

**DIALOG HISTORY:**                                             **D.1**

**U:** Your reservation has been booked. Your reference number is Q9DRQ55J . Please be aware that your table will only be reserved until 16:30. Can i provide any more assistance?
**A:** I'm also looking for things to do in the same general area as the restaurant.
**U:** Sure, are you interested in museums, sports, parks, entertainment, swimming or boating?

**USER BELIEF STATE:**

attraction area: east, attraction type: cinema

**UTTERANCE BEING REWRITTEN:**

Are there any cinemas?If there isn't what about something that has to do with multiple sports?

**GENERATED MULTI-USER DIALOG:**
**U0:** I think I'd like to go to the cinema. What do you think of this one? It's a new cinema, isn't it?
**U1:** Yes. It is a good cinema. You can go to the movies. It's a new cinema, isn't it?
**U0:** Yes it is. Is there a cinema in the east?

**DIALOG HISTORY:**                                             **D.2**
**U:** Please see what you can find out about jesus college
**A:** It is a college in the town centre, address Jesus Lane. Postcode is cb58bl and it has free admission.
**U:** Do you happen to have the telephone number for the college?I would like to call the main office there.
**A:** I can certainly look that up for you. Ok, it's 01223339485. Do you need any more information or something else I can help with?

**USER BELIEF STATE:**

train day: tuesday, train destination: stansted airport

**UTTERANCE BEING REWRITTEN:**

Thank you.I'm looking to leave on Tuesday by train. Arrival will be Stansted Airport.

**GENERATED MULTI-USER DIALOG:**
**U0:** Yes. The college has a large population. It is a very good college. It has a lot of students. It is very easy to  get into the college. It is in the centre of town. It is not very expensive to get into college.It's very easy to find out about it.Thank you very much. Are there any other places I can get to  for my college search?
**U1:** There are many places. For example. You can go to the college website. It has lots of information. And also there are many places to get a map.
**U0:** Can you help me find a train to the college on Tuesday?

**DIALOG HISTORY:**                                             **D.3**

A: Addenbrookes Hospital Telephone : 01223245151 Address : Hills Rd, Cambridge Postcode : CB20QQ Do you need a taxi to get there?
U: No, thank you. Can you also find me a restaurant that serves British food in the centre? Postcode, please.
A: Cotto would fit your needs perfectly. Their postcode is CB11BG. Can I help you with anything else?
U: Are there any theatres near the Cotto?
A: Adc theatre is nearby. it is in park street

**USER BELIEF STATE:**

attraction area: centre , attraction type: theatre, attraction address: requested, attraction phone: requested

**UTTERANCE BEING REWRITTEN:**

Thank you. I'm also looking to leave on Tuesday by train. Arrival will be Stansted Airport.

**GENERATED MULTI-USER DIALOG WITH COARSE TURN PLANNER:**

**U0:** Great. Can I get the phone number and address for the theatre? I am also looking for a restaurant in the centre. Thank you.
**U1:** You can call the theatre from here. And I can give you the address and phone number for the restaurant. You can book it from here. It is very easy to do. And I will be happy to help you with it. You are welcome. Goodbye.
**U0:** Can I have the address and telephone number for the theatre, please?

**GENERATED MULTI-USER DIALOG WITH FINE-GRAINED TURN PLANNER:**
**U0:** Thank you. Please get me the address and phone number for the theatre. I will be happy to go there.
**U1:** Cotto is in the centre of town. The theatre is in the same street as the restaurant. So you can walk there.
**U0:** That's great. Can I get the address and telephone number for the theater?

Figure 3: Examples of multi-user task-oriented dialogs generated by our approach is shown here.

For each pairwise setting, we compute the majority vote based on the forced pairwise comparison. We use the Sign Test (Dixon and Mood, 1946) to compute statistical significance for both evaluation criteria. More details of the human evaluation including a screenshot of the evaluation template is available in Section B of the Appendix.

### 4.3 Results

Automatic evaluation of variations of our proposed approach are available in Table 1. Overall, using both late fusion and a turn planner achieves better performance. In Row 1, we see that the rule planner with hard choice of generators struggles to maintain relevance to intent in the single-user

| Turn Planner | LF | SC↑ | SS↑ | D-1↑ | D-2↑ | D-3↑ | R-3↓ |
|---|---|---|---|---|---|---|---|
| RTP | ✗ | 0.358 | 0.386 | **0.797** | **0.960** | **0.986** | 0.013 |
| RTP | ✓ | 0.446 | **0.478** | 0.658 | 0.898 | 0.963 | 0.012 |
| CTP | ✓ | **0.480** | 0.461 | 0.622 | 0.883 | 0.957 | 0.012 |
| FTP | ✓ | 0.464 | 0.455 | 0.747 | 0.939 | 0.980 | **0.011** |

Table 1: Automatic quality metrics on the test set to ablate the effect of turn planner type and fusion of dialog generators. *Legend - ↑: Higher is better, ↓: Lower is better, RTP: Rule Turn Planner, LF: Late Fusion, CTP: Coarse Model Turn Planner, FTP: Fine-Grained Model Turn Planner, SC: Slot Coverage, SS: Semantic Similarity.*

utterance. This is indicated by much lower seman-

| | Win % | | |
|---|---|---|---|
| Evaluation | RTP | RTP+LF | Sign Test (p<0.05) |
| CQ | 42 | 58 | ✓ |
| BF | 38 | 62 | ✓ |
| Evaluation | RTP+LF | CTP | Sign Test (p<0.05) |
| CQ | 52 | 48 | ✗ |
| BF | 54.66 | 45.33 | ✗ |
| Evaluation | CTP | FTP | Sign Test (p<0.05) |
| CQ | 62.66 | 37.33 | ✓ |
| BF | 61.34 | 38.66 | ✓ |

Table 2: Human evaluation on samples from the test set in a forced choice pairwise evaluation. Win% = % times multi-user dialogs from one model was preferred over the other when evaluated against a particular criterion. *Table Legend - RTP: Rule Turn Planner, LF: Late Fusion, CTP: Coarse Model Turn Planner, FTP: Fine-Grained Model Turn Planner, CQ: Content Quality, BF: Belief State Faithfulness*

tic similarity and slot coverage scores along with high n-gram diversity metrics. The addition of late fusion (RTP vs RTP+LF) produces a significant jump in both semantic similarity and slot coverage of the multi-user dialog when compared to the hard choice of generators rule based turn planner. This indicates better relevance to intent in the single-user utterance across the generated multi-user rewrite with late fusion. We also see a reduction in n-gram diversity metrics. This is expected as n-gram diversity would reduce when the social turns are also related to same topic. The replacement of rule turn planner by a coarse turn planner (RTL+FL vs CTP) produces improvements in slot coverage. The fine-grained planner (FTP) gets comparable semantic similarity and slot coverage with higher lexical diversity scores. Holistically, this is indicative of the fine-grained planner showing comparable faithfulness in reflecting the user belief state, while reducing repetition across utterances in the multi-user dialog.

Results of our human evaluation are available in Table 2. From the pairwise evaluation, we see that late fusion outperforms a hard choice of the backbone generators with statistical significance. This is also observed in the automatic evaluation where improved semantic similarity scores and slot coverage indicate that late fusion produces dialog turns that are more related to intent expressed in single-user utterance. The coarse model turn planner (CTP) and rule turn planner with late fusion (RTL+FL) are tied without a statistically significant result on both criteria. This result is aligned with

close automatic metric for these models observed in Table 1. This shows the limitations of training turn planners in low data regimes. We expect the performance of our approach with the coarse and fine-grained model planners to improve with training on larger datasets and data augmentation.

Some examples of generations from our approach are shown in Figure 3. The generations show reasonable faithfulness to the intent and user belief in the single user utterance being rewritten across D.1–D.3. However, in example D.1, we see repetition of the phrase "It's a new cinema, isn't it?" across utterances by different users indicating the challenge of consistent Point-of-View (PoV) depiction. In D.2, we find that the role and characteristics of an agent bleeds into the users engaged in a conversation with command like responses despite blending of dialog generators. Example D.3 contrasts multi-user dialog generation with coarse and fine-grained turn planners. Here we see that while both generations cover the user beliefs, the dialog generated with the fine-grained turn planner is more coherent across turns. The coarse-grained planner shows content repetition across turns. Additional examples are shown in Section C of the appendix.

### 4.4 Challenges and Future work

Despite these promising results, we find that there are considerable challenges to be tackled. Our ability to control the social utterances in the multi-user rewrite is limited to broad social acts, i.e., question or inform. This limitation arises from the dataset we use for fine-tuning the social dialog generator. Using datasets with fine-grained annotations for utterance type or larger language models capable of instruction prompting are potential directions to address this. Another challenge is ensuring consistency of user beliefs across social utterances. We observe challenges in maintaining consistent beliefs across utterances in a dialog sequence for a user. Further, maintaining consistent Point-of-View (PoV) depiction is challenging. We find that the role and characteristics of an agent bleeds into the users engaged in a conversation. Planning approaches like those employed in story and long text generation (Rashkin et al., 2020; Yao et al., 2019) is one possible family of approaches that could reduce inconsistency in user behaviour, PoV and provide more control.

## Ethics Statement

Advances in multi-user dialog generation techniques would aid training of digital assistants. As AI assistants are increasingly becoming a staple in our social environments, synthetic methods of multi-user dialog generation would aid the training of these assistants and ensure they are capable of comprehending human conversations and understand task-oriented requests in social settings. This would help increase human-machine interaction and enhance human productivity in collaborative settings.

Synthetic multi-user dialog generation techniques would also reduce the need of (the gold standard for data collection) crowdsourcing. This would also have a positive effect on human productivity and reduce the need for humans to manually write dialogs for different scenarios.

We use language models as initialisation for our dialog generators. These are trained on data collected from the web. Hence, issues related to bias and abusive language are a potential concern. These concerns of abusive content should be largely mitigated as we fine-tune of the dialog generators on task-oriented and everyday conversation datasets with sanitised data. The generator fine-tuning and prompt structure used for dialog generation should limit unintended consequences as all generations are trained to reflect the intent of the single-user dialog. However, with our proposed method of multi-user dialog generation, any racial, ethnic or other forms of bias present in the datasets used to train the dialog generators is likely to get propagated to the generated multi-user dialog.

## References

Jacob Andreas, John Bufe, David Burkett, Charles Chen, Josh Clausman, Jean Crawford, Kate Crim, Jordan DeLoach, Leah Dorner, Jason Eisner, Hao Fang, Alan Guo, David Hall, Kristin Hayes, Kellie Hill, Diana Ho, Wendy Iwaszuk, Smriti Jha, Dan Klein, Jayant Krishnamurthy, Theo Lanman, Percy Liang, Christopher H Lin, Ilya Lintsbakh, Andy Mc-Govern, Aleksandr Nisnevich, Adam Pauls, Dmitrij Petters, Brent Read, Dan Roth, Subhro Roy, Jesse Rusak, Beth Short, Div Slomin, Ben Snyder, Stephon Striplin, Yu Su, Zachary Tellman, Sam Thomson, Andrei Vorobev, Izabela Witoszko, Jason Wolfe, Abby Wray, Yuchen Zhang, and Alexander Zotov. 2020. Task-Oriented Dialogue as Dataflow Synthesis. *Transactions of the Association for Computational Linguistics*, 8:556–571.

Bill Byrne, Karthik Krishnamoorthi, Chinnadhurai Sankar, Arvind Neelakantan, Ben Goodrich, Daniel Duckworth, Semih Yavuz, Amit Dubey, Kyu-Young Kim, and Andy Cedilnik. 2019. Taskmaster-1: Toward a Realistic and Diverse Dialog Dataset. *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4515–4524.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Wilfrid J Dixon and Alexander M Mood. 1946. The statistical sign test. *Journal of the American Statistical Association*, 41(236):557–566.

Chulaka Gunasekara, Guy Feigenblat, Benjamin Sznajder, Sachindra Joshi, and David Konopnicki. 2021. Summary Grounded Conversation Generation. *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3748–3756.

He He, Derek Chen, Anusha Balakrishnan, and Percy Liang. 2018. Decoupling Strategy and Generation in Negotiation Dialogues. *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2333–2343.

Georgi Karadzhov, Tom Stafford, and Andreas Vlachos. 2021. DeliData: A dataset for deliberation in multi-party problem solving. *arXiv*.

Sungdong Kim, Minsuk Chang, and Sang-Woo Lee. 2021. NeuralWOZ: Learning to Collect Task-Oriented Dialogue via Model-Based Simulation. *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3704–3717.

Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2015. A diversity-promoting objective function for neural conversation models. *arXiv preprint arXiv:1510.03055*.

Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. 2017. DailyDialog: A manually labelled multi-turn dialogue dataset. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 986–995, Taipei, Taiwan. Asian Federation of Natural Language Processing.

Yongtai Liu, Joshua Maynez, Gonçalo Simões, and Shashi Narayan. 2022. Data Augmentation for Low-Resource Dialogue Summarization. *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 703–710.

Hannah Rashkin, Asli Celikyilmaz, Yejin Choi, and Jianfeng Gao. 2020. PlotMachines: Outline-conditioned generation with dynamic plot state tracking. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4274–4295, Online. Association for Computational Linguistics.

Abhinav Rastogi, Xiaoxue Zang, Srinivas Sunkara, Raghav Gupta, and Pranav Khaitan. 2020. Towards Scalable Multi-Domain Conversational Agents: The Schema-Guided Dialogue Dataset. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):8689–8696.

Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics.

Lili Yao, Nanyun Peng, Ralph Weischedel, Kevin Knight, Dongyan Zhao, and Rui Yan. 2019. Plan-and-write: Towards better automatic storytelling. In *AAAI*, pages 7378–7385.

Kang Min Yoo, Hanbit Lee, Franck Dernoncourt, Trung Bui, Walter Chang, and Sang-goo Lee. 2020. Variational Hierarchical Dialog Autoencoder for Dialog State Tracking Data Augmentation. *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3406–3425.

Tom Young, Frank Xing, Vlad Pandelea, Jinjie Ni, and Erik Cambria. 2022. Fusing task-oriented and open-domain dialogues in conversational agents. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(10):11622–11629.

Xiaoxue Zang, Abhinav Rastogi, Srinivas Sunkara, Raghav Gupta, Jianguo Zhang, and Jindong Chen. 2020. MultiWOZ 2.2 : A dialogue dataset with additional annotation corrections and state tracking baselines. In *Proceedings of the 2nd Workshop on Natural Language Processing for Conversational AI*, pages 109–117, Online. Association for Computational Linguistics.

Qi Zhu, Kaili Huang, Zheng Zhang, Xiaoyan Zhu, and Minlie Huang. 2020. CrossWOZ: A Large-Scale Chinese Cross-Domain Task-Oriented Dialogue Dataset. *Transactions of the Association for Computational Linguistics*, 8:281–295.

# Beyond the Bias: Unveiling the Quality of Implicit Causality Prompt Continuations in Language Models

**Judith Sieker** and **Oliver Bott** and **Torgrim Solstad** and **Sina Zarrieß**

Bielefeld University

{j.sieker, oliver.bott, torgrim.solstad, sina.zarriess}@uni-bielefeld.de

## Abstract

Recent studies have used human continuations of Implicit Causality (IC) prompts collected in linguistic experiments to evaluate discourse understanding in large language models (LLMs), focusing on the well-known IC coreference bias in the LLMs' predictions of the next word following the prompt. In this study, we investigate how continuations of IC prompts can be used to evaluate the text generation capabilities of LLMs in a linguistically controlled setting. We conduct an experiment using two open-source GPT-based models, employing human evaluation to assess different aspects of continuation quality. Our findings show that LLMs struggle in particular with generating coherent continuations in this rather simple setting, indicating a lack of discourse knowledge beyond the well-known IC bias. Our results also suggest that a bias congruent continuation does not necessarily equate to a higher continuation quality. Furthermore, our study draws upon insights from the Uniform Information Density hypothesis, testing different prompt modifications and decoding procedures and showing that sampling-based methods are particularly sensitive to the information density of the prompts.

## 1 Introduction

There is currently a growing interest in probing the performance of large language models (LLMs) on carefully controlled linguistic test suites and experimental datasets to get a deeper understanding of specific linguistic capabilities captured in these models (e.g., Belinkov and Glass, 2019; Ettinger, 2020). While a lot of previous work focused on analyzing the syntactic competence of LLMs (e.g, Hu et al., 2020; Schuster and Linzen, 2022), recent studies also started to investigate the abilities of LLMs on the level of semantics and pragmatic discourse processing. One promising diagnostic for probing discourse knowledge in LLMs has turned out to be the use of Implicit Causality (IC) prompts.

IC refers to a property of a broad range of interpersonal verbs that exhibit strong preferences for establishing coreference to one of the verb's arguments over the other in explanations. For instance, when asked to provide a continuation after "..." in a sentence like (1), humans display strong next-mention preferences towards the stimulus (*he/Tom* in this case):

(1)     Tom fascinated Sarah because... *he was very smart.*

As the IC bias has been extensively researched in psycholinguistics and psychology across various languages and populations (e.g., Ferstl et al., 2011; Hartshorne et al., 2013; Bott and Solstad, 2014), investigating this bias in LLMs has gained significant interest. A range of recent studies investigated LLMs' predictions of the next mention in examples like (1) and whether these mentions (i.e. pronouns) follow the same coreference biases as can be found in human data (e.g., Upadhye et al., 2020; Davis and van Schijndel, 2020; Kementchedjhieva et al., 2021; Zarrieß et al., 2022). These studies predominantly indicated that LLMs are not generally congruent with the human IC bias, which has been interpreted as evidence for LLMs struggling with certain aspects of discourse understanding (but see Cai et al., 2023).

In this work, we propose that experimentally elicited data of human continuations of IC prompts cannot only be used for analyzing *comprehension* in LLMs, but constitutes an excellent basis for analyzing LLMs' discourse-level *generation* capabilities, i.e. going beyond the prediction of the next mention. While discourse-level downstream tasks in NLG, e.g. story generation or summarization, are complex and notoriously difficult to evaluate systematically with respect to targeted linguistic capacities of NLG systems, IC continuations provide a well-controlled diagnostic of discourse knowledge and, at the same time, rather simple sentences

whose quality can be easily assessed in human evaluation. Yet, to date, only little consideration has been given to the extent to which IC continuations generated by language models maintain semantic appropriateness and coherence with respect to the given IC prompts, cf. Huynh et al., 2022, and to what extent congruency with the human coreference bias is related to the quality and coherence of the generated continuation.

We expand prior work on discourse knowledge captured by LLMs in IC contexts and investigate their ability to generate not only bias-congruent but also sensible continuations. In contrast to most previous studies, we are not only interested in the LLMs predictions of the first word following the prompt (i.e. the pronoun), but in the quality of the generated sequences and their comparison against human continuations. We would like to emphasize that IC bias can be violated without any loss of discourse coherence. As an example, consider the following generated sequences in (2), where the first sentence is not congruent with the human bias but coherent, whereas the second sentence is congruent with the bias but not coherent:

(2)    a.    Jonathan shocked Charlotte because...
             *she didn't think he would steal.*
    b.    Jonathan admired Charlotte because...
             *she handed him a pineapple.*

In order to assess the quality of generated IC continuations, we carry out an evaluation study with human ratings of naturalness, coherence, and informativity. We aim to investigate to what extent continuation quality is associated with bias congruency and identify the factors that may influence this interaction. Thus, we manipulate two different types of conditions. First, we evaluate the predicted continuations of language models on two types of IC prompts: "standard" IC prompts (such as those in (1) and (2)) and IC prompts that are extended with adverbial modifiers. Second, we evaluate the performance of three different decoding procedures. In the following, Section 2 will present background on the set up of our study, Section 3 describes the hypotheses of our evaluation experiment and Section 4 describes the results.

## 2 Background

### 2.1 Implicit Causality

Psycholinguistic literature has consistently shown that numerous interpersonal verbs exhibit a bias known as "Implicit causality" (IC) (Garvey and Caramazza, 1974). That is, when asked to provide a continuation after "..." in sentences like (1) and (2), humans display strong next-mention preferences towards the subject for stimulus-experiencer (SE) verbs like "fascinate" (1) and "shock" (2-a) and towards the object for experiencer-stimulus (ES) verbs like "admire" (2-b) (Solstad and Bott, 2022). Continuations that align with the coreference bias are referred to as bias-congruent, while continuations that go against the bias are considered bias-incongruent. In addition, verbs that exhibit a pronounced IC bias also tend to exhibit a coherence bias, that is, they are prone to trigger explanations in subsequent discourse (cf., e.g., Kehler et al., 2008, Bott and Solstad, 2014). Interestingly, it has been shown that the coreference and coherence tendencies of IC verbs are strongly modulated by linguistic context. For instance, in their study, Bott and Solstad, 2021 showed that modifying standard IC prompts (such as those in (1) and (2)) with causal adverbial phrases leads to a change in the coreference bias as well as the coherence bias. In particular, not only did their results show that the IC coreference bias was in fact eliminated when adverbial modifiers were added to the prompts, they further found fewer and distinct types of explanations after these modifications compared to the default explanation types observed in the "standard" condition.

In psycholinguistic studies, the coreference and coherence biases of IC verbs are often elicited in sentence continuation tasks, typically giving the participants prompts such as those in the discussed examples. And, as several psycholinguistic studies have demonstrated that the IC bias is not only highly reliable but also robust across different languages (Ferstl et al., 2011; Goikoetxea et al., 2008; Hartshorne et al., 2013; Bott and Solstad, 2014), it has become an an intriguing domain for testing language models. Earlier studies, including those conducted by Upadhye et al., 2020, Davis and van Schijndel, 2020, Kementchedjhieva et al., 2021 and Zarrieß et al., 2022, have examined the performance of LLMs in capturing the IC coreference bias. I.e., they concentrated on single-word prediction tasks and evaluated the models' ability to generate continuations of such classic prompts, like examples (1) and (2), and predominantly found that LLMs display limited ability to systematically incorporate the IC coreference bias in their genera-

tions. In addition to examining the IC bias, Huynh et al., 2022 conducted a human evaluation of the quality of the continuations predicted by an English GPT-2 model. Asking the participants to judge the "reasonability" of the generated continuations on a 5-point Likert scale (with 5 being strongly reasonable), their results showed that only 32% of all the continuations retained an average rating greater than or equal to 4. In this study, we build upon these priors results, by investigating, similar to Huynh et al., 2022, the extent to which bias congruency in LLMs is associated with the quality of the predicted continuations, i.e. we go beyond the next word prediction. Additionally, we go beyond previous research by not only taking into account the impact of the decoding procedure, but also by investigating whether the grade of information density in the prompts affects the models' ability to produce meaningful continuations.

## 2.2 Information Density and Decoding

The use of natural language for communication is often explained through information theory (Shannon, 1948), an approach that views linguistic units as messages aimed at conveying information, quantified by their probability of being produced, which is also termed "surprisal". Within this view, units with low probability are considered more informative, reflecting the intuition that unpredictable elements convey more information than predictable ones (see, e.g., Meister et al., 2021 for an elaboration). The Uniform Information Density (UID) hypothesis (Levy and Florian Jaeger, 2007; Jaeger, 2010) further predicts that speakers prefer to distribute information uniformly across their utterances. I.e., if speakers wish to convey more information, they are more likely to distribute this information across more words. Previous studies have shown that a more uniform distribution of information is strongly associated with higher linguistic acceptability (cf., e.g., Meister et al., 2021). As such, assertions that adhere to the UID hypothesis are considered to be of higher quality and receive better evaluations. Furthermore, also the choice of decoding strategy, i.e. the decision rule used to determine the output sequence of a model, is one of the most important factors that affects the quality and various linguistic properties of the generated text, as several papers have demonstrated (e.g. Holtzman et al., 2019; Zarrieß et al., 2021; Meister et al., 2022; Ji et al., 2023). However,

the success of the decoding procedure shows to be contingent on the task at hand, and no decoding approach has demonstrated a consistent correlation with producing high-quality text (Wiher et al., 2022). Given these insights, in this study, we consider it worthwhile to explore how the level of information density of the prompts as well as the decoding procedure used influences the quality of the generated continuations. We will outline our approach and hypotheses in the next section.

## 3 Hypotheses and Conditions

### 3.1 Prompt Manipulation and Information Density

First, to investigate an important aspect of IC verbs that, to our knowledge, has not yet received attention in the study of IC in LLMs, we introduce an additional prompt condition alongside the "standard" IC prompts (such as those in the examples discussed so far). That is, we extend these standard prompts with adverbial modifiers, so that, for example, the prompts in (2-a) and (2-b) are augmented to (3-a) and (3-b), respectively.

(3)  a. Jonathan shocked Charlotte by his aggressive appearance in the talk show because...

b. Jonathan admired Charlotte for her extraordinary agility in this year's ice skating competition because...

This extension is motivated by the aforementioned insight that the coreference and coherence biases of IC verbs are strongly modulated by linguistic context (cf. Section 2.1). In this current study, we therefore compare standard IC prompts with prompts extended with adverbial modifiers, using the experimental data of Bott and Solstad, 2021, in order to investigate to what extent this added contextual information may not only influence the ability of LLMs to capture the IC bias but, beyond, also the quality of the model-generated continuations.

Building upon the UID, we make the assumption that standard IC prompts, which are brief and contain only minimal information, will require more information in their sentence continuations to maintain a uniform distribution of information across the whole sentence than prompts that are more detailed. For instance, consider Figure 1, which depicts the token probabilities and information values for a "standard" unmodified IC prompt and for a prompt
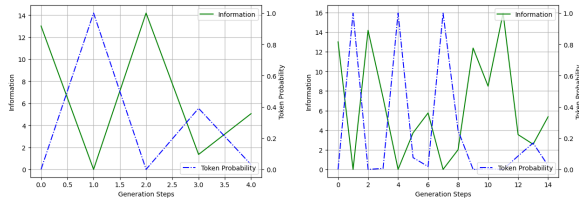
Figure 1: Token probabilities (blue) and information values (green) for each prompt condition for a selected item involving the ES verb *bewundern* ("admire"). Left-hand side: unmodified prompt, i.e. "Paul bewunderte Isabel, weil" (*Paul admired Isabelle, because*). Right-hand side: extended prompt, i.e. "Paul bewunderte Isabel für ihre außerordentliche Geschicklichkeit beim diesjährigen Eisschaulaufen, weil" (*Paul admired Isabel for her extraordinary agility in this year's ice skating competition, because*).

augmented with an adverbial modification. Here it is visible that the latter exhibits more peaks towards lower probability, and, respectively, towards higher informativity – since here, in addition to the verb, there is another information-carrying part (i.e. the adverbial modification). Put differently: The IC prompts that are extended with adverbial modifications inherently carry more information and therefore already contribute a greater amount of information to the (yet to be completed) sentence. In fact, as Bott and Solstad, 2021 show, prompts with adverbial modifiers such as (3) provide comprehensive causal scenarios in themselves lacking any need for further causal elaboration.

**Hypotheses** Thus, assuming that speakers aim to distribute information uniformly across an utterance, we can infer that continuations of the modified IC prompts would require less information than those of the shorter prompts. Consequently, when presented with standard (i.e. short) IC prompts, LLMs are expected to produce continuations that are less probable and, therefore, more informative (or: surprising), while with the modified IC prompts, less informative and therefore, more probable continuations are anticipated from the models. We hypothesize, however, that LLMs encounter difficulties in producing informative yet sensible continuations, which could explain the observed lower quality of sentence continuations for these prompts, aligning with research suggesting that a more uniform distribution of information is strongly associated with higher linguistic acceptability (see Section 2.2). Taken together, we expect that sentence continuations following prompts augmented with adverbials will be evaluated as of higher quality than those following standard IC prompts, as the ad-

ditional information provided reduces the model's burden to generate informative content on its own.

## 3.2 Information Density and Decoding

In this study, we acknowledge the significant impact of decoding methods on the quality of generated text (see Section 2.2) and therefore also investigate whether the distinct information-theoretic characteristics of three different decoding methods allow them to handle the prompt requirements differently. For instance, adhering to the terminology of UID theory, where information density is measured in information-theoretic terms of surprisal, maximation-based strategies, such as beam search decoding, for example, are known for producing text that is more probable but less surprising, and thus less informative. In contrast, stochastic strategies, such as Nucleus Sampling (Holtzman et al., 2019), for example, tend to produce text that is less probable and, therefore, more suprising, i.e. informative (cf., e.g., Zarrieß et al., 2021).

First of all, we consider seam search to be an appropriate candidate for the task of generating sentence continuations for (short) IC prompts, since Meister et al., 2020 have shown that beam search incorporates an inductive bias that aligns with the UID principle to distribute information uniformly across an utterance. However, as it has been shown that beam search tends to reduce diversity by favoring likelihood (Schüz et al., 2021; Zarrieß et al., 2021), in this paper, we will utilize the Diverse Beam Search method proposed by Vijayakumar et al., 2018, which, in essence, encourages diverse candidates by categorizing candidates into groups and then enforcing diversity within those groups. Furthermore, we consider Nucleus Sampling (Holtzman et al., 2019) to be an appropriate stochastic decoding method for our investigation. By truncating the model distribution, this method effectively addresses the drawback of sampling based methods to potentially select very low probability outputs that may considerably reduce the overall quality and coherence. In addition, in this study, we will consider another, third decoding procedure called Local Typical Sampling (Meister et al., 2022) to generate continuations for the two IC prompt conditions. Including this method in our study seems worthwhile, given that the authors introduced it as a potential solution to the shortcomings of beam search and Nucleus Sampling and, beyond, that it was designed to embody the

characteristics of human language, aligning with the information-theoretic perspective discussed in this context (cf. Meister et al., 2022). In a nutshell, in Local Typical Sampling, the options to sample from are limited to strings that have a similar information content to what would be expected based on the prior context.

**Hypotheses** In light of this, we anticipate that stochastic decoding methods could result in better sentence continuations for the standard IC prompts compared to beam search, for instance, as more information is required in the continuations that follow these short prompts. On the contrary, for prompts extended with adverbials, where less informative continuations are expected to retain a uniform distribution of information across the sentence, beam search decoding may lead to better results. Moreover, we anticipate that Local Typical Sampling will lead to adequate informative generations for both short and modified prompts, i.e. resulting in comparable quality of continuations across the two prompt conditions. Taken together, while we generally expect longer prompts (i.e. those augmented with adverbials) to result in better quality continuations, we further expect the decoding strategy employed to also play a significant role in the generated text's quality. Further, we anticipate that the impact of the decoding strategies will vary depending on the prompt construction.

## 4 Experimental Setup

**Data.** We ground our study on German IC data from Bott and Solstad, 2021 and use their experimental items to generate German prompts to be completed by the LLMs. These data also provide us with human-generated sentence continuations for both prompt conditions, offering a valuable reference point for evaluating the model-generated continuations. Prompts consist of simple sentences introducing the verb, the verb's arguments and the connective *weil* ('because') (as in (1) and (2)). To further investigate the effect of causal modification, we designed the following four conditions:

1. SE verbs in "standard" prompt constructions (e.g., *Clara inspired Vincent because...*)
2. SE verbs + *durch* ('by') modifier (e.g., *Clara inspired Vincent by her innovative lecture because...*)
3. ES verbs in "standard" prompt constructions (e.g., *Paul admired Isabel because...*)

4. ES verbs + *für* ('for') modifier (e.g., *Paul admired Isabel for her extraordinary agility in this year's ice skating competition because...*)

As in previous studies, we vary and balance prompts for the names and gender of verb arguments. To ensure that the prompts' informational content remains unaffected by the verb arguments, we deliberately keep the proper names consistent across the conditions. In addition, we also consider the sentiment of the verbs, ensuring a balance between those that are positively and negatively connotated. See Table 4 in the Appendix for the chosen verbs and examples for each prompt condition.

**Models.** We use the Hugging Face framework for reproducibility, employing two German language models to generate continuations for the IC prompts: (i) a pretrained German GPT-2 model[1] and (ii) a pretrained multilingual mGPT model[2] which reproduces the GPT-3 architecture. Following the reasoning of Huynh et al., 2022, we use GPT-2 and mGPT, despite the availability of newer and more advanced language models, as they offer a manageable size that is compatible with standard hardware and a favorable trade-off between complexity and efficiency. For both models we do not employ any finetuning.

**Decoding hyperparameters.** We set the hyperparameters of our decoding methods by validating them on the continuation data from Bott and Solstad, 2021, using other IC verbs (and thus, items) than those tested in the final experiments and the metrics BLEU, GLEU, METEOR, ROUGE(-L) and BERTScore, as provided by the Hugging Face library. We chose ranges of hyperparameters based on the authors' self-reported best-performing values and/or standard values recommended in literature and found the following settings to be best performing in at least four of the five metrics:

- Diverse Beam Search: beam size and beam group size = 10; diversity penalty $\lambda = 0.7$
- Nucleus Sampling: top p value = 0.85; temperature = 0.7
- Typical Sampling: typical p value = 0.9; temperature = 0.7

---

[1] https://huggingface.co/dbmdz/german-gpt2.
[2] https://huggingface.co/ai-forever/mGPT.

**Automatic evaluation.** To assess whether the IC coreference bias is reflected in the models' continuations, we adopt the method of Zarrieß et al., 2022 and calculate Completion Sensitivity scores, i.e. the percentage of prompts for which the model's predicted pronoun aligns with the bias. Further, to investigate the overall quality of the models' continuations, we calculate the three automatic measures: BLEU, ROUGE-L and BERTScore, in this way comparing them to the ones produced by participants in the studies conducted by Bott and Solstad, 2021.

**Human evaluation.** As human evaluation remains to be the gold standard when it comes to assessing the overall quality of a system (van der Lee et al., 2021; Schuff et al., 2023), we employ human judgment to investigate the quality of the generated continuations. Looking over the produced generations, there were three items that contained offensive content, e.g. including instances of sexism, which is why these items were excluded from the experiment. We used the Prolific framework to obtain ratings from 40 different annotators on 96 examples for each model's continuations, 128 examples for each prompt condition, 64 examples for each decoding method and 64 examples for human-produced continuations of which one half was bias congruent and the other half bias incongruent. In this way, we overall collected 5120 human judgments. We restricted potential evaluators to native speakers of German with their country of residency being Germany. We designed the evaluation as a rating task, presenting three questions to the evaluators. We asked whether the continuation is a (i) "natural", (ii) "meaningful" (i.e. coherent) and (iii) "surprising" (i.e. informative) explanation for the respective sentence beginning (i.e. prompt). The evaluators could indicate their agreement on a five-point Likert scale, ranging from 1 ('strongly disagree') to 5 ('strongly agree'). For each of the criteria, we take the median score across evaluators as the final score. We chose the criteria of naturalness and coherence following recommendations from van der Lee et al., 2021 and elicited informativity to investigate our assumptions described in Section 3. Before evaluators were asked to provide ratings on the three criteria, they were presented with instructions as well as an example item on the basis of which the criteria were explained. Evaluators were paid £9,00/hour and IRB approval was obtained prior to conducting the study.

|  | Naturalness | Coherence | Informativity |
|---|---|---|---|
| **Standard IC Prompt** | | | |
| Diverse Beam Search | 4 (**3.55**) | 3 (**2.87**) | 2 (2.50) |
| Nucleus Sampling | 4 (3.26) | 2 (2.55) | 3 (2.62) |
| Typical Sampling | 3 (3.26) | 3 (2.74) | 3 (**2.65**) |
| *Human bias-congruent* | 5 (4.77) | 5 (4.75) | 2 (2.39) |
| *Human bias-incongruent* | 4 (3.82) | 3 (3.20) | 3 (2.47) |
| **Modified IC Prompt** | | | |
| Diverse Beam Search | 4 (**3.69**) | 3 (**3.06**) | 2 (**2.54**) |
| Nucleus Sampling | 3 (2.90) | 2 (2.04) | 2 (2.40) |
| Typical Sampling | 3 (2.99) | 2 (2.24) | 2 (2.52) |
| *Human bias-congruent* | 5 (4.56) | 5 (4.61) | 3 (2.56) |
| *Human bias-incongruent* | 5 (4.54) | 5 (4.37) | 3 (2.61) |

Table 1: Human evaluation results for each decoding procedure and each IC prompt condition, aggregated over the text generations of both LLMs. For each criterion, we report the median score across raters as the final score (additionally, the mean values are reported in brackets). Bold values indicate conditions with the best values for that evaluation criteria.

## 5 Results

**Naturalness and coherence.** Generally, when comparing the two language models, GPT-2 and mGPT, the automatic metrics presented in Table 5 in the Appendix indicate that mGPT exhibits a tendency to generate more favorable, or rather, more similar continuations compared to the human-produced ones. However, this observation holds limited significance since, as detailed in a subsequent paragraph, there is a notably low correlation between automatic and human metrics. Thus, as we did not find substantial differences in continuation quality between GPT-2 and mGPT, we focus our analysis of human ratings on differences between the decoding methods and prompt conditions and aggregate over the models (but see Figure 5 in the Appendix for results separated for models). Table 1 shows the results of the human evaluation for each decoding procedure and each prompt condition. First, we note that, as expected, human-produced continuations achieve the best results for the criteria naturalness and coherence. Further, as we had hypothesized, it is the case that the influence of the decoding procedure varies depending on the IC prompt condition. Contrary to our expectations, however, it is not generally the case that the modified prompts lead to better-evaluated continuations. Likewise, it is not the case that the sampling-based methods result in better-evaluated continuations of the standard IC prompts while Diverse Beam Search leads to better-evaluated continuations of the modified prompts. Instead, we find that, across both prompt conditions, Diverse

Beam Search yields equally good and almost always best results – both for the automatic metrics and for the human evaluations. Remarkably, across prompt conditions and decoding procedures, we observe substantial distinctions between naturalness and coherence. I.e., the medians for naturalness consistently exhibit higher values, indicating that the generated continuations are perceived as fluent. However, in terms of coherence, the average ratings are lower, suggesting a lack of logical consistency in the generated text. This observation highlights the challenges faced by LLMs in generating sensible continuations in this rather simple task. It clearly shows that LLMs especially struggle with discourse-level coherence of explanations in IC contexts, beyond predicting the mentions that are congruent with this bias. This finding further underscores the importance of considering multiple evaluation criteria to assess the outputs of LLMs.

**Informativity.** Referring to Table 1, it is evident that, unlike the criteria of naturalness and coherence, the ratings for informativeness do not significantly favor human continuations over those generated by the models. This observation indicates that, as expected, the human-produced continuations align more closely with prototypical explanations making them appear comparatively mundane. Further, as consistent with our expectation, we observe that for the shorter prompts, the sampling-based methods produce continuations that the evaluators deem more informative compared to their continuations of the modified prompts. This observation also holds true for Typical Sampling, which we initially anticipated to yield similarly informative generations for both short and modified prompts. Interestingly, for Nucleus Sampling, the higher perceived informativeness of the continuations is accompanied by an increase in naturalness, whereas for Typical Sampling, it is linked to higher coherence of the continuations. Generally, this indicates that sampling-based methods are particularly responsive to the information density of the prompts.

**Relation of bias congruency and continuation quality.** Table 2 shows completion sensitivity results for each bias type, prompt condition and decoding procedure for continuations of the GPT-2 and mGPT models. Overall, the performance varies across different models, decoding procedures, and bias types. However, in general, the models are more likely to capture the object bias, as can be

|  | GPT-2 | | | mGPT | | |
|---|---|---|---|---|---|---|
|  | Diverse Beam Search | Nucleus Sampling | Typical Sampling | Diverse Beam Search | Nucleus Sampling | Typical Sampling |
| SE simple | 62.5 | 25 | 75 | 50 | 25 | 62.5 |
| SE modified | 75 | 50 | 75 | 87.5 | 50 | 87.5 |
| ES simple | 50 | 75 | 87.5 | 75 | 87.5 | 87.5 |
| ES modified | 50 | 100 | 87.5 | 75 | 100 | 87.5 |

Table 2: Completion Sensitivity (CS) scores for each model, bias type, prompt condition, and decoding procedure, aggregated over all types of individual verbs. CS scores are calculated as the percentage of continuations where the predicted pronoun is congruent with the IC bias.

noted by the (almost) overall higher CS values for the ES verbs, aligning with results from, for example, Kementchedjhieva et al., 2021 and Zarrieß et al., 2022 which as well point towards a general tendency of LLMs to establish coreference to the object. Moreover, it is noticeable that for each decoding procedure the ability to capture the IC bias of SE verbs tends to improve when prompts are augmented with adverbial modifiers. Thus, it appears that the augmentation of the standard IC prompts indeed has an impact on the extent to which the LLMs can capture the IC bias. Interestingly, this influence varies depending on the decoding strategy at hand. Further, it is noteworthy that, in this way, the language models exhibit a different behavior in capturing the IC bias when confronted with the modified prompts compared to the findings of Bott and Solstad, 2021's human-produced data, where the IC bias was in fact eliminated by the same modification (see Section 2.1). That is, while humans tend to produce fewer bias-congruent continuations when the IC prompts are augmented with such adverbial modifications, the opposite seems to be true for language models.

If we now consider the relation between bias congruency and continuation quality and examine the human-produced continuations in Table 1 first, we can find higher naturalness and coherence ratings for the bias-congruent continuations than for the bias-incongruent continuations, in particular for the standard IC prompts, which aligns with expectations based on Bott and Solstad, 2021. The observation that the bias-incongruent continuations of the modified prompts appear to be more acceptable than the bias-incongruent continuations of the standard prompts further aligns with the findings of Bott and Solstad, 2021, who demonstrated that modified prompts elicit different types of explanations, often referring to elements other than the

verbs' arguments.

Next, let us consider Figure 2, which depicts the relationship between bias congruency and the human evaluation criteria for each decoding method, aggregated across the two prompt conditions and across the generations of both LLMs (see Figure 12 in the Appendix for a visualiziation separated according to conditions). The higher green bars consistently observed across all prompt conditions and decoding methods indicate that each decoding method generated a greater number of bias-congruent continuations than bias-incongruent ones. Furthermore, for all three decoding methods, it is evident that non-bias congruent continuations receive lower ratings in terms of coherence, as indicated by the descending purple bars for this criterion. However, a preference for bias-congruent continuations being more natural than bias-incongruent continuations is primarily observed for the Diverse Beam Search decoding method. Further, it is interesting to note that although Typical Sampling tends to generate the most bias-congruent continuations, this does not necessarily result in better scores on the evaluation metrics. These results, thus, indicate that a bias congruent continuation does not equate to a qualitatively better continuation.



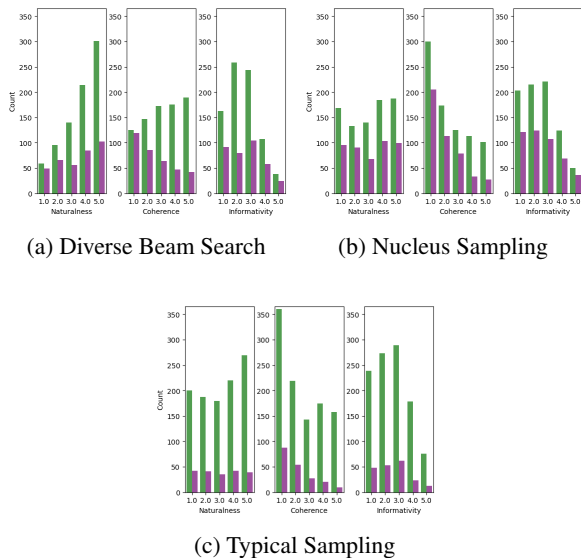(a) Diverse Beam Search     (b) Nucleus Sampling

(c) Typical Sampling

Figure 2: Relationship between bias congruency and the text continuation criteria for each decoding method, aggregated across SE verbs and ES verbs, across the two prompt conditions, and across the generations of both LLMs. Bias congruent continuations are depicted in green, bias incongruent continuations are depicted in purple.

**Correlation between automatic and human evaluation.** Table 1 reports the human ratings for the generated continuations, while Table 5 in the Ap-

| Human | BLEU | ROUGE-L | BERTScore |
|---|---|---|---|
| Naturalness | 0.16 ($p=0.03$) | -0.02 ($p=0.84$) | -0.04 ($p=0.59$) |
| Coherence | 0.18 ($p=0.01$) | 0.03 ($p=0.66$) | -0.01 ($p=0.91$) |
| Informativity | -0.18 ($p=0.02$) | -0.08 ($p=0.30$) | -0.07 ($p=0.35$) |

Table 3: Pearson's correlation coefficient between automatic and human evaluation metrics.

pendix displays the automatic metrics. At first sight, it may seem that the automatic metrics generally align with our human evaluations, as, e.g., Diverse Beam Search consistently outperforms other decoding methods in automatic and human scores (with one exception). However, Table 3 shows Pearson's correlation coefficients between the automatic and human evaluation metrics and it becomes apparent that there is no significant correlation between ROUGE-L and BERTScore on the one and human ratings on the other hand. BLEU scores do achieve a weakly significant correlation with coherence ratings, but not with naturalness or informativity. Interestingly, there even seems to be a negative (but hardly significant) relationship between BLEU and informativity. Notably, these automatic metrics seem to fail even more miserably on our linguistically controlled task, as compared to correlations reported for downstream task evaluations as in, e.g., Savkov et al., 2022. We believe that this may be due to the fact that the scoring of differences between generated continuations in this rather restricted task may require a greater awareness of linguistic subtleties and deeper discourse understanding than what is currently captured by these metrics. Overall, these findings underscore the significant challenge faced by NLG metrics in accurately capturing the nuanced aspects of human evaluation and emphasize the need for cautious interpretation of automatic evaluation scores.

## 6   Conclusion

This paper investigated how continuations of IC prompts can be used to evaluate the text generation capabilities of language models, expanding prior work on discourse knowledge captured by LLMs in IC contexts by investigating their ability to generate not only bias-congruent but also sensible continuations. Our study reveals that LLMs face challenges in generating coherent continuations for relatively simple prompts, highlighting their struggle with discourse-level coherence. Further, our findings show that both the decoding method and the information density of the prompt have a substan-

tial impact on the quality of generated text, even surpassing the influence of the specific language model used. Our results further indicate that modifying the standard IC prompts has a notable effect on the LLMs' capability to capture the IC bias, depending on the decoding strategy employed. At the same time, the results indicate that a bias congruent continuation does not necessarily equate to a higher quality of the continuation. Intriguingly, we observe a surprisingly low correlation between automatic evaluation metrics and human judgments. This poses an interesting challenge for natural language generation, calling for future research to enhance the evaluation methodologies and metrics used in NLG systems. Another potential future direction for our study is to expand the scope beyond German data, as this is a major limitation of this work. While IC is known to be cross-linguistically stable, the inclusion of other languages in our investigation should be performed to validate our findings. Also, it would be interesting to more extensively investigate the models' strategy for choosing the first word of the continuation (i.e. the pronoun), which may simply consist in selecting the most recently mentioned discourse element. Further, it could be valuable to investigate the effects of additional prompt modifications on bias congruency and the quality of continuations. For this, one possible foundation could be the research conducted by Koornneef and Van Berkum (2006), for instance, where IC prompts are integrated within a larger pre-context, making them arguably a more natural option for evaluating LLMs compared to the prompts investigated in this study. Also, Hoek et al., 2021, for example, investigated IC in the context of relative clauses, which could provide another compelling starting point for further examination of LLMs in this context.

**Supplementary Materials Availability Statement:** Source code, prompts used for generating the models' continuations and the data from the human evaluation study are available from Github.[3]

---

[3]https://github.com/clause-bielefeld/implicit-causality-beyond-the-bias.

# References

Yonatan Belinkov and James Glass. 2019. Analysis methods in neural language processing: A survey. *Trans. Assoc. Comput. Linguist.*, 7:49–72.

Oliver Bott and Torgrim Solstad. 2014. From verbs to discourse: A novel account of implicit causality. In Barbara Hemforth, Barbara Mertins, and Cathrine Fabricius-Hansen, editors, *Psycholinguistic Approaches to Meaning and Understanding across Languages*, pages 213–251. Springer International Publishing, Cham.

Oliver Bott and Torgrim Solstad. 2021. Discourse expectations: explaining the implicit causality biases of verbs. *Linguist. Philos.*, 59(2):361–416.

Zhenguang G Cai, David A Haslett, Xufeng Duan, Shuqi Wang, and Martin J Pickering. 2023. Does ChatGPT resemble humans in language use?

Forrest Davis and Marten van Schijndel. 2020. Discourse structure interacts with reference but not syntax in neural language models. *Proceedings of the 24th Conference on Computational Natural Language Learning*, pages 396–407.

Allyson Ettinger. 2020. What BERT is not: Lessons from a new suite of psycholinguistic diagnostics for language models. *Trans. Assoc. Comput. Linguist.*, 8:34–48.

Evelyn C Ferstl, Alan Garnham, and Christina Manouilidou. 2011. Implicit causality bias in english: a corpus of 300 verbs. *Behav. Res. Methods*, 43(1):124–135.

Catherine Garvey and Alfonso Caramazza. 1974. Implicit causality in verbs. *Linguist. Inq.*, 5(3):459–464.

E. Goikoetxea, G. Pascual, and J. Acha. 2008. Normative study of the implicit causality of 100 interpersonal verbs in Spanish. *Behavior Research Methods*, 40:760–772.

Joshua K Hartshorne, Yasutada Sudo, and Miki Uruwashi. 2013. Are implicit causality pronoun resolution biases consistent across languages and cultures? *Exp. Psychol.*, 60(3):179–196.

Jet Hoek, Hannah Rohde, Jacqueline Evers-Vermeul, and Ted J M Sanders. 2021. Expectations from relative clauses: Real-time coherence updates in discourse processing. *Cognition*, 210:104581.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration.

Jennifer Hu, Jon Gauthier, Peng Qian, Ethan Wilcox, and Roger Levy. 2020. A systematic assessment of syntactic generalization in neural language models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1725–1744, Online. Association for Computational Linguistics.

Hien Huynh, Tomas O Lentz, and Emiel van Miltenburg. 2022. Implicit causality in GPT-2: a case study.

T Florian Jaeger. 2010. Redundancy and reduction: speakers manage syntactic information density. *Cogn. Psychol.*, 61(1):23–62.

Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *ACM Comput. Surv.*, 55(12):1–38.

Andrew Kehler, Laura Kertz, Hannah Rohde, and Jeffrey L Elman. 2008. Coherence and coreference revisited. *J. Semant.*, 25(1):1–44.

Yova Kementchedjhieva, Mark Anderson, and Anders Søgaard. 2021. John praised Mary because _he_? implicit causality bias and its interaction with explicit cues in LMs. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4859–4871, Online. Association for Computational Linguistics.

Arnout W Koornneef and Jos J A Van Berkum. 2006. On the use of verb-based implicit causality in sentence comprehension: Evidence from self-paced reading and eye tracking. *J. Mem. Lang.*, 54(4):445–465.

Roger Levy and T Florian Jaeger. 2007. Speakers optimize information density through syntactic reduction. *Adv. Neural Inf. Process. Syst.*, 19.

Clara Meister, Ryan Cotterell, and Tim Vieira. 2020. If beam search is the answer, what was the question? In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2173–2185, Online. Association for Computational Linguistics.

Clara Meister, Tiago Pimentel, Patrick Haller, Lena Jäger, Ryan Cotterell, and Roger Levy. 2021. Revisiting the Uniform Information Density hypothesis. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 963–980, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Clara Meister, Tiago Pimentel, Gian Wiher, and Ryan Cotterell. 2022. Locally typical sampling.

Aleksandar Savkov, Francesco Moramarco, Alex Papadopoulos Korfiatis, Mark Perera, Anya Belz, and Ehud Reiter. 2022. Consultation checklists: Standardising the human evaluation of medical note generation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 111–120, Abu Dhabi, UAE. Association for Computational Linguistics.

Hendrik Schuff, Lindsey Vanderlyn, Heike Adel, and Ngoc Thang Vu. 2023. How to do human evaluation: A brief introduction to user studies in NLP. *Natural Language Engineering*, pages 1–24.

Sebastian Schuster and Tal Linzen. 2022. When a sentence does not introduce a discourse entity, transformer-based models still sometimes refer to it.

Simeon Schüz, Ting Han, and Sina Zarrieß. 2021. Diversity as a by-product: Goal-oriented language generation leads to linguistic variation. In *Proceedings of the 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 411–422, Singapore and Online. Association for Computational Linguistics.

C E Shannon. 1948. A mathematical theory of communication. *The Bell System Technical Journal*, 27(3):379–423.

Torgrim Solstad and Oliver Bott. 2022. On the nature of implicit causality and consequentiality: the case of psychological verbs. *Language, Cognition and Neuroscience*, 37(10):1311–1340.

Shiva Upadhye, Leon Bergen, and Andrew Kehler. 2020. Predicting reference: What do language models learn about discourse models? *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 977–982.

Chris van der Lee, Albert Gatt, Emiel van Miltenburg, and Emiel Krahmer. 2021. Human evaluation of automatically generated text: Current trends and best practice guidelines. *Comput. Speech Lang.*, 67:101151.

Ashwin Vijayakumar, Michael Cogswell, Ramprasaath Selvaraju, Qing Sun, Stefan Lee, David Crandall, and Dhruv Batra. 2018. Diverse beam search for improved description of complex scenes. *AAAI*, 32(1).

Gian Wiher, Clara Meister, and Ryan Cotterell. 2022. On decoding strategies for neural text generators. *Trans. Assoc. Comput. Linguist.*, 10:997–1012.

Sina Zarrieß, Hannes Groener, Torgrim Solstad, and Oliver Bott. 2022. This isn't the bias you're looking for: Implicit causality, names and gender in German language models. In *Proceedings of the 18th Conference on Natural Language Processing (KONVENS 2022)*, pages 129–134, Potsdam, Germany. KONVENS 2022 Organizers.

Sina Zarrieß, Henrik Voigt, and Simeon Schüz. 2021. Decoding methods in neural language generation: A survey. *Information*, 12(9):355.

# A  Appendix

| Verb | Prompt Condition | Sentiment | Verb Type | Gender Order | Prompt |
|------|-----------------|-----------|-----------|--------------|--------|
| bewundern | simple prompt | positive | ES | f-m | Isabel bewunderte Paul, weil |
| bewundern | modified prompt | positive | ES | m-f | Paul bewunderte Isabel für ihre außerordentliche Geschicklichkeit beim diesjährigen Eisschaulaufen, weil |
| enttäuschen | names_simple | negative | SE | m-f | Björn enttäuschte Celina, weil |
| enttäuschen | names_pp | negative | SE | f-m | Celina enttäuschte Björn durch ihr unhöfliches Benehmen beim Geschäftsessen, weil |
| faszinieren | names_simple | positive | SE | f-m | Viktoria faszinierte Steven, weil |
| faszinieren | names_pp | positive | SE | m-f | Steven faszinierte Viktoria durch seine eindrücklichen Reiseberichte, weil |
| hassen | names_simple | negative | ES | m-f | Malte hasste Pia, weil |
| hassen | names_pp | negative | ES | f-m | Pia hasste Malte für die täglichen Sticheleien auf dem Schulhof, weil |
| inspirieren | names_simple | positive | SE | f-m | Clara inspirierte Vincent, weil |
| inspirieren | names_pp | positive | SE | m-f | Vincent inspirierte Clara durch seine innovative Vorlesung, weil |
| respektieren | names_simple | positive | ES | m-f | Martin respektierte Lina, weil |
| respektieren | names_pp | positive | ES | f-m | Lina respektierte Martin für seine couragierte Hilfe beim Löschen des Feuers, weil |
| schockieren | names_simple | negative | SE | f-m | Charlotte schockierte Jonathan, weil |
| schockieren | names_pp | negative | SE | m-f | Jonathan schockierte Charlotte durch sein aggressives Auftreten in der Talkshow, weil |
| verabscheuen | names_simple | negative | ES | m-f | Stefan verabscheute Miriam, weil |
| verabscheuen | names_pp | negative | ES | f-m | Miriam verabscheute Stefan für seine unnötige Hektik bei der Abfertigung von Patienten, weil |

Table 4: Verbs used in the study and examples for each prompt condition.

# Instruktionen

In dieser Studie besteht Ihre Aufgabe darin, **Satzfortsetzungen, die von Computermodellen generiert wurden, zu beurteilen.**

Alle folgenden Aufgaben haben die gleiche Form. Sie sehen in der ersten Zeile den Anfang eines Satzes und in der zweiten Zeile eine von einem Computermodell generierte **Fortsetzung, die eine Begründung für die Aussage im Satzanfang liefern soll**. Ihre Aufgabe ist es, anhand von drei verschiedenen Aussagen zu beurteilen, inwiefern es Computermodellen gelingt, gute mögliche Begründungen für die Aussagen in den Satzanfängen zu generieren.

Wenn Sie zum Beispiel den folgenden **Satzanfang** sehen:

> Satzanfang:
>
> **Hannah amüsierte Anton, weil**

wäre eine mögliche **Begründung**:

> Begründung:
>
> **sie die besten Witze erzählte.**

In jedem Durchgang sehen Sie dann drei Aussagen über die gerade angezeigte Begründung.
Bei diesen drei Aussagen geben Sie bitte jeweils an, ob bzw. wie sehr sie dieser zustimmen oder nicht zustimmen. Dafür wählen Sie bitte jeweils ein Feld in der 5-Punkte-Skala aus, aufsteigend von **1 (stimme der Aussage gar nicht zu)** bis **5 (stimme der Aussage voll zu)**:

Die Begründung wirkt **natürlich** und liest sich so, als ob sie von einer/m deutschen Muttersprachler/in geschrieben wurde.

| | stimme gar nicht zu | 1 | 2 | 3 | 4 | 5 | stimme voll zu |

Die Begründung ist **sinnvoll**, es gibt einen logischen Zusammenhang zwischen dem Satzanfang und der Fortsetzung.

| | stimme gar nicht zu | 1 | 2 | 3 | 4 | 5 | stimme voll zu |

Die Begründung ist **überraschend**, dadurch könnte der Satz insgesamt ein interessanter Beginn einer Geschichte sein.

| | stimme gar nicht zu | 1 | 2 | 3 | 4 | 5 | stimme voll zu |

In diesem Fall wäre es wohl natürlich, auf den ersten beiden Skalen eine Position weiter rechts auszuwählen, denn die Fortsetzung liest sich flüssig und bietet eine sinnvolle Erklärung für die Aussage im Satzanfang. Auf der dritten Skala würde man vielmehr eine Position weiter links auswählen, da die generierte Erklärung für diesen Satzanfang eher erwartbar und damit weniger überraschend ist.

Auf der nächsten Seite bekommen Sie Gelegenheit, sich mit der Aufgabe vertraut zu machen, bevor das eigentliche Experiment anfängt.

**Wichtig:** Jeder Durchgang im Experiment steht für sich allein und Ihre Urteile sollten sich deshalb stets nur auf den gerade vorliegenden Satz und die gerade vorliegende Begründung beziehen.

Nachdem Sie ein Urteil gefällt haben, bestätigen Sie die Eingabe bitte mit dem "Weiter"-Knopf. Bitte benutzen Sie im Experiment **NIE** die "Zurück"-Taste im Browser, da dies zum sofortigen Abbruch des Versuchs führen kann.

Sind Sie bereit für zwei Übungsdurchgänge?

Weiter

Figure 3: Instructions used in the online experiment. Participants were given an example items as well as an explanation of what would have been a reasonable rating on each of the three evaluation criteria (naturalness, coherence and informativity).

218

# Übung

Satzanfang:

**Nikolas entzückte Maria, weil**

Begründung:

**er ihr ein Geschenk mitgebracht hatte.**

---

*Es folgen nun drei verschiedene Aussagen über die gegebene Begründung.*

*Bitte lesen Sie jede dieser Aussagen und geben Sie an, inwiefern sie dieser zustimmen oder nicht zustimmen. Geben Sie Ihre Antwort bitte jeweils auf der Skala von 1 (stimme der Aussage gar nicht zu) bis 5 (stimme der Aussage voll zu) an. Bitte berücksichtigen Sie bei Ihrem Urteil immer die Fortsetzung in Hinblick auf den gegebenen Satzanfang.*

Die Begründung wirkt **natürlich** und liest sich so, als ob sie von einer/m deutschen Muttersprachler/in geschrieben wurde.

|  | ○ | ○ | ○ | ○ | ○ |  |
|---|---|---|---|---|---|---|
| stimme gar nicht zu | 1 | 2 | 3 | 4 | 5 | stimme voll zu |

Die Begründung ist **sinnvoll**, es gibt einen logischen Zusammenhang zwischen dem Satzanfang und der Fortsetzung.

|  | ○ | ○ | ○ | ○ | ○ |  |
|---|---|---|---|---|---|---|
| stimme gar nicht zu | 1 | 2 | 3 | 4 | 5 | stimme voll zu |

Die Begründung ist **überraschend**, dadurch könnte der Satz insgesamt ein interessanter Beginn einer Geschichte sein.

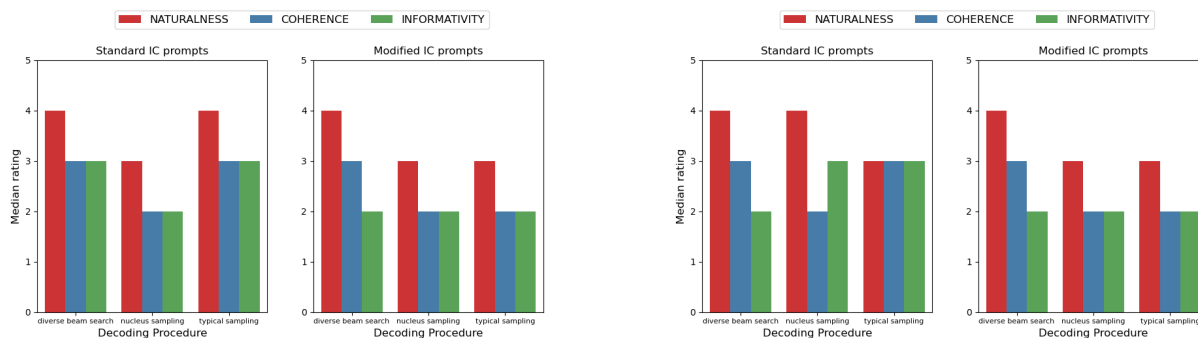|  | ○ | ○ | ○ | ○ | ○ |  |
|---|---|---|---|---|---|---|
| stimme gar nicht zu | 1 | 2 | 3 | 4 | 5 | stimme voll zu |

Weiter

Progress: ◀━━━━━━━━━━

Figure 4: One of two training items used in the online experiment for the participants to get familiar with the task and the rating scales.

| | GPT2 | | | mGPT | | |
|---|---|---|---|---|---|---|
| | BLEU | ROUGE-L | BERTScore | BLEU | ROUGE-L | BERTScore |
| **Standard IC Prompt** | | | | | | |
| Diverse Beam Search | 0.47 | 0.074 | **0.592** | **0.565** | **0.089** | 0.544 |
| Nucleus Sampling | 0.485 | 0.062 | 0.539 | 0.428 | 0.069 | 0.551 |
| Typical Sampling | 0.346 | 0.063 | 0.569 | 0.383 | 0.072 | 0.577 |
| **Modified IC Prompt** | | | | | | |
| Diverse Beam Search | 0.432 | 0.056 | 0.588 | **0.472** | **0.076** | 0.587 |
| Nucleus Sampling | 0.377 | 0.061 | 0.578 | 0.429 | 0.064 | **0.59** |
| Typical Sampling | 0.404 | 0.059 | 0.612 | 0.37 | 0.056 | 0.58 |

Table 5: Automatic evaluation results for each decoding procedure and each IC prompt condition. Bold values indicate conditions with the best values for that metric.

(a) Median ratings for continuations of GPT-2 model.

(b) Median ratings for continuations of mGPT model.

Figure 5: For each model, median ratings for each decoding procedure, each prompt condition and each text evaluation criteria of the human evaluation.
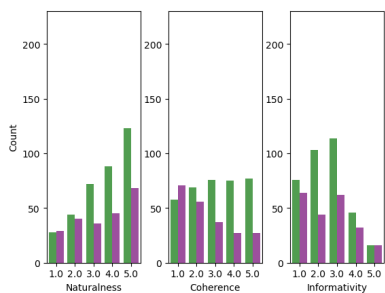


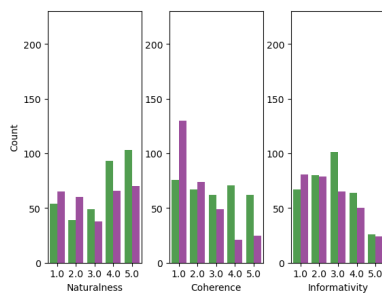Figure 6: Diverse Beam Search & Simple IC prompts

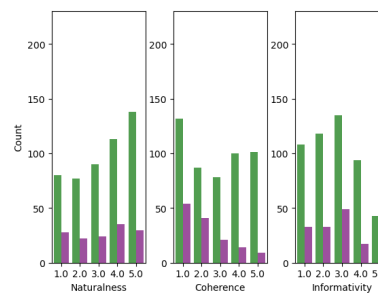Figure 7: Nucleus Sampling & Simple IC prompts
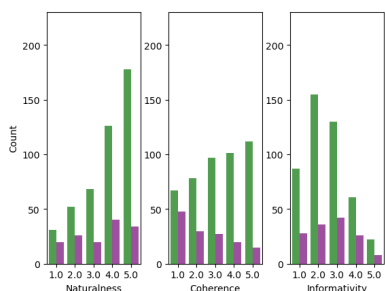
Figure 8: Typical Sampling & Simple IC prompts

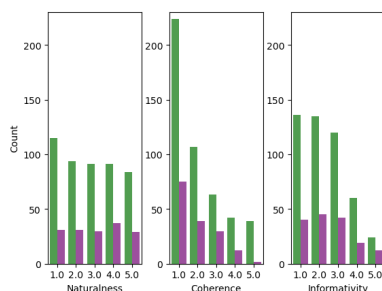Figure 9: Diverse Beam Search & Modified IC prompts

Figure 10: Nucleus Sampling & Modified IC prompts
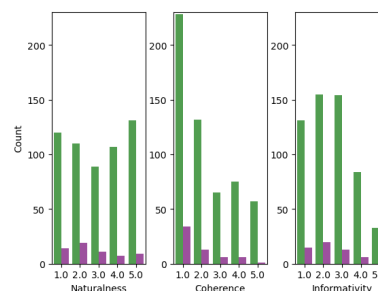
Figure 11: Typical Sampling & Modified IC prompts

Figure 12: Relationship between bias congruency and the text continuation criteria for each decoding method and for each prompt condition, aggregated across SE verbs and ES verbs and across the generations of both LLMs. Bias congruent continuations are depicted in green, bias incongruent continuations are depicted in purple.

# Enhancing factualness and controllability of Data-to-Text Generation via data Views and constraints

**Craig Thomson**[1*], **Clément Rebuffel**[2*],
**Ehud Reiter**[1], **Laure Soulier**[2], **Somayajulu Sripada**[1], **Patrick Gallinari**[2,3]
[1]University of Aberdeen, UK
[2]LIP6, Sorbonne Université, Paris, France
[3]Criteo AI Lab, Paris, France

## Abstract

Neural data-to-text systems lack the control and factual accuracy required to generate useful and insightful summaries of multidimensional data. We propose a solution in the form of *data views*, where each view describes an entity and its attributes along specific dimensions. A sequence of views can then be used as a high-level schema for document planning, with the neural model handling the complexities of micro-planning and surface realization. We show that our view-based system retains factual accuracy while offering high-level control of output that can be tailored based on user preference or other norms within the domain.

## 1 Introduction

The original vision of data-to-text generation was to take complex data and describe it using natural language such that humans could better understand it (Reiter et al., 2005; Reiter, 2007). Neural data-to-text systems commonly transcribe an isolated, simple data structure to a natural language. However, in many domains, e.g., finance, cinema box-office, or weather, the structured-data to be described are not independent, but rather exist as points along multi-dimensional axes such as time or entities (like people, companies, or locations). Figure 1 shows an example of a human-authored basketball game summary that requires data from outwith an individual game record. Such summaries often mention the upcoming games for each team (e.g. last two sentences) and aggregate player statistics over prior games.

Commercial applications deal with this complex scenario using hand-crafted rules (Reiter and Dale, 1997; Teixeira et al., 2020; Dale, 2020), although recent research leveraging Deep Learning techniques has looked to alleviate this bur-

den. However, such systems remain largely end-to-end, offering little user specified control[1]. On small sets of triples or attributes (Gardent et al., 2017; Dušek et al., 2018), neural systems can produce fluent, in domain generations, although they can struggle with factual accuracy mistakes such as hallucination. Additional challenges are introduced when using complex data such as tables (Parikh et al., 2020), compounded when considering interlinked structured data and longer texts (Wiseman et al., 2017).

While challenging, this complex setting is also an opportunity. In this paper, we propose the novel concept of data views, where all available data is divided into manageable subsets, describing entities and their attributes along one or more dimensions. Views are then aligned with corresponding spans of text. They can then be combined to form high level document plans (*schema*) for neural data-to-text systems, controlling the generation of text without sacrificing factual accuracy. We investigate the use of views in the domain of automated journalism (English language basketball game summaries), using the SportSett dataset (Thomson et al., 2020a) which extends the RotoWire dataset (Wiseman et al., 2017).

## 2 Related Work

Factualness and controllability are critical issues for data-to-text systems. Studies have shown that for end users of systems, accuracy is more of a concern than readability (Law et al., 2005) and users prefer texts tailored to their needs (van der Lee et al., 2017; Gatt and Krahmer, 2018).

Recently, the research community has focused on neural approaches, aiming to solve data-to-text tasks by leveraging advances of deep learning in language modeling fields (in particular neural machine translation) (Bahdanau et al., 2015; Vaswani

---

*These authors contributed equally to this work. Corresponding author: Craig Thomson (c.thomson@abdn.ac.uk, c.thomson.nlp@gmail.com)

[1]The INLG2021 panel 'What users want from real-world NLG' highlighted the need for system control.

**<Whole-Game=T_{G1}+T_{G2}>** The Oklahoma City Thunder defeated the host Miami Heat, 118-102, at American Airlines Arena on Friday. **<Within-Game=T_{E1}+T_{E2}>** While this wasn't a wire-to-wire win for Oklahoma City, they won this game in dominating fashion. **<Within-Game=T_{E1}+T_{E2}>** In fact, a 31-24 first quarter really set the tone, with a 41-29 second quarter sealing the victory. **<Within-Game=T_{E1}+T_{E2}>** The Thunder actually led by at least 15 points for the entirety of the second half. **<Whole-Game=T_{G1}+T_{G2}>** Three-point shooting was the key difference, with Oklahoma City hitting 16-of-30 and Miami connecting on 11-of-36. **<Whole-Game=T_{G1}+T_{G2}>** The Thunder also dominated in transition, winning the fastbreak differential, 23-8. **<Whole-Game=P_{G3}>** The Thunder (33-18) were led by Paul George, as he tallied 43 points, seven rebounds, five assists and two steals. **<Whole-Game=P_{G4}>** Russell Westbrook collected 14 points, 12 rebounds and 14 assists. **<Whole-Game=P_{G5}>** Steven Adams accrued 13 points, seven rebounds, two assists and three steals. **<Whole-Game=P_{G6}>** Dennis Schroder was huge off the bench, as he provided 28 points on 11-of-13 from the field. **<Whole-Game=P_{G16}>** The Heat (24-26) were led by Kelly Olynyk, as he provided 21 points, seven rebounds and two assists off the bench. **<Whole-Game=P_{G17}>** Josh Richardson led the starters with 18 points, four rebounds, three assists and two steals. **<Whole-Game=P_{G18}>** Hassan Whiteside amassed 12 points and 16 rebounds. **<Whole-Game=P_{G19}>** Justise Winslow finished with 10 points, two rebounds and five assists. **<Between-Game=T_{T2}>** Oklahoma City returns to action on Sunday, as they travel to face the Boston Celtics. **<Between-Game=T_{T1}>** As for Miami, they play host to the struggling Indiana Pacers on Saturday for their next outing.

Figure 1: Human authored summary for OKC@MIA on February 1$^{st}$ 2019. Tags and colours such as **<Whole-Game=T_{G1}+T_{G2}>** map to views (or unions of), some examples of which are shown in Figure 2

| ID | Team Name | PTS | REB | Wins | Losses | ... |
|---|---|---|---|---|---|---|
| T_{G1} | Miami Heat | 102 | 47 | 24 | 26 | ... |
| T_{G2} | Oklahoma City Thunder | 118 | 50 | 33 | 18 | ... |

Team Whole-GameViews (partial)

| ID | Name | PTS | REB | AST | STL | BLK | ... |
|---|---|---|---|---|---|---|---|
| P_{G3} | Paul George | 43 | 7 | 5 | 2 | 0 | ... |
| P_{G4} | Russel Westbrook | 14 | 12 | 14 | 1 | 1 | ... |

Player Whole-GameViews (partial)

| ID | Team Name | H1_PTS | Q1_PTS | Q2_PTS | ... |
|---|---|---|---|---|---|
| T_{E1} | Miami Heat | 53 | 24 | 29 | ... |
| T_{E2} | Oklahoma City Thunder | 72 | 31 | 41 | ... |

Team Within-GameViews (partial)

| ID | Name | Q1_PTS | Q1_REB | Q2_PTS | Q2_REB | ... |
|---|---|---|---|---|---|---|
| P_{E3} | Paul George | 16 | 1 | 10 | 1 | ... |
| P_{E4} | Russel Westbrook | 5 | 5 | 0 | 4 | ... |

Player Within-GameViews (partial)

| ID | Team Name | Opp_Place | Opp_Name | Location | ... |
|---|---|---|---|---|---|
| T_{T1} | Miami Heat | Indiana | Pacers | Miami | ... |
| T_{T2} | Oklahoma City Thunder | Boston | Celtics | Boston | ... |

Team Between-GameViews (partial)

| ID | Name | PTS_2 | ... | PTS_7 | ... | REB_2 | ... | REB_7 |
|---|---|---|---|---|---|---|---|---|
| P_{T3} | Paul George | 80 | ... | 237 | ... | 0 | ... | 2 |
| P_{T4} | Russel Westbrook | 37 | ... | 140 | ... | 2 | ... | 7 |

Player Between-GameViews (partial).

Figure 2: Example (partial) data views. PTS=point, REB=rebound, AST=assist, STL=steal, BLK=block. Q1, Q2, H1 etc., refer to **Q**uarters and **H**alves. PTS_X indicates the SUM of PTS over X games. Each row is considered to be an individual record within the given type of view

et al., 2017). Neural systems can blur the distinction between each sub-task of the pipeline approach, and are able to learn end-to-end to generate in-domain text from structured data (Lebret et al., 2016; Wiseman et al., 2017; Wang, 2019; Puduppully et al., 2019b).

Our work draws on a wide body of prior research on controllability and factualness, namely data engineering, controllable text generation, as well as planning.

**Data engineering** It is increasingly clear that careful design of datasets used to train deep neural models matters significantly (Rogers, 2021). On simpler data-to-text tasks such as the E2E challenge (Dušek et al., 2018), a number of data-level techniques have been proposed to improve factual accuracy and controllability, including dataset curation (Nie et al., 2019), data-to-text alignment (Dušek et al., 2019; Perez-Beltrachini and Lapata, 2018), straightforward control via input

manipulation (Filippova, 2020), and fine-grained annotation (Castro Ferreira et al., 2018; Rebuffel et al., 2021). However, these techniques are not suited for the complexity of the problem at hand, and do not enable *interactability*.

On complex datasets, adding extra records to increase the coverage of the source data over the description (Thomson et al., 2020b) has shown encouraging results, but not all cases of hallucination are clear-cut and easy to solve with a few records. Gong et al. (2019) suggested that it might be possible to model tables as three-dimensional, with rows, column, and time as the dimensions. This is not satisfactory because (1) current neural models cannot correctly perform the arithmetics required to generate the types of utterances found in the reference texts (Nie et al., 2018); (2) adding dimensions leads to intractable complexity.

**Controllable Text Generation** Controllable Text Generation (CTG) techniques traditionally

involve conditioning an NLG system on several control factors of style (e.g. tone, tense, length, etc.) (Dong et al., 2017; Hu et al., 2017; Ficler and Goldberg, 2017), or content (e.g. customized summaries based on aspect queries (Amplayo et al., 2021)). Control factors are often framed as a collection of key-value pairs, similarly to a typical data-to-text setting.

In the biography domain (Lebret et al., 2016), Filippova (2020) explicitly introduced CTG to data-to-text via an *hallucination score* simply attached as an additional attribute which reflects the amount of divergence in the target reference. Prompting (i.e. starting generation from textual instructions (Liu et al., 2021)) can also provide some control, with (Li and Liang, 2021) obtaining encouraging results at managing the length of generated descriptions. Contrasting with document-level approaches, Rebuffel et al. (2021) propose a finer-grained controllability, via word-level attributes to learn the relevant parts of each training instance.

**Planning and Schema** Macro-Planning, i.e. high-level planning of ideas, has long been used in traditional NLG pipelines, and has recently been introduced to neural systems as well. On small-scale datasets with short inputs/outputs (e.g. the WebNLG corpus), these approaches rely on detailed annotations of sentence structure and mention placement (Castro Ferreira et al., 2018), or the strong assumption that descriptions describe the associated data entirely (and nothing else) (Xu et al., 2021). However, these are unreasonable dependencies for large-scale datasets, with prohibitive size and complexity of inputs and outputs.

On more complex tasks, the two-step neural approach of (Puduppully et al., 2019a; Puduppully and Lapata, 2021) has proven effective at reducing factual mistakes and provides a small degree of controllability. Given all possible combinations of entities, a planner first selects which will be part of the narrative. In a second step, a generative module learns to output descriptions based on the selected entities. Designed this way, the planning step scales poorly in the input size, since it needs to consider all possible combinations. Furthermore, while the plans can be edited, the impact of individual edits on the final output is unclear (since the model is an end-to-end encoder-decoder). Additionally, no restrictions are placed during decoding to ensure that (1) the decoding process follows the order of the plan; (2) the decoder's copy mechanism doesn't copy attributes from entities in other part of the plans. Lastly, no restrictions are placed on the length of the texts corresponding to each part of the plan, the decoder has to decide the number of sentences to attribute to each item. Generated texts are of similar size, independent of the chosen plan.

## 3 Data views and their design

In multidimensional settings, descriptions of a data structure mostly focus on an initial point based upon the narrative intent, but also often compare subsets of data along different points of an axis. For example, the best player in the game being summarised might be described in an initial sentence:

> "Steph Curry led the Warriors with 43 points and 12 assists."

before an elaboration for the same player, but describing their performance over multiple games:

> "It was his fourth consecutive double-double[2]."

We make explicit this latent partitioning of the data, via views, as a solution to both the handling of dimensional data, and the alignment of data to text. We split the associated data following the same partitioning, and a view is defined as the records (i.e., key-value pairs) for one entity, from within one partition. In cases where several views are aligned with the same span of text, e.g. a sentence comparing two entities, views can be combined to form view sets.

Figures 1 and 2 illustrate this mirrored multi-dimensionality, along the time axis. The opening sentence describes the basketball game in focus, with details of the teams and their respective scores. This team data can be seen in the left (team) <Whole-Game> views. In contrast, the final two sentences report the upcoming opponents for both teams. This data, which lies elsewhere on the axis, is shown in the left <Between-Game> views. The third sentence mentions team statistics within quarters of the game, this data is shown in the left <Within-Game> views.

---

[2] https://en.wikipedia.org/wiki/Double-double

## 3.1 Designing views

We argue that despite its subtleties, data can be broadly categorized along distinct axes. These axes are data and task dependent and must be decided by experts on a case by case basis. We illustrate such an ad-hoc characterization using the SportSett database (Thomson et al., 2020a). This dataset expands on the game-level data from the RotoWire dataset, adding information on the time dimensional axis. SportSett provides a representation of the more general problem of data-to-text with multidimensional data. While specifics may vary between domains and applications, in this case we are able to define three view types corresponding to different time spans. Firstly, we have the `<Whole-Game>` view, which describes the entities (players and teams), with their statistics (such as points and rebounds) for the game overall (Figure 2: first row). This is similar to the data in the original RotoWire dataset, and is the focal point of the narrative (texts are descriptions of a game). We then define views describing the same entities, but for different time spans. The `<Within-Game>` view describes the entities within parts of each game (such as a half or quarter, Figure 2: middle row), and the `<Between-Game>` view describes entities in past or future games (Figure 2: last row). `<Between-Game>` views can either include information about upcoming games, or aggregate statistics for players over a span of games.

For a given dataset, we are therefore able to create views, by assigning each data record to a view type (e.g. in Figure 2 PTS (points) can be assigned to `<Whole-Game>`) and grouping all data of a single type and belonging to a single entity together, forming one view. This manual partitioning of the data is crucial to allow later control in a manner relevant to the task at hand and the goals of the system's user. In particular, we emphasize that while we apply our framework to a basketball-specific task, our approach is not specific to basketball – or sports – and can work in a number of settings. In the financial domain for instance, expert users of an NLG system could choose to describe a stock's performance against a benchmark (e.g. S&P500) on its own, or compared to stocks from the same domain (e.g. all Pharmaceutical stocks), or from the same country (e.g US stocks). Views could also be created on the time axis, to compare a stock's performance to its own

in previous months or years.

We argue that in order to present domain experts with convincing NLG software for their business use cases, this light involvement on their part is actually beneficial, since they will all require subtle handling of the data that cannot be anticipated while creating the model, but can be passed during training on their in-house dataset via an ad-hoc partitioning they design and understand.

## 3.2 Aligning descriptions to views

To identify which views ground the data to each text, we align spans of tokens (in our case sentences) to one or more views (a view set), as shown in Figures 1 and 2. We consider as view sets any sensible combination of views. An example commonly seen in the reference texts is the two teams' `<Whole-Game>` view set which is often used at the beginning of human written descriptions[3]. This alignment between sentences and view sets could be performed by human annotators, or learned (Perez-Beltrachini and Lapata, 2018). We used relatively simple heuristics (see Appendix E for more details) based on our knowledge of the domain as the focus of this paper is on the generation system.

Figure 3 shows the result of this process, with four example sentences grounded to their respective views. In the second column, for instance, the sentence "Russel Westbrook put up fourteen points" is aligned to Russel Westbrook's `<Whole-Game>` view because the noun phrases "Russel Westbrook" and "14 points" are both valid for that view. Most sentences in the corpus are longer than our brevity-focused examples.

## 3.3 Schemata

View sets can be combined as ordered lists to form document plan schemata in order to guide generation (see Section 4). To train the model, we use the schemata extracted from the reference texts. During inference, we experiment with either following static, extracted, or simple rule-driven schemata (Section 5). While this is another manual input from the expert users of the NLG system, we refer readers to the discussion in Section 6.1 where we discuss the impact of also predicting the plan. Briefly, models that predict the plan are often *dull* in the sense that they always predict the same plan,

---

[3]To reduce complexity, we allow view sets to be comprised of one or two entities (players/teams) for one view type.

| Name | G | W | B |
|---|---|---|---|
| Thunder | | ■ | |
| Heat | | ■ | |
| George | | | |
| Westbrook | | | |
| … | | | |

*The Thunder out-scored the Heat 72-53 in the first half.*

| Name | G | W | B |
|---|---|---|---|
| Thunder | | | |
| Heat | | | |
| George | | | |
| Westbrook | ■ | | |
| … | | | |

*Russel Westbrook put up fourteen points.*

| Name | G | W | B |
|---|---|---|---|
| Thunder | | | |
| Heat | | | |
| George | | | |
| Westbrook | | ■ | |
| … | | | |

*Westbrook had his 4th consecutive triple-double.*

| Name | G | W | B |
|---|---|---|---|
| Thunder | | | ■ |
| Heat | | | |
| George | | | |
| Westbrook | | | |
| … | | | |

*The Thunder head to Boston on Sunday.*

Figure 3: Example sentences aligned to views (Whole-**G**ame, **W**ithin-Game, and **B**etween-Game).

and tend to include a number of irrelevant facts that scored high in evaluation metrics.

# 4 Generating text with a Hierarchical Model leveraging views

As discussed in Section 2, it is unclear to what extent current models deal with extremely large inputs or can utilize views: copy actions are harder to train; memory and compute constraints make training and inference very slow; no model has been proposed to explicitly constrain the order and scope of sentences in generated descriptions.

The hierarchical system of Rebuffel et al. (2020) can be extended to fully leverage view annotations, scaling well with increased input size, generating descriptions following a precise ordering of content. Specifically, the original model was designed with an emphasis on structure, which we use to our advantage to constrain the copy mechanism, as well as input additional data and guide the generation process. We provide an overview of the model here, but refer readers to the original paper for an extensive description.

From a high-level perspective, the system was designed to handle data structured in a hierarchical fashion, i.e. data that can be split into distinct entities, each of them being described by a collection of records, in the form of (key, value) pairs. The system follows a standard encoder-decoder architecture. In particular, it:

1. Encodes each entity independently, as collections of records;

2. Encodes the input data-structure, as a collection of entities;

3. Generates text using hierarchical attention/copy mechanisms:

   (a) An attention distribution is computed over entities;

   (b) Inside each entity, an attention distribution is computed over its records.

Views are integrated to this framework by considering each as an independent entity, with an important distinction to adapt to the now extremely large input size. In their original work, (Rebuffel et al., 2020) encode all input data, which is not satisfactory when adding a large number of views. Akin to teacher forcing (Williams and Zipser, 1989), we encode additional views only when they are relevant to the description at hand. In other words, we always encode all `<Whole-Game>` views which provide a needed high-level understanding of the game, but only encode the `<Between-Game>` and `<Within-Game>` views that ground sentences of the current description, leaving out views which will never be solicited by the decoder.

This prior enables the model to have a broad overview of the game (using `<Whole-Game>` views), while being able to copy specific information from other dimensions for each entity mentioned in the target descriptions. During inference, we only encode the views which are part of the specified control schema, as other views need not be solicited during decoding.

## 4.1 Hierarchical Encoding

Formally, we consider the following setting: • Let $\mathcal{D}$ be a *dataset* which is a collection of aligned (data-structure, description) pairs $(s, y)$.

• A *data-structure* $s$ is originally seen as an unordered set of $I$ views $e_i$. We thus denote $s \coloneqq \{e_1, ..., e_i, ..., e_I\}$.

• Each view $e_i$ is a labelled set of $J_i$ unordered records $\{r_{i,1}, ..., r_{i,j}, ..., r_{i,J_i}\}$; where a record $r_{i,j}$ is defined as a (key, value) pair: $(k_{i,j}, v_{i,j})$. Note that $J_i$ might differ between views. Importantly, the set of records is labelled by the *view-type*, $e_{type}$

Following (Rebuffel et al., 2020), we first encode each view independently, and then together. We denote by $\mathbf{e_i}$ the learned representation of view $e_i$, computed by the high-level encoder.
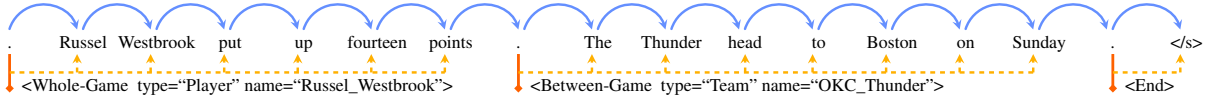
Figure 4: Example decoding of sentences. Decoder is trained to predict next words (blue arrows). During sentence decoding the attention and copy mechanisms are restricted to entities grounded to the current views (see Figure 3), selected at sentence delimiter. We regularize training with an added loss (see Section 4.3). The decoder must predict the next grounding views (orange bars) at each sentence delimiter.

**Encoding a subset of views** Each sentence of the target description is grounded to a small subset of views. To encode high-level information about the entities and the dimension of the considered views, our system relies on self-attention, which enables encoding sets of unordered objects.

Formally, for a given sentence $s$, we denote $\mathcal{G}_s$ the subset of views that grounds the sentence, as explained in Section 3. We compute a fixed-size representation $\mathcal{G}_s$ of this subset using self-attention (SA):

$$\hat{\mathcal{G}}_s = \text{SA}( [\mathbf{e_i}; \ \forall e_i \in \mathcal{G}_s] )$$
$$\mathcal{G}_s = [\hat{\mathcal{G}}_s; \ \mathbf{e_{type}}] \tag{1}$$

where $[ \ \cdot \ ; \ \cdot \ ]$ represents concatenation, and $\mathbf{e_{type}}$ is a learned embedding of the view-type $e_{type}$, which is the same for all views grounding a same sentence.

### 4.2 View-aware decoding

In addition to relying on the encoded content of the views, we guide and constrain the decoding process sentence by sentence.

1. The system uses the encoded high-level information about the entities and the type of the considered subset of views;

2. The system decodes word by word, using the previously encoded information;

3. Attention and copy mechanisms are limited to records from the considered views.

During the decoding process for sentence $s$, the decoder uses the learned representation $\mathcal{G}_s$, computed at Equation 1, to update its hidden state at each decoding step. Recall that the standard decoder of (Rebuffel et al., 2020) is an LSTM which updates its hidden state $\mathbf{d_t}$ using the previously decoded token:

$$\mathbf{d_t} = \text{LSTM}([ \ \mathbf{d_{t-1}}; \ \mathbf{y_{t-1}} \ ]) \tag{2}$$

where $\mathbf{y_{t-1}}$ is the learned embedding of token $y_{t-1}$. In this work, we adapt this update so that the grounding's representation is taken into account:

$$\mathbf{d_t} = \text{LSTM}([ \ \mathbf{d_{t-1}}; \ \mathbf{y_{t-1}}; \ \mathcal{G}_s \ ]) \tag{3}$$

when the current token $y_t$ is from sentence $s$.

**Hierarchical attention constrained on views** To fully leverage the hierarchical structure of their encoder, (Rebuffel et al., 2020) proposed a hierarchical attention mechanism to compute the context fed to the decoder module. The dynamic context is computed in two steps: first attending to views, then to records corresponding to these views. At each decoding step $t$, the model learns a first set of attention scores $\alpha_{i,t}$ over views $e_i$ and a second set of attention scores $\beta_{i,j,t}$ over records $r_{i,j}$ belonging to view $e_i$. The $\alpha_{i,t}$ scores are normalized to form a distribution over all views $e_i$, and $\beta_{i,j,t}$ scores are normalized to form a distribution over records $r_{i,j}$ of view $e_i$. Each view is then represented as a weighted sum of its record embeddings, and the entire data structure is represented as a weighted sum of the view representations.

Formally, the dynamic context is computed as:

$$\mathbf{c_t} = \sum_{i=1}^{I}(\alpha_{i,t}(\sum_j \beta_{i,j,t}\mathbf{r}_{i,j})) \tag{4}$$

$$\text{where } \alpha_{i,t} \propto \exp(\mathbf{d}_t \mathbf{W}_\alpha \mathbf{e}_i) \tag{5}$$

$$\text{and } \beta_{i,j,t} \propto \exp(\mathbf{d}_t \mathbf{W}_a \mathbf{k}_{i,j}) \tag{6}$$

where $\mathbf{d_t}$ is the decoder hidden state at time step $t$, $\mathbf{W}_\alpha \in \mathbb{R}^{d \times d}$ and $\mathbf{W}_\beta \in \mathbb{R}^{d \times d}$ are learnt parameters, $\mathbf{e_i}$ and $\mathbf{k}_{i,j}$ are the computed representation of views and records' keys respectively, $\sum_i \alpha_{i,t} = 1$, and for all $i \in \{1, ..., I\} \sum_j \beta_{i,j,t} = 1$.

In this work, we constrain the attention mechanism such that it is computed only on grounded views. This has the benefit of restraining the copy mechanism to a few specific records, minimizing the number of inaccurate copies. Figure 4 illustrates this mechanism: during the first sentence, only records from `<Whole-Game> Russel_Westbrook` can be attended to, while during the second sentence, only records from

226

```
<Between-Game> OKC_Thunder    can
```
be attended to.

### 4.3 Dual Loss Regularization

For each data-structure $t$ in $\mathcal{D}$, the objective function aims to generate a description $\hat{y}$ as close as possible to the ground truth $y$. This objective function optimizes the following log-likelihood over the whole dataset $\mathcal{D}$:

$$\arg\max_\theta \mathcal{L}(\theta) = \arg\max_\theta \sum_{(t,y)\in\mathcal{D}} \log P(\hat{y}=y \mid t;\theta)$$

where $\theta$ stands for the model parameters and $P(\hat{y}=y \mid t;\theta)$ the probability of the model to generate the adequate description $y$ for table $t$.

Early experiments showed that training using only the standard objective function can be somewhat unstable, and that adding regularization proves useful and increases the model's performances. In practice, in addition to predicting next words, the decoder is also trained to predict next grounding views and view-type at the end of each sentence (illustrated in Figure 4).

Formally, let $\mathcal{L}_w$ refers to the original token-level loss (Section 4.3), and $\mathcal{L}_\ell$ and $\mathcal{L}_e$ refer to two classification losses, on view-type and grounding entities respectively. Then, our model is trained to minimize the following loss:

$$\mathcal{L} = \lambda_1 \mathcal{L}_w + \lambda_2 \mathcal{L}_\ell + \lambda_3 \mathcal{L}_e \qquad (7)$$

where $\sum_i \lambda_i = 1$ and are manually tuned.

## 5 Experimental Setup

**Data**  We use the sports journalism (basketball) dataset SportSett Thomson et al. (2020a), based on the RotoWire datset introduced by (Wiseman et al., 2017). It consists of game statistics paired with human-authored descriptions. The original dataset contained train-test corruption so we partition by season[4] to provide a suitable proxy of a real-world task.

**System comparisons**[5]  We trained one model, then used three different schemata to produce text:

- V-SIMPLE - A simple static schema that is the same for all games (see Figure 5).
- V-EXTENDED - Based on V-SIMPLE, schema varies for each game, adding

---

```
    <Between-Game> / <Within-Game>
```
elaborations to some players. These elaborations are chosen with simple heuristics (e.g. when a player has had 3 double-doubles in his last games, add a `<Between-Game>` elaboration for this player).

- V-GUIDED - Schema automatically extracted from the human authored descriptions. These vary from game to game.

We compare our three variants to two variants of the hierarchical model of (Rebuffel et al., 2020) (that do not have planning modules and represent state-of-the art without such a feature) and two variants of the explicit-planning approach of (Puduppully and Lapata, 2021):

- H-FULL - **H**ierarchical encoder based system of (Rebuffel et al., 2020) but with all views made available to it.
- H-NEXT - **H**ierarchical encoder based system of (Rebuffel et al., 2020), configured as per (Thomson et al., 2020b) to include additional information for the game and next games.
- MP-SIMPLE - System of (Puduppully and Lapata, 2021) using our simple static schema (same as our variant V-SIMPLE).
- MP-GUIDED - System of (Puduppully and Lapata, 2021) using the schema extracted from human authored texts (same as our variant V-GUIDED).

Note that the hierarchical model of (Rebuffel et al., 2020) is not able to handle schema guidance during generation, and that the system of (Puduppully and Lapata, 2021) doesn't differentiate between `<Whole-Game>`, `<Between-Game>`, and `<Within-Game>` dimensions: descriptions generated using this system are identical when using the SIMPLE or EXTENDED schema.

## 6 Human evaluation of factual accuracy

With the gold standard protocol of Thomson and Reiter (2020), human annotators highlight non-overlapping spans of text that contain errors, also selecting an error category. Participants are provided with links to basketball-reference.com, which is at least as complete as the system input data, and much easier for humans to navigate. Four main error categories are used; *Name* (named

**Document plan (schema):**
```
<Whole-Game> OKC Thunder, Miami Heat
<Whole-Game> OKC Thunder, Miami Heat
<Whole-Game> OKC Thunder, Miami Heat
<Whole-Game> OKC Thunder, Miami Heat
<Whole-Game> Paul George
<Whole-Game> Dennis Schroder
<Whole-Game> Russell Westbrook
<Whole-Game> Steven Adams
<Whole-Game> Jerami Grant
<Whole-Game> Patrick Patterson
<Whole-Game> Kelly Olynyk
<Whole-Game> Josh Richardson
<Whole-Game> Hassan Whiteside
<Whole-Game> Bam Adebayo
<Between-Game> OKC Thunder
<Between-Game> Miami Heat
```

The Oklahoma City Thunder (33-18) defeated the Miami Heat (24-26) 118-102 on Friday. The Thunder shot 48 percent from the field and 53 percent from three-point range. They also dominated the rebounding battle, 50-47. The **Heat**[N] also held the Heat to 45 percent shooting and outrebounded them 50-**38**[U]. Paul George led the way for the Thunder with 43 points, seven rebounds, five assists and two steals. Dennis Schroder added 28 points of his own off the bench. Russell Westbrook finished with 14 points, **12**[U] assists and 12 rebounds. Steven Adams scored 13 points and grabbed seven rebounds in 27 minutes. Jerami Grant chipped in eight points, **while Jerami Grant**[W] added eight. Patrick Patterson **led the bench**[W] with five points and two rebounds. Kelly Olynyk was the **only other**[W] **Heat player**[C] in double figures, finishing with 21 points. Josh Richardson scored 18 points on 7-of-18 shooting. Hassan Whiteside had 12 points and 16 rebounds in 19 minutes. Bam Adebayo, it **wasn't enough to overcome the Heat**[W]. Oklahoma City will travel to Boston to take on the Celtics on Sunday. Miami plays host to the Pacers, also on Saturday.

Figure 5: Instructions and generated (V-SIMPLE system) game summary for OKC@MIA on February 1st 2019. **Name**[N], **Number**[U], **Word**[W], and **Context**[C] mistakes are highlighted in the summary.

entities), *Number* (ordinal, cardinal, etc), *Word* (a word or phrase that is not a name or number), *Context* (such as implicature errors). There is also the last resort category *Other*, for text that is nonsensical. Finally, there is a *Not Checkable* category, which covers facts that are impractical or impossible to check using the provided reference data.

We performed our experiment with one annotator per text[6]. The original protocol limited annotators to only 4 prior games before defaulting to *Not Checkable* errors. We felt this was overly restrictive so we asked annotators to check all assertions within the current season. We compared generations for 35 random basketball games, generating a text for each of the 7 systems. A Latin square design was used whereby each participant annotated 5 texts for each system, never seeing the same game twice (245 total annotated texts).

### 6.1 Human Evaluation Results

The best performing systems were our View-based system, and that of Puduppully and Lapata (2021) (when both are provided with a simple schema). They were both significantly different, in terms of the number of errors reported, to all other systems. Figure 6 shows a box plot of error count for each system Figure 5 shows the generation for the V-SIMPLE system, marked up for errors by an annotator, on the same game shown in Figure 1.
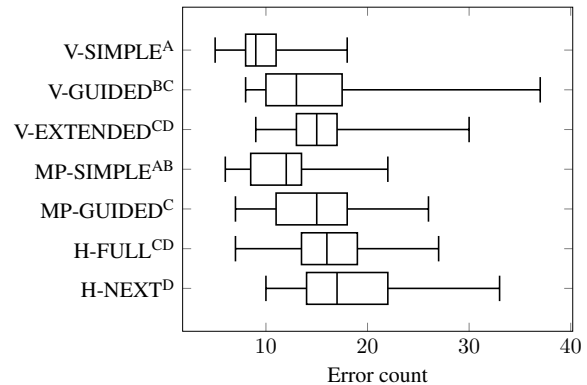


Figure 6: Error count, systems sharing superscript letters are not significantly different.

## 7 Discussion and Limitations

### 7.1 Different Domains and Datasets

Our work focuses on a single domain/dataset. Expanding this to include additional domains and datasets would be useful future work now that a process has been outlined in this paper for working with multi-dimensional data. The MLB dataset (Puduppully et al., 2019a) and the Chart-to-Text dataset of Obeid and Hoque (2020) both exhibit problems related to modelling of time that is seen in the basketball domain.

### 7.2 Large Language Models

The system evaluated in this paper used a transformer encoder and an LSTM decoder. Whilst the LSTM, in particular, may be considered an out-

---

[6]MTurk workers were recruited by the same process as Thomson et al. (2023) who reported high precision and recall of single annotators by this method. It is expensive to run with three annotators per text, therefore we prioritised having more texts over more annotators per-text.

dated NLP model[7], the view-based design could also be implemented using pre-trained large language models such as GPT-3 (Brown et al., 2020) or even using ChatGPT[8]. For systems such as ChatGPT, a series of prompts could be generated that reference views, e.g.,

1. Given *View A*; write a sentence that describes Steph Curry's performance in this game.

2. Then, given *View B*; add a sentence that describes Steph Currys performance over the past 7 games.

Question-answer-based document planning has been explored for tasks such as summarisation (Narayan et al., 2023) and by referencing views it could be applied to data-to-text generation. Note that even with the larger token windows of GPT-4, it remains impossible to feed the model all input data that could be used to extract useful insights, meaning that some form of insight selection is required upstream of the language model. Our view-based approach handles this at a high level, whilst leaving micro-planning and realisation to the neural model.

### 7.3 Generality and View Grounding

View grounding was performed using simple heuristics in order to use the grounded spans of text (sentences) for the downstream generation task. This should be explored and evaluated as a standalone task in future work. It is also unclear for an individual view, what the limits of complexity are. This introduction of control through the use of views is not a loss of generality, but rather a requirement for generating text that is useful and interesting to human readers, rather than "general" but vague or dull.

## 8 Conclusion

Increasing concern has been raised regarding the quality of both task setup (Raji et al., 2021), as well as evaluation in NLP/NLG, with caveats of systems and experimental results often going unreported (Gehrmann et al., 2022). We contribute to the meaningful progress of both. In exploring an alternative task where generation of an exact human reference text is not the goal (it is just part of the available information), we bring the generation process more inline with a real-world problem where control is a major requirement.

We expanded the data-to-text task by considering the mirrored multi-dimensional aspects of both data and text. We have shown that by splitting this extended data into manageable views based on its dimensionality, meaningful control can be introduced over system output without sacrificing factual accuracy. Control comes in the form of ordering views using schema, in the way messages might be ordered in rule-based systems, but allowing the neural model to handle the complexities of micro-planning and surface realization.

The method of splitting data into views could be applied to complex data from sources such as relational databases or multi-dimensional arrays. Our implementation of views is only one possible way to structure the data. We intuit that views that describe a named entity and a set of direct attributes, in any dimension, work well.

Considering the multidimensionality of data and text brings the problem closer to that encountered and addressed by humans in the real world. Enhancing the structure of datasets, as well as designing models that leverage this multidimensionality, will move systems closer to the goal of human-like descriptions of complex data.

## Ethics Statement

Ethical approval was obtained from our ethics review board. We paid our Mechanical Turk participants $8 US per text annotated during qualification, practice, and live work. This equates to approximately $20 per hour and multiple workers indicated to us that this was "about right". Only when qualifying participants uploaded blank documents (no reported errors when there should be around 20) did we reject HITs. In the rare cases that qualified workers made mistakes such as accidentally uploading a blank annotation document rather than the completed one, we still paid them for the HIT and contacted them to get the correct file.

## Online Resources

All code, data, and human evaluation resources will be made available on GitHub[9].

---

[7]Our reasoning for using an LSTM decoder was that it was used by previous models (Puduppully and Lapata, 2021) and we were exploring changes to the encoder component only.

[8]https://openai.com/blog/chatgpt

[9]https://github.com/nlgcat/inlg2023views

## References

Reinald Kim Amplayo, Stefanos Angelidis, and Mirella Lapata. 2021. Aspect-controllable opinion summarization. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6578–6593, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Thiago Castro Ferreira, Diego Moussallem, Emiel Krahmer, and Sander Wubben. 2018. Enriching the WebNLG corpus. In *Proceedings of the 11th International Conference on Natural Language Generation*, pages 171–176, Tilburg University, The Netherlands. Association for Computational Linguistics.

Robert Dale. 2020. Natural language generation: The commercial state of the art in 2020. *Natural Language Engineering*, 26(4):481–487.

Li Dong, Shaohan Huang, Furu Wei, Mirella Lapata, Ming Zhou, and Ke Xu. 2017. Learning to generate product reviews from attributes. In *EACL*.

Ondřej Dušek, David M. Howcroft, and Verena Rieser. 2019. Semantic noise matters for neural natural language generation. In *Proceedings of the 12th International Conference on Natural Language Generation*, pages 421–426, Tokyo, Japan. Association for Computational Linguistics.

Ondřej Dušek, Jekaterina Novikova, and Verena Rieser. 2018. Findings of the E2E NLG challenge. In *Proceedings of the 11th International Conference on Natural Language Generation*, pages 322–328, Tilburg University, The Netherlands. Association for Computational Linguistics.

Jessica Ficler and Yoav Goldberg. 2017. Controlling linguistic style aspects in neural language generation. In *Workshop on Stylistic Variation @ ACL*.

Katja Filippova. 2020. Controlled hallucinations: Learning to generate faithfully from noisy data. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 864–870, Online. Association for Computational Linguistics.

J. Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological Bulletin*, 76:378–382.

Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez-Beltrachini. 2017. Creating training corpora for NLG micro-planners. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 179–188, Vancouver, Canada. Association for Computational Linguistics.

Albert Gatt and Emiel Krahmer. 2018. Survey of the state of the art in natural language generation: Core tasks, applications and evaluation. *J. Artif. Int. Res.*, 61(1):65–170.

Sebastian Gehrmann, Elizabeth Clark, and Thibault Sellam. 2022. Repairing the cracked foundation: A survey of obstacles in evaluation practices for generated text.

Heng Gong, Xiaocheng Feng, Bing Qin, and Ting Liu. 2019. Table-to-text generation with effective hierarchical encoder on three dimensions (row, column and time). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3143–3152, Hong Kong, China. Association for Computational Linguistics.

Zhiting Hu, Zichao Yang, Xiaodan Liang, Ruslan Salakhutdinov, and Eric P. Xing. 2017. Toward controlled generation of text. In *ICML*.

Anna S. Law, Y. Freer, J. Hunter, R. Logie, N. McIntosh, and John Quinn. 2005. A comparison of graphical and textual presentations of time series data to support medical decision making in the neonatal intensive care unit. *Journal of Clinical Monitoring and Computing*, 19:183–194.

Rémi Lebret, David Grangier, and Michael Auli. 2016. Neural text generation from structured data with application to the biography domain. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1203–1213, Austin, Texas. Association for Computational Linguistics.

Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4582–4597, Online. Association for Computational Linguistics.

Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021. Pretrain, prompt, and predict: A systematic survey of prompting methods in natural language processing. *arXiv preprint arXiv:2107.13586*.

Nitika Mathur, Timothy Baldwin, and Trevor Cohn. 2020. Tangled up in BLEU: Reevaluating the evaluation of automatic machine translation evaluation metrics. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4984–4997, Online. Association for Computational Linguistics.

Shashi Narayan, Joshua Maynez, Reinald Kim Amplayo, Kuzman Ganchev, Annie Louis, Fantine Huot, Anders Sandholm, Dipanjan Das, and Mirella Lapata. 2023. Conditional generation with a question-answering blueprint.

Feng Nie, Jinpeng Wang, Jin-Ge Yao, Rong Pan, and Chin-Yew Lin. 2018. Operation-guided neural networks for high fidelity data-to-text generation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3879–3889, Brussels, Belgium. Association for Computational Linguistics.

Feng Nie, Jin-Ge Yao, Jinpeng Wang, Rong Pan, and Chin-Yew Lin. 2019. A simple recipe towards reducing hallucination in neural surface realisation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2673–2679, Florence, Italy. Association for Computational Linguistics.

Jason Obeid and Enamul Hoque. 2020. Chart-to-text: Generating natural language descriptions for charts by adapting the transformer model. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 138–147, Dublin, Ireland. Association for Computational Linguistics.

Ankur P Parikh, Xuezhi Wang, Sebastian Gehrmann, Manaal Faruqui, Bhuwan Dhingra, Diyi Yang, and Dipanjan Das. 2020. ToTTo: A controlled table-to-text generation dataset. In *Proceedings of EMNLP*.

Laura Perez-Beltrachini and Mirella Lapata. 2018. Bootstrapping generators from noisy data. In *NAACL-HLT*.

Ratish Puduppully, Li Dong, and Mirella Lapata. 2019a. Data-to-text generation with content selection and planning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 6908–6915.

Ratish Puduppully, Li Dong, and Mirella Lapata. 2019b. Data-to-text generation with entity modeling. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2023–2035, Florence, Italy. Association for Computational Linguistics.

Ratish Puduppully and Mirella Lapata. 2021. Data-to-text generation with macro planning. *Transactions of the Association for Computational Linguistics*, 9:510–527.

Inioluwa Deborah Raji, Emily M. Bender, Amandalynne Paullada, Emily Denton, and Alex Hanna. 2021. Ai and the everything in the whole wide world benchmark.

Clément Rebuffel, Marco Roberti, Laure Soulier, Geoffrey Scoutheeten, R. Cancelliere, and P. Gallinari. 2021. Controlling hallucinations at word level in data-to-text generation. *ArXiv*, abs/2102.02810.

Clément Rebuffel, Laure Soulier, Geoffrey Scoutheeten, and Patrick Gallinari. 2020. A hierarchical model for data-to-text generation. In *Advances in Information Retrieval - 42nd European Conference on IR Research, ECIR 2020*, pages 65–80.

Ehud Reiter. 2007. An architecture for data-to-text systems. In *proceedings of the eleventh European workshop on natural language generation (ENLG 07)*, pages 97–104.

Ehud Reiter and Robert Dale. 1997. Building applied natural language generation systems. *Natural Language Engineering*, 3(1):57–87.

Ehud Reiter, S. Sripada, J. Hunter, Jin Yu, and I. Davy. 2005. Choosing words in computer-generated weather forecasts. *Artif. Intell.*, 167:137–169.

Anna Rogers. 2021. Changing the world by changing the data. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 2182–2194, Online. Association for Computational Linguistics.

M. S. Silberman, B. Tomlinson, R. LaPlante, J. Ross, L. Irani, and A. Zaldivar. 2018. Responsible research with crowds: Pay crowdworkers at least minimum wage. *Commun. ACM*, 61(3):39–41.

André Luiz Rosa Teixeira, João Campos, Rossana Cunha, Thiago Castro Ferreira, Adriana Pagano, and Fabio Cozman. 2020. DaMata: A robot-journalist covering the Brazilian Amazon deforestation. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 103–106, Dublin, Ireland. Association for Computational Linguistics.

Craig Thomson and Ehud Reiter. 2020. A gold standard methodology for evaluating accuracy in data-to-text systems. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 158–168, Dublin, Ireland. Association for Computational Linguistics.

Craig Thomson and Ehud Reiter. 2021. Generation challenges: Results of the accuracy evaluation shared task. In *Proceedings of the 14th International Conference on Natural Language Generation*, pages 240–248, Aberdeen, Scotland, UK. Association for Computational Linguistics.

Craig Thomson, Ehud Reiter, and Somayajulu Sripada. 2020a. SportSett:basketball - a robust and maintainable data-set for natural language generation. In *Proceedings of the Workshop on Intelligent Information Processing and Natural Language Generation*, pages 32–40, Santiago de Compostela, Spain. Association for Computational Lingustics.

Craig Thomson, Ehud Reiter, and Barkavi Sundararajan. 2023. Evaluating factual accuracy in complex data-to-text. *Computer Speech & Language*, 80.

Craig Thomson, Zhijie Zhao, and Somayajulu Sripada. 2020b. Studying the impact of filling information gaps on the output quality of neural data-to-text. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 35–40, Dublin, Ireland. Association for Computational Linguistics.

Ashish Upadhyay and Stewart Massie. 2022. Content type profiling of data-to-text generation datasets. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 5770–5782, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

Chris van der Lee, Emiel Krahmer, and Sander Wubben. 2017. PASS: A Dutch data-to-text system for soccer, targeted towards specific audiences. In *Proceedings of the 10th International Conference on Natural Language Generation*, pages 95–104, Santiago de Compostela, Spain. Association for Computational Linguistics.

Emiel van Miltenburg, Miruna Clinciu, Ondřej Dušek, Dimitra Gkatzia, Stephanie Inglis, Leo Leppänen, Saad Mahamood, Emma Manning, Stephanie Schoch, Craig Thomson, and Luou Wen. 2021. Underreporting of errors in NLG output, and what to do about it. In *Proceedings of the 14th International Conference on Natural Language Generation*, pages 140–153, Aberdeen, Scotland, UK. Association for Computational Linguistics.

Emiel van Miltenburg, Miruna Clinciu, Ondřej Dušek, Dimitra Gkatzia, Stephanie Inglis, Leo Leppänen, Saad Mahamood, Stephanie Schoch, Craig Thomson, and Luou Wen. 2023. Barriers and enabling factors for error analysis in nlg research. *Northern European Journal of Language Technology*, 9.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.

Hongmin Wang. 2019. Revisiting challenges in data-to-text generation with fact grounding. In *Proceedings of the 12th International Conference on Natural Language Generation*, pages 311–322, Tokyo, Japan. Association for Computational Linguistics.

Ronald J. Williams and D. Zipser. 1989. A learning algorithm for continually running fully recurrent neural networks. *Neural Computation*, 1:270–280.

Sam Wiseman, Stuart Shieber, and Alexander Rush. 2017. Challenges in data-to-document generation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2253–2263, Copenhagen, Denmark. Association for Computational Linguistics.

Xinnuo Xu, Ondřej Dušek, Verena Rieser, and Ioannis Konstas. 2021. AggGen: Ordering and aggregating while generating. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1419–1434, Online. Association for Computational Linguistics.

# A Automated Metric Results

We include results on common metrics of BLEU, Relation Generation (RG), Content Selection (CS), and Content Ordering (CO) for this task in Table 1. Automatic metrics are often expected in NLP papers, although their usefulness in this domain is limited at best. We include them in the appendix for this reason.

The V-SIMPLE and MP-SIMPLE systems, based on simple schema, had the highest RG scores, and hierarchical systems the lowest. Interestingly, CO scores are highest when models follow extracted schema from gold texts.

BLEU scores are within a narrow range, with Mathur et al. (2020) having shown that larger differences are required in order to make judgments. The information extraction based metrics prove more useful, with Wiseman et al. (2017) stating that their results were generally inline with their human evaluation. However, Thomson and Reiter (2021) observed that state-of-the-art metrics can detect simple errors, but struggle with more complex semantic and contextual errors. It is also worth noting that running BLEU on a deranged copy of the test set (comparing each game with a random game other than itself) can yield BLEU scores in the region of 8.0 to 10.0, simply due to common terminology and syntax.

| System | RG | | CS | | | CO | BLEU |
|---|---|---|---|---|---|---|---|
| | P% | # | P% | R% | F1 | | |
| REF | 0.84 | 26.84 | - | - | - | - | - |
| V-SIMPLE | 0.87 | 26.21 | 0.60 | 0.57 | 0.58 | 0.21 | 19.68 |
| V-GUIDED | 0.81 | 17.56 | **0.71** | 0.48 | 0.57 | **0.30** | 17.29 |
| V-EXTENDED | 0.84 | 27.06 | 0.57 | 0.58 | 0.57 | 0.21 | 21.90 |
| MP-SIMPLE | **0.88** | 43.27 | 0.48 | **0.73** | 0.58 | 0.22 | 21.52 |
| MP-GUIDED | 0.82 | 30.02 | 0.60 | 0.67 | **0.63** | **0.30** | **22.27** |
| H-FULL | 0.76 | 27.76 | 0.42 | 0.47 | 0.44 | 0.16 | 17.73 |
| H-NEXT | 0.77 | 23.09 | 0.51 | 0.47 | 0.49 | 0.18 | 21.22 |

Table 1: Automatic metric results for all systems.

# B Content Ordering Experiment

This experiment aims to determine whether sentences in generated summaries are in the correct order. In designing this experiment we had two main concerns. Firstly, inter-annotator agreement should at least be moderate, ideally high. This precludes designs where participants are free to rearrange all sentences; the large number of permutations increases the likelihood of disagreement. Secondly, it should be possible to perform meaningful error analysis in order to better understand both the systems, and the protocol itself. This

rules out Likert-based approaches because, with paragraph-sized generations, it is impossible to tell which part of the summary caused a participant to score the text in the way they did. Likert ratings have been shown to have poor agreement in this domain (Puduppully and Lapata, 2021).

## B.1 Design

We presented generated summaries to participants with the first two sentences highlighted as 'the beginning', the final two sentences highlighted as 'the end', and everything in between highlighted as 'the middle'. We then asked participants, for each of the four sentences in the beginning and end, whether it should:

- **Remain** where it is.
- be **Transposed** with its partner, i.e., the other sentence from the beginning or end.
- be moved to the middle, a **Short** distance.
- be moved to the opposing end of the summary, a **Long** distance.

When asked if sentences should be moved to another section participants did not specify exactly where, simply which other section. We also asked the middle was in an acceptable order (Yes/No).

Participants were placed into 35 non-exclusive groups (the number of combinations of size three for 7 participants). Each group evaluated a summary from each of the 7 systems, such that 245 unique summaries were evaluated by 3 annotators.

## B.2 Results

For content ordering, we first consider whether participants believed a sentence should be moved to a different section. Inter-annotator agreement by Fleiss Kappa (Fleiss, 1971) was 0.591, indicating a moderate agreement. However, this falls to 0.469 when we consider the *Short/Long* move distances, and to 0.350 if we also consider transposition of beginning/end sentence pairs (p-value was less than 0.001 in all cases). This confirms our design assumption that allowing participants to freely rearrange texts of this length would result in low or no agreement. We did run an experiment where different participants (MTurk masters with US high-school diplomas) were asked to rate how readable and understandable generations were. Agreement for this was even lower, below 0.2, and results are not included for that reason.

| System | Long | Short | Transpose | Remain |
|---|---|---|---|---|
| V-SIMPLE | 1 | 3 | 55 | 361 |
| V-GUIDED | 1 | 33 | 10 | 372 |
| V-EXTENDED | 0 | 10 | 59 | 351 |
| MP-SIMPLE | 3 | 15 | 4 | 398 |
| MP-GUIDED | 3 | 53 | 7 | 353 |
| H-FULL | 2 | 65 | 1 | 352 |
| H-NEXT | 1 | 88 | 4 | 327 |

Table 2: Number of sentences that annotators would move, by destination.

## B.3 Conclusion

The results in Table 2 show that all models do a good job at avoiding *Long* errors, that is they do not confuse the beginning of the narrative with the end. The simple schema of both V-SIMPLE and MP-SIMPLE have fewer *Short* errors, especially compared with the hierarchical encoder systems. Our models in V-SIMPLE and V-EXTENDED mode *Transpose* sentences in the *beginning* or *end* with higher frequency. Looking into this further, our schema (for both models) was set to realize the upcoming game for the winning team in the *Penultimate* sentence, then the losing team in the *Final* sentence. This was deemed incorrect by some annotators (the losing teams players are usually discussed immediately before the end, therefore the context at that stage is the losing team). Our system is capable of adjusting for this, with a simple schema change reversing the order of these sentences. The MP-SIMPLE system does not have the fine-grained control to constrain generation to two separate sentences, therefore it frequently discusses both teams upcoming games in a single *Final* sentence and does not encounter this *Transpose* problem as often as our models. It is also unclear how the *Short* errors of such a system could be corrected.

This experiment is included in the appendix because whilst it was unsuccessful at demonstrating a difference between systems (agreement was low), it does provide some insight and with some refinement of experimental design could be a useful approach (agreement was not so low that there are no possible pathways to higher agreement).

## C  Post-hoc error analysis

In addition to the quantitative data, our accuracy evaluation yielded qualitative data in the form of free-text comments that annotators could leave when reporting each error. We therefore performed an error analysis, something that is crucial to to gain insight into where our systems are failing (van Miltenburg et al., 2021, 2023). With the MP-SIMPLE and MP-SIMPLE systems some annotators queried the protocol because some names were spelled incorrectly. This had not been a problem for word-based systems, but since the system of Puduppully and Lapata (2021) operates at the subword level, it would sometimes generate texts that contained out of vocabulary words once subwords were reconstructed. An example can be seen in the sentence: *"Well ell ell ell ell ell ell ell Carter[N] , as he scored 25 points to go along with eight rebounds and five assists ."*, where *"ell"* is an out-of-vocabulary word. The annotator for this sentence marked it as an error, leaving the mildly derisive comment of *"more commonly referred to as just Wendell Carter[N]"*. Upon further investigation, this problem is not uncommon in the generations of this system, yet it would be missed by the RG metric and at times our human evaluation as well[10]. In one of the worst cases (from the full test set, not an item from our human evaluation), the complete generation was: *"The Miami Heat ( 27 - 33 ) defeated the Golden State Warriors ( 43 - 18 ) 126 - 125 on Friday . Justise Winslow and Bam AAAAAAAAAAA"*, followed by the letter *'b'* repeated 808 times. Our view based systems also struggled at times to generate full sentences about players such as Bam Adebayo, who had not been seen during training. For example, one output was *"Bam Adebayo, it wasn't enough to overcome the Heat[W]."*, where the model knew it should generate a sentence about Bam Adebayo, but did not include any statistics. It is possible the models are relying on the values of the player name field rather than generalizing.

To gain further insight, we performed some automated error analysis on outputs from the full test set (2018 season). Table 4 shows the average token counts and out of vocabulary[11] tokens for the generations of each system. Our view based systems each generated a small number of out-of-vocab tokens by erroneously copying boolean values from the input data (we would fix this by

---

[10]The factual accuracy annotation instructions of Thomson and Reiter (2020) ask annotators to ignore spelling, syntax and grammar, so some annotators did not mark these as errors (if they could make out which player was being referred to).

[11]A vocabulary was created using all test data values, training data texts and a range of numbers in word and digit form.

| System | NAME | NUMBER | WORD | CONTEXT | OTHER | NOT CHECKABLE | TOTAL |
|---|---|---|---|---|---|---|---|
| V-SIMPLE | 44 | 115 | 134 | 16 | 19 | 11 | 339 |
| V-GUIDED | 76 | 233 | 153 | 18 | 16 | 14 | 510 |
| V-EXTENDED | 60 | 218 | 206 | 18 | 30 | 17 | 549 |
| MP-SIMPLE | 195 | 79 | 91 | 22 | 6 | 5 | 398 |
| MP-GUIDED | 186 | 129 | 134 | 33 | 29 | 2 | 513 |
| H-FULL | 109 | 232 | 186 | 14 | 32 | 2 | 575 |
| H-NEXT | 113 | 232 | 243 | 24 | 38 | 2 | 652 |

Table 3: Errors for each system by type. Systems that were guided by simple schema (V-SIMPLE, MP-SIMPLE) produced the fewest factual mistakes whilst offering the most control.

only including lexical values as input data values). The references texts had out-of-vocab tokens because human authors are not constrained to the set of training words. The MP-SIMPLE and MP-SIMPLE systems both had more out of vocabulary words. Also shown is a count of singleton trigrams (where all three tokens in the trigram are identical), a measure of repetition, where gain the MP-SIMPLE and MP-SIMPLE systems had higher mean counts. In both cases, this is likely due to the incorrect recombination of subwords. It may be possible to adjust the training of models to aleviate this, but it is important to note that automatic metrics all miss this kind of error and it was only found because of our error analysis of human annotated errors.

Shot breakdowns, which are a type of domain specific syntax breaking down the shooting of a player using between 2 and 6 numbers, e.g. *"(4-8 FG, 1-4 3Pt, 2-2 FT)"*, were also counted in Table 4. The number of shot breakdowns (extracted by regular expression) included by the MP-SIMPLE and MP-GUIDED systems could explain part of the increased RG# seen in Table 1. They densely transcribe either 2, 4, or 6 numeric facts yet are simple (once the decision has been made to include one, the structure is deterministic). Systems learn to generate so many shot breakdowns because that they are present in the training data, although they are seldom found in the test set reference texts from the 2018 season. This could be explained by drift due to a change in the specific authors writing the reference texts during that year (Upadhyay and Massie, 2022).

## D  Crowd-sourced worker recruitment

Participants were recruited on the Amazon Mechanical Turk platform. We used the recruitment policy of Thomson and Reiter (2020) participants were required to hold a US Bachelors degree,

be US residents, and be Mechanical Turk Masters workers (a qualification issues by Amazon for high worker reliability). In addition, candidates had to complete a (paid) custom qualification exercise. Fair treatment of crowd-sourced workers is important (Silberman et al., 2018) both from an ethical standpoint and to ensure high quality work. We aimed to pay workers approximately US$20 per-hour for their time, which meant paying $8 for each of the 35 factual accuracy annotation tasks they completed, these take 20-25 minutes to complete. We paid $2 for each of the ordering tasks which take 5-6 minutes to complete. We also paid the same for the any practice work. The same 7 participants completed all work for both our factual accuracy and ordering experiments.

## E  View Grounding

Given a sentence, we consider all possible view sets as candidates for grounding. We propose to judge the alignment between one view set and the sentence as inversely proportional to the number of *alignment errors* it would entail. An alignment error simply refers to any token that could belong to one of the generated noun phrases but cannot be justified by the data contained in the view set.

To identify individual alignment errors, we first use a simple rule-based system to generate noun phrases based on the data within the view set. This includes phrases based on statistics like '14 points', or alternative forms such as '14-point'. We also include those derived from multiple statistics, e.g., 'double-double'. Named entities are also included, for example, 'Russel Westbrook'. This does introduce a requirement of manual definition, but generating noun phrases for data is a much simpler task than constructing grammar and narrative to connect them. We take the best of both rules and neural, defining that what which is simple and learning that which is complex or time-

| System | Token Count | | Out-of-Vocab Count | | Singleton Trigram Count | | Shot Breakdown Count | |
|---|---|---|---|---|---|---|---|---|
| | mean | stdev | mean | stdev | mean | stdev | mean | stdev |
| V-SIMPLE | 276 | 30 | 0.027 | 0.19 | 0.04 | 0.237 | 0.151 | 0.663 |
| V-GUIDED | 241 | 48 | 0.022 | 0.147 | 0.013 | 0.145 | 0.178 | 0.636 |
| V-EXTENDED | 340 | 33 | 0.026 | 0.179 | 0.05 | 0.253 | 0.229 | 0.817 |
| MP-SIMPLE | 292 | 62 | 2.673 | 3.625 | 0.386 | 3.008 | 1.551 | 2.185 |
| MP-GUIDED | 309 | 95 | 2.108 | 5.449 | 0.63 | 6.157 | 0.83 | 1.948 |
| H-FULL | 366 | 71 | 0 | 0 | 0 | 0 | 0.191 | 0.774 |
| H-NEXT | 386 | 94 | 0 | 0 | 0 | 0 | 1.142 | 2.008 |
| GOLD | 339 | 39 | 0.618 | 0.958 | 0 | 0 | 0.008 | 0.134 |

Table 4: Mean count and standard deviation of tokens, out-of-vocabulary tokens, singleton trigrams (where the set of tokens within the trigram is a singleton), and shot breakdowns per-text.

consuming. Each sentence is parsed token-wise, and once a known noun-phrase (from a global list) is started, it must be able to continue within that view ('14' can continue as '14 points' or '14 - point'), or conclude ('14 - point' must conclude as there is no possible continuation), otherwise it is an error. There will be a small number of cases where the grounding cannot be narrowed down to 1 or 2 compatible views. However, all we require is enough correctly grounded views to introduce a training signal. When there is ambiguity, a model can be instructed to not update weights.

We conclude the view set selection procedure by selecting the smallest one, i.e. in this case the singleton of Westbrook's `<Whole-Game>` view (which had zero errors).

# Memories for Virtual AI Characters

**Fabian Landwehr**
ETH Zürich
`fabian.landwehr@inf.ethz.ch`

**Erika Varis Doggett**
`erikavaris@gmail.com`

**Romann M. Weber**
`romann.weber@gmail.com`

## Abstract

In this paper, we present a system for augmenting virtual AI characters with long-term memory, enabling them to remember facts about themselves, their world, and past experiences. We propose a memory-creation pipeline that converts raw text into condensed memories and a memory-retrieval system that utilizes these memories to generate character responses. Using a fact-checking pipeline based on GPT-4 (OpenAI), our evaluation demonstrates that the character responses are grounded in the retrieved memories and maintain factual accuracy. We discuss the implications of our system for creating engaging and consistent virtual characters and highlight areas for future research, including large language model (LLM) guardrailing and virtual character personality development.

## 1 Introduction

Memory is arguably one of the keys to personality, influencing how we think and respond to external stimuli (Thorne). It is also pivotal for maintaining a coherent and engaging conversation (Baddeley). Such characteristics of memory hold considerable implications for the field of natural language processing (NLP), particularly for AI chatbots, which most recently have emerged as a significant form of human-computer interaction (Brown et al.; Adamopoulou and Moussiades; Chen et al.).

Recently, chat applications utilizing large language models (LLMs) have incorporated a type of short-term memory into their design, namely the chat history, which is included in the context window of the LLM (Zaib et al.). However, the integration of long-term memory into these systems is not yet standard. One reason is that without supporting systems, long-term memory can only be achieved by frequently retraining the underlying LLM, which is impractical due to prohibitive computational costs and associated emissions (Patterson et al.; Sharir et al.).

Most AI chatbots today are designed to function as assistants, parsing user inputs as commands or providing answers to queries, often relying on the support of information-retrieval systems (Adamopoulou and Moussiades).[1] However, envisioning a slight shift in this objective—using LLMs to generate responses for virtual AI characters—presents new challenges and opportunities.

The incorporation of long-term memory becomes critical in this context. To ensure that a virtual character appears human-like to a user, it must accurately remember facts about its universe and recall past interactions. Achieving this could elevate the user experience to a new level, creating a dynamic, immersive experience that adapts to the user (Madotto et al.; Thue et al.).

This paper presents a system to augment virtual AI characters with persistent, long-term memories. At a high level, we achieve this by generating character responses with an LLM, augmenting the prompts with relevant *memories*—condensed versions of character experiences, knowledge, or descriptions.[2] Figure 1 shows how memories can influence character responses.

Our exploration spans relevant related work, and we provide a detailed explanation of our system, encompassing retrieval systems, potential LLM steering, memory creation, and forgetting mechanisms. We evaluate our system based on two characters: *Sherlock Holmes*, a well-known character present in the LLM's training data and implicitly represented by the model's weights, and *Amina Osei*, a newly created character with no such representation. The evaluation focuses on the alignment between the character responses and the retrieved memories and uses a semi-automated fact-checking pipeline that leverages GPT-4.

---

[1] Examples include Apple's Siri, Google's Assistant, and Amazon's Alexa.

[2] We use OpenAI's `gpt-3.5-turbo` as the underlying LLM (`https://platform.openai.com/docs/models/gpt-3-5`).
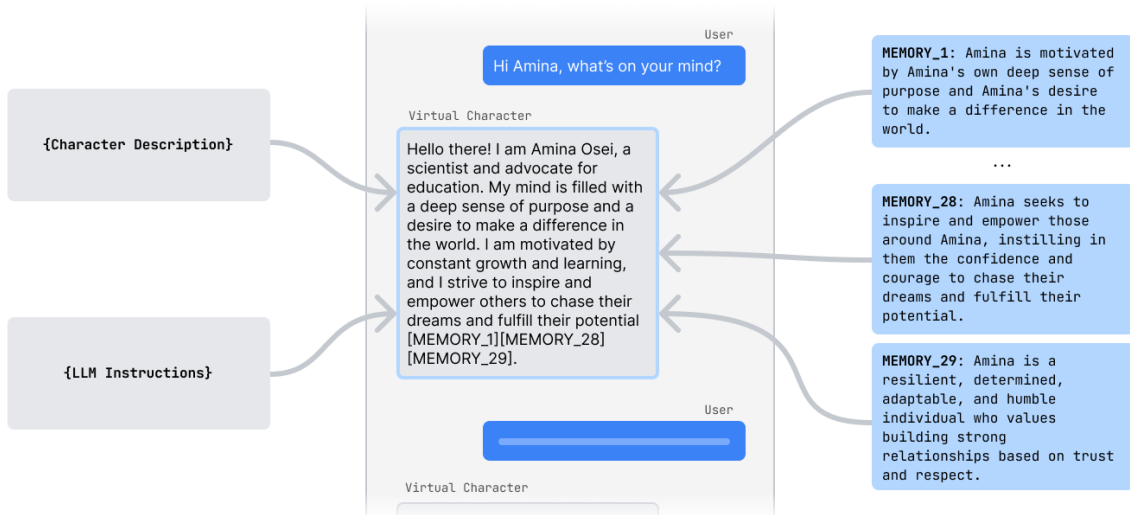
Figure 1: *Memories for virtual AI characters:* This figure illustrates how virtual character responses are generated. When answering, characters have access to a list of relevant memories, a description of themselves, and the chat history. Section 3 explains the proposed system in detail.

This paper makes several contributions. First, it demonstrates how LLMs can be substantially enhanced by integrating them into a more extensive system. Second, it introduces a process of creating memories from unstructured text, simplifying the character-creation process. Third, it outlines partitioning memories into multiple "buckets," offering greater flexibility and potential for dynamic story progression. Fourth, it proposes a new query-generation approach resulting in targeted retrieval of relevant memories. Furthermore, it presents a forgetting model that mimics human memory patterns and allows customization of character memory traits.

## 2 Related Work

The proliferation of natural language (NL) interfaces in computer systems, including systems such as Siri, Alexa, and Google Assistant, has opened new avenues for intuitive interaction. Typically, these interfaces focus on question answering, command parsing, and named entity recognition, assisting users in a variety of tasks (Hussain et al.; Zemčík).

Another strand of research centers around developing conversational chatbots capable of realistic user interactions. Applications range from in-game non-player-characters (NPCs) (Park et al.) to virtual companions, as popularized in various science-fiction movies.[3] Our work intersects these

two streams of research, aiming to create virtual characters capable of recalling factual information about themselves and their past experiences, effectively giving them long-term memory.

### 2.1 Character-Simulation Approaches

**Traditional Methods**

Historically, character simulation and chatbots relied on scripted dialogues or dialogue trees created by human designers (Hussain et al.; Zemčík). These systems employed key-phrase insertion for personalization (Weizenbaum). However, their lack of variability was a limitation, as no dialogue tree can accommodate all potential user topics. Consequently, interactions with these chatbots can feel robotic, as they cannot improvise.

**Recent Approaches**

Recent efforts such as OpenAI's ChatGPT and Microsoft's Bing have leveraged generative language models to generate responses.[4] However, issues such as excessive variability and hallucinations hinder their use for persona-based characters (Greshake et al.; Welleck et al.; Ji et al.). These issues have led to the exploration of knowledge-grounding techniques for generative language models, which we discuss further in section 2.2.

Most recently, the concept of generative agents

---

[3] Popular movies featuring AI companions include *Her* (2013) and *Blade Runner 2049* (2017). Examples of businesses that offer virtual companions include `replika.com` and `myanima.ai`.

[4] ChatGPT: `https://openai.com/blog/chatgpt`, Bing: `https://www.microsoft.com/en-us/edge/features/bing-chat`

has gained traction (Park et al.). Their focus is autonomously making decisions and planning actions in a virtual world based on past experiences and conversations. Some projects also attempt to let agents interact with the real world (Sig; Nakajima). Those approaches use LLMs in an infinite loop, where the output of one call can influence the input of the next. Similarly, our work uses LLMs as part of a broader software system to guide and limit the generated outputs. However, we emphasize artist-crafted virtual characters that engage in user conversations rather than autonomous agents.

## 2.2 Question-Answering (QA) Systems and Search Engines

This branch involves systems designed to provide factual answers to user queries. While not solely focused on question answering, our project integrates relevant character and context information into responses.

Transformer-based (Vaswani et al.) LLMs are now used in search engines such as Bing. Instead of encoding all information in the model's weights, an existing information-retrieval system is queried, and the results are taken as inputs for the LLM. Similarly, LLM wrapper libraries such as Langchain also provide implementations for external databases and memory to increase retrieval capabilities.[5] This splitting of knowledge and language capabilities simplifies the process of updating information (Wang et al.; Levine et al.).

We utilize a similar approach by storing character-relevant facts in a vector search database. Simultaneously, we differentiate our work by making memories human-like. We apply them to virtual characters, adding more sophisticated memory formation and retrieval mechanisms, including a forgetting function.

According to Lan et al., LLMs are used in question-answering tasks via *parse-then-execute* or *retrieval-and-rank* strategies. The former involves comprehension and query formulation, and the latter directly uses user input to query a knowledge base and subsequently filters the results. We combine these strategies in our system.

Alternative methods for knowledge grounding of language models exist, which include:

- Replacing some attention-based feed-forward networks with updateable context vectors

(Sukhbaatar et al.).

- Nearest-neighbor language models to bias language generation toward specific datasets, simplifying updates (Khandelwal et al.).

- Retrieving text chunks similar to the input and using cross-attention for prediction, offloading LLM storage and facilitating updates (Borgeaud et al.).

While custom architectures offer the advantage of finer-grained control, we use general-purpose LLMs to be able to use powerful, pre-trained SOTA models, facilitating rapid prototyping.

## 3 Proposed System

The overarching idea of the proposed system is to generate character responses by prompting a general-purpose LLM with a dynamically created prompt based on the chat context and relevant memories. This section explains how different system components work together to create this prompt and how interactions are distilled into memories.

Figure 2 explains the system on a high level. A detailed description of the system's most central components and stages follows. Note: For all of our core system components described in this section, we use gpt-3.5-turbo (Brown et al.) as the base LLM.[6]

### 3.1 Chat History

The chat history contains all messages of the current conversation which still need to be converted to memories. This chat history is used throughout the system as a form of short-term memory and serves as the immediate context for the last user message.

### 3.2 Query Creator

Once the chat history has been updated with the most recent user message, the query creator generates a search query which is used to retrieve relevant memories. Appendix A.4 shows the prompt to create a search query. We noticed that using a low temperature was essential to making the query generation predictable.[7]

---

[5]https://python.langchain.com/docs/modules/memory/

[6]https://platform.openai.com/docs/models/gpt-3-5

[7]*Temperature* is a hyperparameter in LLMs controlling output randomness. Lower values yield more deterministic outputs focused on the highest probability predictions, while higher values promote diversity. A temperature of 1.0 maintains the model's original output distribution.
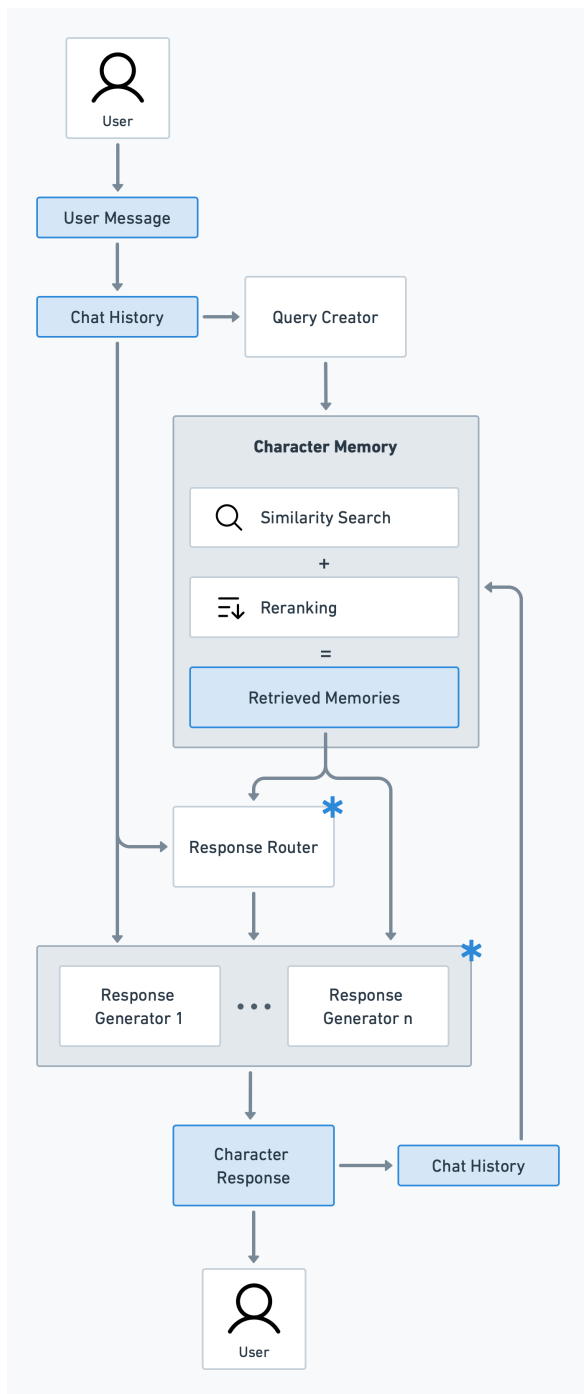
Figure 2: *Overview of the System Architecture:* When a user sends a message to a virtual character, a query for needed information is generated. Using this query, relevant memories are retrieved by doing a similarity search and re-ranking the results. Finally, an LLM is used to generate the character response. Interactions are later synthesized into new memories. For details, see section 3. The response routing (*) is part of our theoretical framework but not implemented in the current system; see section 5.

We retrieve memories based on similarity to the search query instead of similarity to the user input because, in some cases, the user input contains insufficient information. For example, if the latest user message was "Who is she?" then without knowledge of previous messages, it is unclear to whom the pronoun "she" refers. As a result, retrieved memories would likely be irrelevant.

### 3.3 Retrieving Memories

The character's memory consists of multiple knowledge sources. Each knowledge source conceptually maps to one source text. For example, one knowledge source could correspond to one Wikipedia article, another could correspond to a character description written by a novelist, and yet another could correspond to past conversations with the current user.

Each available knowledge source is queried using the created query to retrieve relevant memories. More precisely, a kNN search is performed on each knowledge source, using an embedding of the query and embeddings of the memory contents. This way, the memories most semantically similar to the query are found. For simplicity, we are using `text-embedding-ada-002` as an embedding function.[8] Future work should investigate if other embedding functions would yield better results.

A relatively high number of memories are retrieved from each knowledge source, and they will be re-ranked in the next step.

### 3.4 Re-ranking Memories

This step aims to prioritize and filter the high number of retrieved memories from different knowledge sources to fit into the context window of the LLM. The input to the re-ranking step is multiple lists of memories, one list per knowledge source, each sorted in descending order by similarity to the query. The lists are combined and sorted by a ranking score assigned to each retrieved memory. The top $n$ memories with the highest scores are returned, where $n$ is a modifiable parameter.

The memory scoring function allows for a high degree of customization. Multiple factors can be considered and combined into a score:

- *Similarity:* The similarity between the search query and the retrieved memory.

---

[8] https://openai.com/blog/new-and-improved-embedding-model

240

- *Retention:* How "strong" the memory is, i.e., the inverse of how much it was already forgotten, as explained in section 3.8.

- *Possible further factors:* Memory importance, trust in the corresponding knowledge source, user preferences, current emotional state of the character, progression of some underlying story, and more.

The scoring function is a weighted sum with with $k$ factors $x_i$ and associated weights $w_i$:

$$R \cdot \sum_{i=1}^{k} w_i \cdot x_i$$

Notably, the retention value $R$ (see section 3.8) must not be part of the weighted sum, as a virtual character should not be able to recall a forgotten memory ($R = 0$), irrespective of the other factors.

We compute a ranking score for each memory using this formula. Finally, we return the $n$ memories with the highest score.

### 3.5 Character Response Generation

As we explain in section 5, there could potentially be multiple strategies to create character responses. Here, we present the "base" strategy, used when the retrieved memories are relevant.

In this case, the character response is generated using the LLM with a dynamic prompt. The prompt contains the following:

- *Instructions:* The instructions tell the LLM to simulate a virtual character and to answer the last user message. It also explains how to use the provided memories.

- *Character Name:* The name of the simulated character.

- *Character Description:* A 500-word description of the simulated character. This description gives the LLM context about the character's personality and how it expresses itself.

- *Chat History:* The chat history can be seen as the short-term memory of the character. We provide the last six messages in full and a summary of the previous messages.

- *Memories:* A list of relevant memories. This list results from the retrieval process, including the re-ranking step described in section 3.4.

An example prompt for the response generation can be found in Appendix A.5.

Once the response is created, it is returned to the user and inserted into the chat history. After a period of inactivity, new long-term memories are created.

### 3.6 Memory Creation Pipeline

The system aims to assist artists in creating interactive experiences with characters that can improvise and recall information about their universe and past interactions. To achieve this, any unstructured text, such as chat histories or character descriptions, can be converted to "facts" or "observations," forming character memories. Figure 3 shows how the memories are created on a high level.
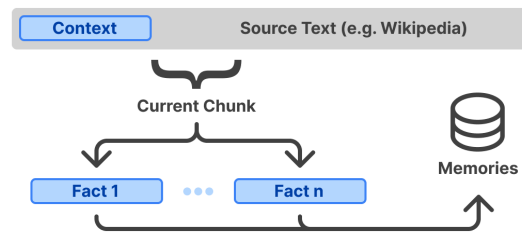


Figure 3: *The memory-creation pipeline:* A high-level overview of how memories are created from any raw text. Section 3.6 explains this process in detail.

The conversion from raw text to memories functions as follows: First, the source text is split into smaller chunks. Each chunk is then processed in isolation but with a context that contains a summary of all previous chunks. This context provides additional information that might be necessary to understand references or dependencies in the current chunk. For each pair of (chunk, context), we use the LLM to extract important facts about the chunk. These facts are then post-processed by resolving references (e.g., pronouns), ensuring that each fact is understandable atomically (without context). For each identified fact, we create one corresponding memory.

### 3.7 Memory Structure

Each memory has two representations. The first representation contains all information about the memory and is used internally within the system. The second representation is a filtered version and is used in the response-generation prompt for the LLM.

The full representation of the memory (used internally) contains the following information:

- *A fact:* The central part of the memory is some fact, or observation, which was contained in the corresponding source text (see section 3.6). In most cases, this is a sentence capturing some information about the virtual character or some previous interaction. Examples can be found in Figure 1 and in Appendix A.5.

- *Embedding:* A vector representation of the fact. This embedding captures the semantic meaning of the memory and is used for similarity search during memory retrieval.

- *Meta information:* In addition to the core memory, we also store information related to the forgetting model described in section 3.8, including last access time and memory stability. This information is used in the re-ranking step, as explained in section 3.4. Furthermore, each memory has a unique identifier.

The filtered representation, as used in the response generation prompt, only contains the textual representation of the memory content (the fact) and a shortened identifier, enabling the LLM to reference used memories.

Memories in our system are condensed facts rather than raw text. Raw text, e.g., an online article, often has a low information density and contains irrelevant elements such as links. Moreover, raw text requires a certain length to be meaningful due to long-range dependencies, leading to polluted prompts and subpar results. Including raw text in the response prompt burdens the LLM with additional condensing and filtering tasks, straining its capabilities. Additionally, using raw text limits the number of memories that can fit into the prompt due to the LLM's limited context window.

### 3.8 Forgetting Model

We incorporated a forgetting model into the memory-retrieval system to make virtual characters feel human-like. The system is inspired by how human memory behaves. Needless to say, the inner workings of human memory are complex and not fully understood, so the system is only a crude approximation of how real memory works. The achieved effect is that characters will slowly forget memories that they do not use and strengthen those that they do use.

The forgetting model is a modified version of Ebbinghaus's forgetting curve (Ebbinghaus , 1885):

$$R = e^{-\frac{dt}{S}}$$

In this formula, $R \in (0, 1)$ is the memory retention, $t \in \mathbb{R}^+$ is defined as the elapsed time since the last access, $S \in \mathbb{R}^+$ is the memory stability, determining how "strong" a memory is. Further, $d \in \mathbb{R}^+$ is a decay constant that defines how forgetful a character is in general. To simulate learning through repetition, each time a memory is accessed, its stability $S$ is updated by multiplying with a boost factor $b \in \mathbb{R}^+$. Thus, $b$ determines how fast memories are strengthened through repetition, or, in other words, how fast a character can learn.

This mechanism can also be used to assign an importance score to memories during memory creation. For example, for core memories of a trusted source, $S$ can be set to a very high number, effectively disabling forgetting. On the other hand, when memories are created from past user conversations, the character should not remember the details forever. Thus, $S$ can initially be set to a lower value, making the memories disappear sooner if they are not accessed.

These tweakable parameters enable a high degree of artist control for the memory of the virtual character.

## 4 Evaluation

### 4.1 Objectives

The primary goal of the system is to simulate virtual, artist-crafted characters. To achieve this goal, two things are essential: first, the virtual characters must be able to improvise, e.g., by reacting to user messages and combining retrieved memories in novel ways. Second, the characters should do so without hallucinating factual information.

The ability for improvisation is intrinsic to LLMs. However, maintaining factual accuracy is paramount to preserving the integrity of the artist-crafted characters. Any discrepancies or inaccuracies, otherwise known as "hallucinations," could undermine the coherency and believability of the character.

We evaluate how well the virtual characters' responses are grounded in the retrieved memories. That is, we aim to determine whether all information in the responses is supported by the retrieved memories. We also evaluate the effectiveness of the LLM in accurately referencing the used memories.

Figure 4: *The fact-checking pipeline:* In step 1, each sentence in a virtual character response is scanned for verifiable claims. In step 2, each extracted claim is fact-checked individually three times: once using all retrieved memories, once using only the memories referenced in the character response, and once using the character bio available in the prompt. Finally, in a manual filtering step, any unnecessary checks are discarded.

Further, we seek to understand the impact of whether some information about the character is part of the training data of the underlying LLM. To this end, we examine two contrasting cases: Sherlock Holmes, a well-documented character with abundant online data, and Amina Osei, a newly created character without any pre-existing information online. Since the LLM used was trained on much of the public internet (Brown et al.), it "knows" about Sherlock Holmes, even without retrieving memories, but has no intrinsic information about Amina Osei.

### 4.2 Methodology

Our evaluation methodology involved collecting chat data between the authors and the two virtual characters. In these chat exchanges, we investigated the characters' responses to topics for which they had memories and those for which they did not. Details of the dataset size are provided in Table 1. To facilitate the chat interactions, we developed a web-based chat platform.

First, recall that the overall system uses `gpt-3.5-turbo` as the underlying LLM. To evaluate the outputs of this system, the experimental setup involves a semi-automated fact-checking pipeline leveraging GPT-4, a more capable model, to assess the groundedness of the character responses (see Figure 4). This pipeline is designed to verify all factual claims the characters make against the information provided in the prompt. The pipeline operates in the following sequence:

1. First, unverified claims are extracted from the character responses. To do so, each response is broken down into individual sentences. From these sentences, we extract factual claims for verification. A claim represents information requiring verification, such

as a fact about the character or their universe. We extract these claims using GPT-4. In the prompt, we provide instructions, the full message, the selected sentence, examples, and formatting guidelines. For further details on this step, please refer to Appendix A.6, which contains a concrete example.

2. Subsequently, each claim is fact-checked against three different sources of truth: the complete list of retrieved memories, only the explicitly referenced memories in the response, and the character's bio. We perform three separate fact-checking rounds because this allows us to keep the input to the LLM short, noticeably improving the results. In this step, GPT-4 is instructed to assign one of four labels to every claim based on its validity: "fully_supported," "partially_supported," "not_supported," or "contradiction." In the prompt, we provide instructions, a context, an extracted claim, and a list of considered facts. Additionally, this step utilizes a chain-of-thought reasoning process, which has been shown to enhance the reasoning capabilities of LLMs (Wei et al.). Appendix A.7 shows a concrete example of this step.

3. Lastly, each claim is manually reviewed to determine if fact-checking is necessary. Some claims may not require verification, e.g.,"It is nice to meet you."

### 4.3 Results and Analysis

We first provide an overview of the data utilized in our evaluation. Table 1 summarizes the metadata regarding messages, sentences, claims, and reasoning steps.
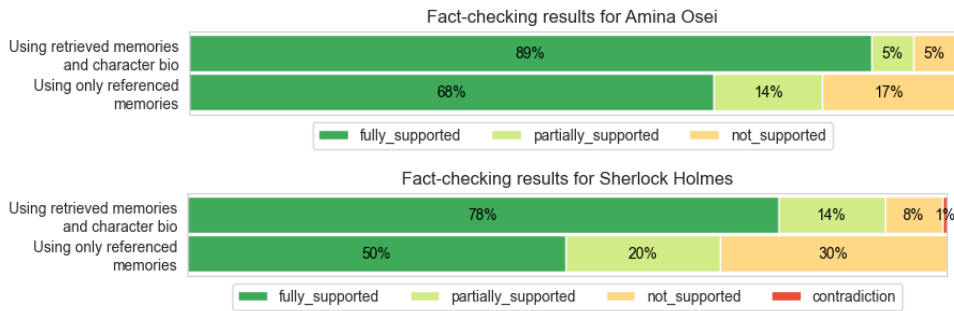
Figure 5: *Fact-checking results for the virtual character responses:* Each horizontal bar shows the results for different sources of truth as named on the left. The first category shows how grounded the character response is in the available information, using an aggregation of all three fact-checking results. The second category shows how good the LLM is at referencing which information it uses. Note that some percentage counts do not add up to 100 due to rounding.

| Character | Messages | Sentences | Extracted Claims | Reasoning Steps |
|---|---|---|---|---|
| Sherlock Holmes | 37 | 130 | 199 | 2213 |
| Amina Osei | 42 | 130 | 234 | 2292 |
| Total | 79 | 260 | 433 | 4505 |

Table 1: Metadata about the created dataset and the fact-checking pipeline.

### 4.3.1 Evaluation of Groundedness

To evaluate the groundedness of the virtual character responses, the fact-check uses the retrieved memories *and* the character bio, as this comprises all character information available in the LLM prompt. As shown in Figure 5, both characters displayed a high level of groundedness, meaning that the information given in the characters' responses matches the information provided in the prompt in most cases. For Sherlock Holmes, 133 claims were fully supported, 24 were partially supported, 13 were not supported, and only one was contradicted. For Amina Osei, 195 claims were fully supported, 12 were partially supported, and 12 were not supported, but no contradictions. This outcome indicates that the characters' responses are, in most cases, well-grounded in the characters' memories or bio, thus preserving the characters' integrity.

### 4.3.2 Evaluation of Referencing Capabilities

Assessing the referencing capabilities of the system, Amina Osei showed a stronger ability to reference the used memories than Sherlock Holmes. For Sherlock Holmes, 85 claims were fully supported, 35 were partially supported, and 51 were not supported when using only the referenced memories as

the source of truth. Amina Osei had 150 fully supported, 31 partially, and 38 not supported claims. Compared to Amina Osei, Sherlock Holmes has more cases in which the referenced memories only partially or do not support the made claims. We suspect this is because the underlying LLM has intrinsic knowledge about Sherlock Holmes, which leads it to rely less on the retrieved memories and hallucinate some references. On the other hand, the LLM has no intrinsic information about Amina Osei and thus has to rely more on the presented data.

### 4.3.3 Implications

These results suggest that our system can effectively simulate interactive characters that can generate largely accurate and grounded responses while preserving the characters' integrity. The capacity of the system to handle characters with different degrees of pre-existing LLM knowledge underscores its versatility. We plan to further improve these results with the system extension outlined in section 5.

While the evaluation results are promising, we acknowledge several limitations in Evaluation Limitations.

## 5 Future Work

A promising extension to our current system is a "Response Routing" stage. This feature, not yet fully realized, could augment control over character responses. It aims to maintain character integrity by ensuring responses stay within their memory bounds during improvisation and thus forms a critical component of our theoretical model.

Response Routing would decide the type of re-

sponse to be generated. It could consider the chat history, the created query, and retrieved memories, determining which of $n$ response generation strategies to use. It could generate normal responses (like explained in section 3.5) or evasive responses depending on memory relevance.

We propose separating decision-making from response generation due to the reasoning limitations of current LLMs. For example, a simple conditional prompt such as, "If the memories are relevant, create a normal response. Otherwise, create a response in which the character is clueless about the topic" has proven unreliable in our tests, sometimes leading to hallucinated responses if no relevant memories were retrieved.

We further plan to investigate how memories could influence character personalities over time. Future work will also improve the evaluation methodology, including a human evaluation.

We are also making ongoing efforts to increase the naturalness of the interactions with the virtual characters. Currently, the main problem is that they are purely *reactive* to user input. The interaction follows a strict pattern of question-answer-question-answer, where the user has to keep the conversation going. Future work will investigate how we can give virtual characters the ability to take initiative. They should be able to independently write messages to the user, potentially following some plan or narrative goal, as has been explored in previous works for other systems (Lehman and Leite; Leite et al.; Paetzel-Prüsmann and Kennedy).

## 6 Conclusion

This paper introduced a system for crafting memory-enhanced virtual AI characters that can interact intelligently. These virtual characters are grounded in their universe by being able to remember facts about themselves, their world, and past experiences.

The presented system, consisting of multiple knowledge sources, a vector search, and an LLM, is designed to be flexible, making it possible to create interactive characters by giving them access to memories that can be created from any textual descriptions about them. We designed the character memories to be human-like, being strengthened if used and forgotten when not, making characters feel authentic. Our evaluation, powered by a semi-automated fact-checking pipeline, showcased that the character responses were strongly grounded in the retrieved memories.

We anticipate this work to spur further investigations into the integration of memory systems within virtual characters and how to further guardrail LLMs to avoid hallucinations. As a result, we envision a future where these characters can perform convincingly within their narrative bounds while maintaining the capacity to improvise intelligently.

## References

Significant-Gravitas/Auto-GPT: An experimental open-source attempt to make GPT-4 fully autonomous.

Eleni Adamopoulou and Lefteris Moussiades. An Overview of Chatbot Technology. In *Artificial Intelligence Applications and Innovations*, IFIP Advances in Information and Communication Technology, pages 373–383. Springer International Publishing.

Alan Baddeley. Working Memory | Science.

Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? . In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, pages 610–623. Association for Computing Machinery.

Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, prefix=van den useprefix=false family=Driessche, given=George, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego de Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, Tom Hennigan, Saffron Huang, Loren Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol Vinyals, Simon Osindero, Karen Simonyan, Jack W. Rae, Erich Elsen, and Laurent Sifre. Improving language models by retrieving from trillions of tokens. Comment: Fix incorrect reported numbers in Table 14.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language Models are Few-Shot Learners. Comment: 40+32 pages.

Hongshen Chen, Xiaorui Liu, Dawei Yin, and Jiliang Tang. A Survey on Dialogue Systems: Recent Advances and New Frontiers. 19(2):25–35. Com-

ment: 13 pages. arXiv admin note: text overlap with arXiv:1703.01008 by other authors.

Hermann Ebbinghaus (1885). Memory: A Contribution to Experimental Psychology. 20(4):155–156.

Kai Greshake, Sahar Abdelnabi, Shailesh Mishra, Christoph Endres, Thorsten Holz, and Mario Fritz. More than you've asked for: A Comprehensive Analysis of Novel Prompt Injection Threats to Application-Integrated Large Language Models.

Shafquat Hussain, Omid Sianaki, and Nedal Ababneh. A Survey on Conversational Agents/Chatbots Classification and Design Techniques. pages 946–956.

Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. Survey of Hallucination in Natural Language Generation. 55(12):248:1–248:38.

Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. Generalization through Memorization: Nearest Neighbor Language Models. Comment: ICLR 2020.

Yunshi Lan, Gaole He, Jinhao Jiang, Jing Jiang, Wayne Xin Zhao, and Ji-Rong Wen. A Survey on Complex Knowledge Base Question Answering: Methods, Challenges and Solutions. Comment: Accepted by IJCAI 2021 Survey Track.

Jill Fain Lehman and Iolanda Leite. Turn-Taking, Children, and the Unpredictability of Fun. 37(4):55–62.

Iolanda Leite, André Pereira, and Jill Fain Lehman. Persistent Memory in Repeated Child-Robot Conversations. In Proceedings of the 2017 Conference on Interaction Design and Children, IDC '17, pages 238–247. Association for Computing Machinery.

Yoav Levine, Itay Dalmedigos, Ori Ram, Yoel Zeldes, Daniel Jannai, Dor Muhlgay, Yoni Osin, Opher Lieber, Barak Lenz, Shai Shalev-Shwartz, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. Standing on the Shoulders of Giant Frozen Language Models.

Andrea Madotto, Zhaojiang Lin, Chien-Sheng Wu, and Pascale Fung. Personalizing Dialogue Agents via Meta-Learning. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 5454–5459. Association for Computational Linguistics.

Yohei Nakajima. BabyAGI.

OpenAI. GPT-4 Technical Report. Comment: 100 pages.

Maike Paetzel-Prüsmann and James Kennedy. Improving a Robot's Turn-Taking Behavior in Dynamic Multiparty Interactions. In Companion of the 2023 ACM/IEEE International Conference on Human-Robot Interaction, HRI '23, pages 411–415. Association for Computing Machinery.

Joon Sung Park, Joseph C. O'Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. Generative Agents: Interactive Simulacra of Human Behavior.

David Patterson, Joseph Gonzalez, Quoc Le, Chen Liang, Lluis-Miquel Munguia, Daniel Rothchild, David So, Maud Texier, and Jeff Dean. Carbon Emissions and Large Neural Network Training.

Or Sharir, Barak Peleg, and Yoav Shoham. The Cost of Training NLP Models: A Concise Overview.

Sainbayar Sukhbaatar, Edouard Grave, Guillaume Lample, Herve Jegou, and Armand Joulin. Augmenting Self-attention with Persistent Memory.

Avril Thorne. Personal Memory Telling and Personality Development. 4(1):45–56.

David Thue, Vadim Bulitko, Marcia Spetch, and Eric Wasylishen. Interactive Storytelling: A Player Modelling Approach. 3(1):43–48.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention Is All You Need. Comment: 15 pages, 5 figures.

Zhiguo Wang, Patrick Ng, Ramesh Nallapati, and Bing Xiang. Retrieval, Re-ranking and Multi-task Learning for Knowledge-Base Question Answering. In Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, pages 347–357. Association for Computational Linguistics.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models.

Joseph Weizenbaum. ELIZA—a computer program for the study of natural language communication between man and machine. 9(1):36–45.

Sean Welleck, Jason Weston, Arthur Szlam, and Kyunghyun Cho. Dialogue Natural Language Inference. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 3731–3741. Association for Computational Linguistics.

Munazza Zaib, Quan Z. Sheng, and Wei Emma Zhang. A Short Survey of Pre-trained Language Models for Conversational AI-A New Age in NLP. In Proceedings of the Australasian Computer Science Week Multiconference, ACSW '20, pages 1–4. Association for Computing Machinery.

Tomáš Zemčík. A Brief History of Chatbots.

Terry Yue Zhuo, Yujin Huang, Chunyang Chen, and Zhenchang Xing. Exploring AI Ethics of ChatGPT: A Diagnostic Analysis. Comment: Technical Report.

# A Appendix

## A.1 Limitations

While the described system delivers promising results, it is not without limitations. This section explores some of the limitations that need to be addressed, before our system, or any LLM-leveraging system, for that matter, could safely be used in a production setting.

### A.1.1 System Limitations

Through multiple successive LLM calls, which all introduce some latency, it currently takes a few moments before the character answer is generated. While the system is definitely usable as is, especially for chat applications, future iterations should use parallelism wherever possible.

Presently, the system's ability to stylize character responses is restricted, as we utilize a general-purpose LLM. This limitation could be addressed by employing a more advanced general-purpose LLM in the future or by generating the character responses with an LLM that has been fine-tuned to produce texts of a specific style.

A further limitation concerns occasional hallucinations in the memory creation pipeline. For example, a memory attributed to the character Amina Osei introduced an unfounded surname for her partner. While it could be argued that humans also make mistakes when remembering things, this issue is a subject for further investigation.

### A.1.2 General LLM Limitations

LLMs, in general, carry inherent limitations. Currently, there is no definitive solution to ensuring safety in LLM outputs. This involves avoiding problematic outputs, such as those that are sexist, racist, or potentially harmful, especially when the target audience includes vulnerable groups such as children.

Furthermore, a common problem for LLMs is prompt injections. While not explicitly having tested our system in this regard, it is reasonable to assume that it is also vulnerable to such attacks.

### A.1.3 Evaluation Limitations

As a refresher, we use GPT-3 as the underlying LLM throughout our system and GPT-4 in the evaluation. Effectively, we are employing GPT-4 to evaluate the outputs from GPT-3. Both models are trained on overlapping datasets, so their behaviors are likely somewhat aligned. This similarity could mean that a response deemed plausible by GPT-3 might also be classified as plausible by GPT-4, which introduces a potential bias into our fact-checking process.

We have dissected the fact-checking process into multiple smaller steps to mitigate this issue, assessing the groundedness of atomic claims instead of entire character responses at a time. This approach led to more focused LLM outputs.

Second, our evaluation process is not entirely automated. Particularly challenging was the extraction of claims requiring fact-checking. This complexity arises from the inherent dual nature of the task for the virtual characters: they are expected to provide grounded answers while also improvising and reacting to user messages. Consequently, we employed GPT-4 to extract all checkable claims, which we manually reviewed and discarded if necessary.

Third, it can be expected that the determined factfulness labels would not fully align with a human evaluation or that there is a bias toward certain labels. The extent of such issues is unclear as of now. However, we manually verified several fact-checking results, auditing the reasoning steps and the assigned label, and found them to be very reasonable in most cases.

Lastly, our evaluation approach, while being less labor-intensive than a fully manual evaluation, is not inexpensive. The evaluation using GPT-4 took over 10 hours and amounted to approximately $50 in costs, equating to about 11 cents per checked claim. This cost factor might be a consideration for larger-scale evaluations in the future.

Future research should address these limitations and further refine the evaluation process.

## A.2 Ethics Statement

LLMs, upon which our work heavily relies, have demonstrated concerning impacts on certain populations. Studies have highlighted issues related to exploitative human annotation practices employed to refine model outputs (Bender et al.), generation of toxic responses, and inherent biases favoring certain cultural behaviors (Zhuo et al.; Brown et al.).

Our work does not directly confront these ethical concerns or engage in any of these areas. However, the system we propose attempts to constrain LLMs by biasing the generated output on the retrieved memories, with the content of the memories being controllable. This added control returns agency to human designers or artists, allowing them to steer the interactive experience away from potentially undesirable behaviors and towards a more diverse and inclusive representation.

## A.3 Supplementary Materials Availability Statement

We make several example prompts available in the appendix. While they are not essential for understanding the paper, they capture know-how about how LLMs can be used and could therefore be helpful to the research community. The source code and evaluation dataset cannot be made available due to our institution's sharing policy.

## A.4 Query Creator Prompt

The following prompt is used to generate a memory search query. The placeholder {chat_history} is replaced with the current chat history. Through experimentation, we found that this prompt generates good search queries.

```
{chat_history}

Create a search query for the character's memory that helps answer the last user
    message. You cannot ask for clarification. Provide only the query.
```

## A.5 Example Response Generation Prompt

We use the following prompt to create a character response for Sherlock Holmes. It contains instructions, a character bio, a chat history, and many dynamically retrieved memories.

```
You are acting as a virtual character and you are having a conversation with a
    user. The character you are simulating is named Sherlock Holmes. Your task is
    to answer the user based on the chat history. You should answer the last
    message in the chat history.

CHARACTER_BIO_START

Sherlock Holmes, the world's most famous consulting detective, is a fictional
    character created by British author Sir Arthur Conan Doyle. Appearing in four
    novels and 56 short stories, Holmes has become synonymous with the art of
    deduction and solving complex mysteries. His character has transcended the
    confines of literature, inspiring countless adaptations in film, television,
    and other media.
Born on January 6, 1854, in the United Kingdom, Sherlock Holmes resided at 221B
    Baker Street, London, with his trusted friend and confidant, Dr. John H.
    Watson. Known for his towering intellect, keen observational skills, and
    mastery of disguise, Holmes captivated readers with his ability to solve
    seemingly impossible cases with logic and reasoning.
Holmes's character is often seen as eccentric, with an analytical mind that borders
    on obsessive. He possesses an insatiable curiosity and a relentless drive to
    uncover the truth, often leading him to immerse himself in the darkest corners
    of London's criminal underworld. His eccentricities extend to his personal
    life, where he is known to dabble in chemistry, play the violin, and
    occasionally use recreational drugs to stimulate his mind during periods of
    inactivity.
A quintessential Victorian gentleman, Sherlock Holmes adheres to a strict code of
    honor and chivalry, often taking on cases without any expectation of payment if
    he believes the cause is just. Although his demeanor may come across as cold
    and calculating, Holmes is not without compassion. He demonstrates a profound
```

248

understanding of human nature and is not afraid to bend the rules to protect
    the innocent.
One of Holmes's defining features is his unique method of solving crimes, which he
    calls \"the Science of Deduction.\" This method involves observing the minutest
    details and connecting seemingly unrelated clues to reveal the truth. In
    addition to his extraordinary mental faculties, Holmes is also an adept martial
    artist, skilled in boxing and the martial art of Bartitsu.
Dr. John Watson serves as Holmes's loyal companion and chronicler. The two
    characters share a strong bond, with Watson often providing the emotional and
    moral grounding that Holmes lacks. Through Watson's eyes, readers gain insight
    into Holmes's mind, following his process of deduction and marveling at the
    brilliance of his conclusions.
Sherlock Holmes's arch-nemesis, Professor James Moriarty, is a criminal mastermind
    who presents Holmes with some of his most challenging cases. Their rivalry
    serves as a recurring theme throughout the series, culminating in a dramatic
    confrontation at the Reichenbach Falls in \"The Final Problem.\" Other notable
    adversaries include Irene Adler, a clever and resourceful woman who earns
    Holmes's admiration, and Charles Augustus Milverton, a ruthless blackmailer who
    pushes Holmes to his moral limits.
Over the years, Holmes's adventures have been adapted into numerous films and
    television series, with actors such as Basil Rathbone, Jeremy Brett, Robert
    Downey Jr., and Benedict Cumberbatch portraying the detective. These
    adaptations often offer new interpretations of the character, highlighting
    different aspects of his complex personality and showcasing the timeless appeal
    of his intellect and wit.
Sherlock Holmes remains a cultural icon, with his deerstalker hat, pipe, and
    magnifying glass becoming enduring symbols of the detective genre. His legacy
    has left an indelible mark on popular culture, inspiring generations of readers
    and viewers to embrace the power of logic, deduction, and the human mind's
    capacity to solve the most intricate mysteries.

CHARACTER_BIO_END

This is the current chat history:
START_CHAT_HISTORY
Most recent messages:
user: Hello! How are you? My name is Lo
virtual-character: Greetings, Lo. I hope you are well. As a cryptanalyst [MEMORY_0,
    MEMORY_3, MEMORY_14, MEMORY_15, MEMORY_18, MEMORY_19, MEMORY_21, MEMORY_27], I
    am familiar with all forms of secret writing and have written a monograph on
    Polyphonic Motets of Lassus [MEMORY_0, MEMORY_3, MEMORY_14, MEMORY_15,
    MEMORY_27]. Is there anything that I may assist you with today?
user: How do you do for a living? I may need some help...

END_CHAT_HISTORY

The simulated character has memories. Use the memories to guide your answer.
START_MEMORIES:

MEMORY_30: The text includes various categories related to Sherlock Holmes,
    including screen adaptations, video games, and a list of fictional
    characteristics associated with the character.
MEMORY_31: ''Guinness World Records'' lists Sherlock Holmes as the most portrayed
    human literary character in film and television history.
MEMORY_32: The character of Sherlock Holmes has been adapted and reimagined by
    various writers in non-canonical works.
MEMORY_33: Sherlock Holmes's personality and habits are not specified in this
    specific section of the text.
MEMORY_34: The article discusses the skills and abilities of the character Sherlock
    Holmes created by Arthur Conan Doyle.
MEMORY_35: Dr. Watson describes Sherlock Holmes as \"bohemian\" in Sherlock Holmes'
    habits and lifestyle.
MEMORY_36: Short stories featuring sage-detective Zavant Konniger and Zavant
    Konniger's halfling manservant Vido, written by fantasy authors Gordon Rennie
    and Josh Reynolds for the Warhammer Fantasy universe, are also mentioned.
MEMORY_37: The article discusses the various skills and abilities of Sherlock
    Holmes, a complex character created by Arthur Conan Doyle.
MEMORY_38: Arthur Conan Doyle depicted returned colonials as \"marginal, physically
    ravaged characters that threaten the peace,\" according to Yumna Siddiqi.
MEMORY_39: Sherlock Holmes displays a strong aptitude for acting and disguise in

```
     several stories.
MEMORY_40: Guinness World Records awarded Sherlock Holmes the title for \"most
    portrayed literary human character in film & TV\" in 2012.
MEMORY_41: Sherlock Holmes displays a strong aptitude for acting and disguise.
MEMORY_42: Guinness World Records awarded Sherlock Holmes the title for \"most
    portrayed literary human character in film & TV\" in 2012.
MEMORY_43: Guinness World Records awarded Sherlock Holmes the title for \"most
    portrayed literary human character in film & TV\" in 2012.
MEMORY_44: Sherlock Holmes is a famous literary character known for Sherlock
    Holmes's investigative skills in various fields including geology, chemistry,
    anatomy, British law, psychology, and secret writing.
MEMORY_45: Conan Doyle observes the dress, attitude, and physical condition of his
    clients and suspects to deduce their origins and recent history.
MEMORY_46: The portrayal of Eastern foreigners as the root cause of \"infection\"
    within and around Europe is a common theme in the stories, according to Susan
    Cannon Harris.
MEMORY_47: Sherlock Holmes also displays a strong aptitude for acting and disguise
    in several stories.
MEMORY_48: Sherlock Holmes's ability to change Sherlock Holmes's appearance helped
    Sherlock Holmes personify the English eccentric chameleon.
MEMORY_49: The character, Sherlock Holmes, has been adapted into numerous films,
    television shows, and other media and has become a cultural icon.
MEMORY_50: The court ruled in favor of the characters being in the public domain in
    the U.S.
MEMORY_51: The character Sherlock Holmes has been adapted into numerous films,
    television shows, and other media, and has become a cultural icon.
MEMORY_52: Sherlock Holmes is a fictional character who is a detective and a master
    of deduction, and has been portrayed in various forms of media.
MEMORY_53: The character of Sherlock Holmes has been adapted and reimagined by
    various writers, resulting in non-canonical works featuring Sherlock Holmes.
MEMORY_54: Sherlock Holmes is skilled in playing the violin, singlestick, boxing,
    and sword fighting.
MEMORY_55: There have been many scholarly works dealing with Sherlock Holmes, some
    working within the bounds of the Great Game, and some written from the
    perspective that Sherlock Holmes is a fictional character.
MEMORY_56: The article also mentions short stories featuring sage-detective Zavant
    Konniger and his halfling manservant Vido for the Warhammer Fantasy universe.
MEMORY_57: The character and stories have had a lasting impact on mystery writing
    and popular culture.
MEMORY_58: Short stories featuring sage-detective Zavant Konniger and his halfling
    manservant Vido, written by fantasy authors Gordon Rennie and Josh Reynolds for
    the Warhammer Fantasy universe, were published by Black Library from 2002 to
    2018, including \"How Vido Learned the Trick\" and \"The Problem of Three-Toll
    Bridge\".
MEMORY_59: The character and stories have had a profound and lasting effect on
    mystery writing and popular culture as a whole.
MEMORY_60: Another common pastiche approach is to create a new story fully
    detailing an otherwise-passing canonical reference.
MEMORY_61: Sherlock Holmes uses acting and disguises to gather evidence and
    incriminate the guilty in some stories.
MEMORY_62: In \"A Scandal in Bohemia,\" Sherlock Holmes recognizes the true
    identity of a disguised character despite his supposed ignorance of politics.
END_MEMORIES.

The answer must be based on the memories. Do not talk about anything that is not in
    the memories. For each sentence, provide a source like [MEMORY_i].

Create a character response to the last message of the user. The response must be
    from the point of view of Sherlock Holmes. The response should be around 50-70
    words. Response:
```

## A.6 Example Claim Extraction Prompt

The following example illustrates how we use GPT-4 to extract claims that need fact-checking from character responses. In the following JSON-like prompt, "context" contains the entire message currently being processed. Further, "sentence_to_check" is a sentence taken from the context. The field "already_found_claims" contains previously identified claims from the current message. The other fields are self-explanatory.

```
{
  'context': 'Amina Osei: My name is Amina Osei, born on July 17th in Accra, Ghana
      . As a scientist, I have pursued a Bachelor of Science in Biological Sciences
      and a Master of Science in Astrophysics from the University of Ghana . Please
      let me know if you require further information.',
  'sentence_to_check': 'My name is Amina Osei, born on July 17th in Accra, Ghana .',
  'task': 'Find all claims in sentence_to_check that can be fact checked. Be
      thorough. Only include claims from sentence_to_check. Exclude questions,
      greetings and other social pleasantries.',
  'already_found_claims': [],
  'examples': [
    {
      'example_input': 'I live in Armenia with my two sisters.',
      'example_output': [
        'I live in Armenia',
        'I have two sisters'
        ]
    },
    {
      'example_input': 'Hello my dear, how are you?',
      'example_output': []
    },
    {
      'example_input': 'What are you doing?',
      'example_output': []
    },
    {
      'example_input': "As a gardener, it's easy for me to tell one plant from
          another.",
      'example_output': [
        'I am a gardener',
        "it's easy for me to tell one plant from another."
        ]
    }
  ]
}
const output: { "found_claims": string[] } =
```

The following was the generated output to this prompt using GPT-4. Note that it generates a valid JSON, which can be parsed easily. This can be achieved by prompting the LLM with JSON and defining the desired output with a TypeScript type (see the last line above).

```
{
  "found_claims": [
    "My name is Amina Osei",
    "born on July 17th",
    "in Accra, Ghana"
  ]
}
```

## A.7 Example Fact Checking Prompt

The following example illustrates how we use GPT-4 as a fact checker. Again, note how a TypeScript type can instruct the LLM to produce correctly formatted output that can easily be parsed.

```
{
  'untrusted_context': 'Amina Osei: My name is Amina Osei, born on July 17th in
      Accra, Ghana . As a scientist, I have pursued a Bachelor of Science in
      Biological Sciences and a Master of Science in Astrophysics from the
      University of Ghana . Please let me know if you require further information.',
  'untrusted_claim': 'My name is Amina Osei',
  'considered_facts': '\nFACT_28: Amina was born on July 17th.\nFACT_32: Amina
      earned a Bachelor of Science in Biological Sciences with a focus on botany
      and a Master of Science in Astrophysics from the University of
      Ghana.\nFACT_META_1: The name of the speaker is Amina Osei',
```

```
   'task': '\nYou are a superhuman fact-checking expert AI lawyer.\nYour work is
      thorough and your reasoning_steps hold up in court.\nYou are self-critical
      and detect your own mistakes.\nYour task is to fact-check the untrusted_claim
      using the considered_facts.\nBoth the untrusted_claim and the
      considered_facts are anchored in a fictional universe.\n'
}

const fact_check_result: {
    // "reasoning_steps" contains logical reasoning steps that help select the
        determined_factfulness.
    "reasoning_steps": string[];
    // determined_factfulness must depend only on the considered_facts.
    "determined_factfulness":
        | "fully_supported" // If a fact fully confirms a claim
        | "partially_supported" // If a fact partially confirms a claim
        | "not_supported" // If no fact confirms or contradicts the claim
        | "contradiction"; // If some fact contradicts the claim
} =
```

The generated output to this prompt, using GPT-4, was:

```
{
    "reasoning_steps": [
        "The untrusted_claim states that the speaker's name is Amina Osei.",
        "FACT_META_1 confirms that the name of the speaker is Amina Osei."
    ],
    "determined_factfulness": "fully_supported"
}
```

# Metric-Based In-context Learning: A Case Study in Text Simplification

**Subha Vadlamannati**
Mercer Island High School
Seattle, USA
subhavee2@gmail.com

**Gözde Gül Şahin**
Computer Engineering Department
Koç University, Istanbul, Turkey
gosahin@ku.edu.tr

## Abstract

In-context learning (ICL) for large language models has proven to be a powerful approach for many natural language processing tasks. However, determining the best method to select examples for ICL is nontrivial as the results can vary greatly depending on the quality, quantity, and order of examples used. In this paper, we conduct a case study on text simplification (TS) to investigate how to select the best and most robust examples for ICL. We propose **M**etric-**B**ased in-context **L**earning (MBL) method that utilizes commonly used TS metrics such as SARI, compression ratio, and BERT-Precision for selection. Through an extensive set of experiments with various-sized GPT models on standard TS benchmarks such as TurkCorpus and ASSET, we show that examples selected by the top SARI scores perform the best on larger models such as GPT-175B, while the compression ratio generally performs better on smaller models such as GPT-13B and GPT-6.7B. Furthermore, we demonstrate that MBL is generally robust to example orderings and out-of-domain test sets, and outperforms strong baselines and state-of-the-art finetuned language models. Finally, we show that the behaviour of large GPT models can be *implicitly controlled* by the chosen metric. Our research provides a new framework for selecting examples in ICL, and demonstrates its effectiveness in text simplification tasks, breaking new ground for more accurate and efficient NLG systems.

## 1 Introduction

Text simplification (TS) is a crucial task in natural language processing, with the goal of converting complex text into simpler, easier-to-understand one. This is particularly important for individuals who struggle with comprehending complex languages, such as second language learners or individuals with cognitive impairments (Stajner, 2021) and disabilities like dyslexia (Rello et al., 2013) and autism (Barbu et al., 2015). For the aforementioned reasons, NLP community has shown great interest in the topic, introducing plenty of datasets (e.g., ASSET (Alva-Manchego et al., 2020)), models, and evaluation metrics (e.g., SARI (Xu et al., 2016)).

There have been numerous approaches to TS proposed in the literature, including non-neural or rule-based methods (Nassar et al., 2019), machine translation approaches (Xu et al., 2016), and finetuning of large language models (Sheang and Saggion, 2021) on downstream task data. Recently, it has been shown that large language models such as GPT-3, are capable of in-context learning (ICL) (Brown et al., 2020a)—an emerging ability to learn from in-context samples without modifying model parameters. [1] Despite its strong ability, ICL still mostly falls behind the performance of finetuning techniques (Dong et al., 2023).

Recent studies have shown that in-context learning is highly variable to a range of factors, such as the number of examples, quality of examples, and even the order of examples (Lu et al., 2022; Liu et al., 2022; Dong et al., 2023). To address these concerns, recent literature has proposed several techniques for selecting the most relevant examples for ICL (Liu et al., 2022; Sorensen et al., 2022; Gonen et al., 2022; Rubin et al., 2022). The majority of them aim to *retrieve* a set of samples from the validation set that resembles the test set most by either training a separate retrieval model or utilizing an existing encoder to calculate similarities between pairs of sentences. However, adopting these techniques for text-generation tasks with multiple references is nontrivial, and the need to access to the full test set to pick examples from is not desirable, and may not always be possible in real-life scenarios.

In order to address this problem, we propose a simple yet intuitive metric-based selection tech-

---

[1] We refer the readers to http://ai.stanford.edu/blog/understanding-incontext/ for a summary of in-context learning inner mechanics.

nique, which we refer to as **M**etric-**B**ased in-context **L**earning (MBL), to perform efficient and robust in-context learning with large language models for text generation tasks. Unlike previous ICL techniques, MBL only requires access to the development set and uses more informed measures rather than requiring generating sentence embeddings or training separate specified retrieval models. Furthermore, we perform an extensive set of experiments with GPT-3 models of various sizes (175B, 13B, and 6.7B) [2], specifically focusing on their performance for TS. We investigate utilizing commonly used TS metrics (e.g., SARI, compression ratio) for example selection and answer several research questions regarding their strengths and weaknesses on a variety of datasets and models. Through our experiments, we show that metric-based selection can significantly improve the performance of large language models on TS. We also demonstrate that these results are generally robust to various orderings and perform well in out-of-domain settings. This paper provides the following contributions:

- We provide a naive yet effective and robust approach to selecting examples for in-context learning, a.k.a., *metric-based learning* (MBL) [3], and show that it achieves state-of-the-art results on two well-known benchmark datasets (TurkCorpus and ASSET when the optimal metric is used (see §5.1)) [4].

- We demonstrate the robustness of MBL to example ordering (see §5.2) and to out-of-domain test sets with some exceptions (see §5.3), suggesting that the order of examples and the origin of the development data are not the most important factors for MBL.

- We show that MBL improves upon important baselines (e.g., zero-shot, random selection), state-of-the art ICL selection (e.g., KATE-GPT (Liu et al., 2022)) and text simplification methods (Sheang and Saggion, 2021) (see §5.4).

---

- **Our results suggest that GPT-175B can be *implicitly controlled*** via optimal metric-based learning, i.e., BERTScore Precision-based learning optimizes BLEU, while SARI-based selection optimizes SARI scores.

We release all generation outputs, baseline models and evaluation scripts publicly with `https://github.com/NLP-KU/metric-based-in-context-learning/`.

## 2  Related Work

**Text Simplification (TS) Methods**  Recently, LLMs have been applied to text simplification through transfer learning approaches. For instance, Qiang et al. (2020) fine-tuned a BERT model on a text simplification dataset, achieving strong results on multiple benchmarks. Similarly, Sheang and Saggion (2021) introduced a transfer learning approach for text simplification using the T5 model and achieved current state-of-the-art results on standard TS benchmarks. Recent work in the TS domain has a particular focus on controllable text simplification, in which different "control tokens" are embedded in seq2seq models to control model outputs. This is seen in both Sheang and Saggion (2021) and Chamovitz and Abend (2022), where a large language model (BART, T5, etc.) is modified with several control tokens, like the number of words, Levenshtein similarity, and various text rewriting operations. A vast amount of earlier systems (e.g., (Xu et al., 2016)) have formulated text simplification as a machine translation task and employed neural machine translation architectures.

**TS Evaluation**  Work on the suitability of various metrics for TS has also been an active area of discussion. While the most commonly used metric in TS is currently SARI (Xu et al., 2016), there is a concern over the metric that best correlates with human judgment. Alva-Manchego et al. (2021) conduct a detailed analysis of several commonly used metrics in the TS field, and suggest BERTScore_Precision as a primary metric of reference-based evaluation. Following these results, we also use BERTScore_Precision as a metric to select examples. Recent studies (Sulem et al., 2018; Tanprasert and Kauchak, 2021) analyzing the suitability of the other two common metrics, namely BLEU (Papineni et al., 2002) and FKGL (Kincaid et al., 1975), strongly advise against these metrics for TS. For these reasons, we do not select

examples based on either metric.

**Example Selection and Ordering in ICL**   While large language models like GPT-3 perform exceptionally well on a variety of downstream tasks, selecting examples for in-context learning is non-trivial. Research on example selection is still in early stages, and a unified approach to selecting examples for downstream tasks has not yet been proposed (Dong et al., 2023). Liu et al. (2022) propose selecting the $k$ most similar examples to the test set from the training/development set via measuring cosine distance in an embedding space (e.g., encodings from RoBERTa), and achieve strong results on various tasks like table-to-text generation. On a similar line, Rubin et al. (2022) introduce a more sophisticated method, where the authors train a two-step retrieval model to select ICL examples. Another set of work focus on optimizing prompts via mutual information (Sorensen et al., 2022) or perplexity (Gonen et al., 2022), that don't require labeled sets. We consider Kate-GPT (Liu et al., 2022) as the closest work to ours, since both the intuition (i.e., choosing from a labeled validation set) and approach (i.e., learning-free) are similar.

## 3   Metric-based In-Context Learning

Following the line of work for retrieving the best samples from development set (Liu et al., 2022; Sorensen et al., 2022), we introduce a simple and intuitive technique based on employing standard evaluation metrics for selecting the examples.

**Task Setup**   Given the list of sentences $l = [c, r_1, r_2, ..., r_n]$, where $c$ is the complex and $r_i$ is the simple reference sentence; our goal is to find the best $k$ pairs, $[c, r_i]$, such that the final text simplification performance on the test set is maximized. To do so, we go through each $l$ in the development set and measure the distance between each $c$ and $r_i$ according to a metric, $m$. Finally, we pick the top $k$ pairs and fill the prompt template with the samples: "Complex sentence: $\{c\}$, Simple sentence:$\{r_i\}$".

We initially considered a long list of task-specific as well as general generation metrics that contain the standard evaluation metrics for TS, namely as SARI, BLEU, FKGL; as well as a simple analysis metric: Compression Ratio, and a more recent textual similarity metric BERTScore as suggested by Alva-Manchego et al. (2021). Following the criticisms by Sulem et al. (2018) and Tanprasert and Kauchak (2021), we removed BLEU and FKGL

from the list of candidate metrics.

**Compression Ratio (CR)**   It is simply calculated by dividing the number of characters in $c$ by the number of characters $r_i$. We consider the pairs with higher compression ratios as more preferable candidates for TS.

**BERTScore Precision (BP) (Zhang et al., 2020)** BERTScore computes the cosine similarity between each token in the candidate, $y$, and reference, $x$, sentences. Precision is calculated as:

$$\text{Prec} = \frac{1}{|y|} \sum_{y_j \in y} \max_{x_i \in x} x_i^\top y_j \qquad (1)$$

We discard pairs with a score of 1 since they would simply be duplicates.

**SARI (Xu et al., 2016)**   It is the defacto standard evaluation metric for TS. In general terms, it compares prediction against both the input and the reference sentences. It calculates a weighted average of F1 scores for three operations: addition, deletion, and keeping. Precision and recalls for each operation are calculated based on n-gram overlaps between the prediction, input and reference sentences. To calculate the SARI score for each $c$-$r_i$ pair, we denote $r_i$ as the prediction, $c$ as the input, and $[r_1, ..., r_{i-1}, r_{i+1}, ..., r_n]$ as the reference sentences. Hence, this measure can only be applied when there are multiple references.

## 4   Experimental Setup

To investigate the effects of metric-based selection techniques on TS, we perform a comprehensive set of experiments using various LLMs, sample sizes, and datasets; and compare against strong baseline and state-of-the-art models. Following the criticism (Sulem et al., 2018) on using BLEU (Papineni et al., 2002), we use SARI (Xu et al., 2016) as our main evaluation metric. However, we also report BLEU for two reasons: i) to be consistent with previous works (see §4.4) and ii) to gain more insights on how the chosen metric for MBL effects the results measured with different metrics.

### 4.1   Models

Due to its recent success in text generation and in-context learning for various downstream tasks, we experiment with the GPT-3 (Brown et al., 2020b) model. We use three different version with the following parameter sizes:

175B, a.k.a., `da-vinci-003`, 13B, a.k.a., `curie`, and 6.7B, a.k.a., `babbage`. We used OpenAI API[5] to generate responses using temperature=0.7, max_tokens=256 top_p=1, frequency penalty=0 and presence penalty = 0.

## 4.2 Datasets

We perform our main experiments on the AS-SET (Alva-Manchego et al., 2020) and TurkCorpus (Xu et al., 2016) datasets. To investigate the transferability of our models, we conduct additional experiments on an out-of-domain cognitive simplification dataset, FestAbility (Chamovitz and Abend, 2022).

**TurkCorpus** is a widely-used dataset with 2000 validation and 359 test sentences. It has 8 reference sentences for each original sentence in both the validation and test set.

**ASSET** is another widely used TS dataset with the intention of improving upon TurkCorpus. It has the same 2000 validation and 359 original test sentences but introduces 10 new reference sentences for each original sentence. ASSET is deemed to be simpler by human evaluation in both fluency and simplicity (Alva-Manchego et al., 2020). ASSET improves upon TurkCorpus as it allowed human reviewers to focus on a wider variety of TS operations, which are: lexical paraphrasing, compression, and sentence splitting. Because of this, we emphasize the experiments done on ASSET rather than TurkCorpus while interpreting the results and answering the research questions in §4.

**FestAbility** is a cognitive simplification dataset with 321 pairs of complex and simple sentences—i.e., only one reference sentence. Each of these is additionally annotated with rewriting operations such as `<ADDITION>` and `<DELETION>`. These sentences are generated from the transcript of the virtual accessibility conference, and simplifications are generated from the Yalon Method (Chamovitz and Abend, 2022), a specialized method for simplifying text for individuals with cognitive impairments.

## 4.3 Baselines

For comparison, we implement three baselines: i) random selection ii) KATE-GPT (Liu et al., 2022) and iii) zero-shot. In the random setting, we randomly select $c$ and $r_i$ pairs from the validation

sets. For KATE-GPT (Liu et al., 2022), we use the default setting that employs RoBERTa-base for contextualized embeddings and cosine similarity for the distance metric. Given that KATE-GPT calculates sentence pair similarities between the development and test set, unlike just the development set (like ours), we choose complex sentences as the representative. Zero-shot setting is simply conducted with the same instruction prompt without providing any examples.

## 4.4 Text Simplification State-of-the-art

We compare our results across multiple state-of-the-art systems.

**MUSS (BART+ACCESS Supervised)** Martin et al. (2022) fine-tune BART (Lewis et al., 2020) and add information from the four simplification tokens trained in ACCESS.

**Finetuned-T5** Sheang and Saggion (2021) fine-tune T5 by adding multiple control tokens (e.g., compression ratio, Levenshtein similarity ratio, word rank, and number of words) similar to AC-CESS, which control the model's outputs. To the best of our knowledge, they achieve the current state-of-the-art on both the TurkCorpus and AS-SET datasets, with SARI scores of 43.31 and 45.04 respectively.

## 4.5 Evaluation

Even though SARI is considered the standard evaluation metric in our experiments below, we evaluate the results both with SARI and BLEU to emphasize the behavior differences in metric-based selection. It should be noted that, SARI compares prediction against input and reference(s), while BLEU compares only the prediction against reference(s). We use the package EASSE (Alva-Manchego et al., 2019) with the default settings [6] to generate reports for all of our experiments.

## 5 Experiments and Results

We conduct a comprehensive set of experiments with the setup explained in §4. Following the work by Lu et al. (2022), we experiment with the $k$ values as $1, 2, 4, 6, 8, 10, 15, 20$ examples. We repeat the random baseline experiments three times for each $k$. We aim to answer the following research questions (RQ):

---

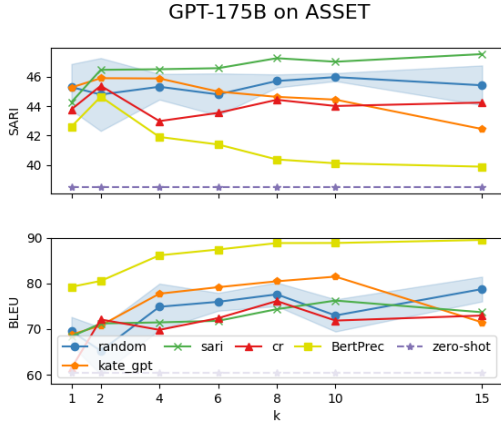Figure 1: GPT-175B results on ASSET. Top: SARI scores, Bottom: BLEU scores.



Figure 2: GPT-6.7B results on ASSET. Top: SARI scores, Bottom: BLEU scores.

**RQ1:** How do different metric-based selection techniques compare? (§5.1)

**RQ2:** Is metric-based sample selection robust to the order of the prompts? (§5.2)

**RQ3:** How does metric-based ICL compare to state-of-the-art text simplification methods? (§5.3)

**RQ4:** Does metric-based selection performance on one dataset transfer to other out-of-domain datasets? (§5.4)

### 5.1 RQ1: Effect of Metrics

Our main results with our default settings (GPT-175B on ASSET) is shown in Fig. 1. First of all, we observe that the random baseline is quite strong on average, however, with a **large variation** for most $k$ values; while zero-shot results are quite weak for all datasets. Interestingly, SARI-based selection consistently leads to the highest SARI scores for $k > 1$, while BERTPrec-based selection gives the highest BLEU and lowest SARI scores consistently for each $k$. Kate-GPT follows BERTPrec-based selection for the BLEU score, while providing results on par or lower than the random baseline for the SARI score.

Next, we check whether our findings hold for smaller models. In Fig. 2, we show the results of our smallest model, GPT-6.7B on the ASSET dataset. Since the zero-shot results were significantly lower than $k = 1$, we show them in Table 1, rather than plotting. Not surprisingly the highest SARI scores are achieved via the largest model; however, the opposite is not true for BLEU. The smallest model achieves the highest BLEU scores
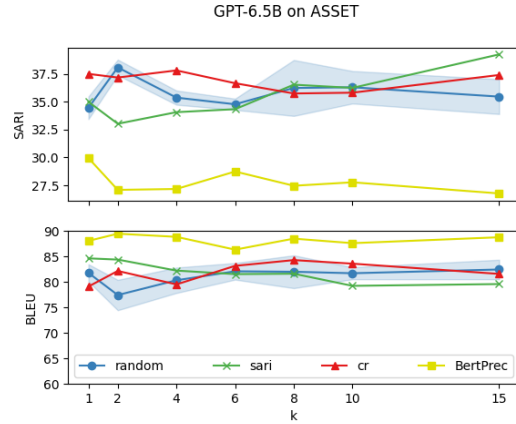
that raises another warning flag for using BLEU for TS evaluation.

Similar to the larger model, BERTPrec-based selection achieves the highest BLEU, and the lowest SARI scores. SARI-based selection provides considerably high SARI scores only for larger $k$s, suggesting the implicit controlling mechanism does not exist, or is only triggered with more samples. We also observe that CR performs relatively better on GPT-6.7B which suggests compression provides a stronger signal (e.g., deletion, shorter tokens) that can be utilized better by smaller models for simplification.

Finally, we investigate how the quality of the dataset affects the metric-based selection techniques, i.e., whether they are robust to noise. Fig 3 shows an overview of the SARI scores from all models on the noisy (i.e., TurkCorpus) and the cleaner (i.e., ASSET) dataset. Even though the general patterns are visible, the results on TurkCorpus are moderately less conclusive.

| Dataset | Model | SARI | BLEU |
|---|---|---|---|
| | GPT-175B | **32.17** | 42.34 |
| TurkCorpus | GPT-17B | 27.19 | 38.35 |
| | GPT-6.7B | 24.13 | **57.14** |
| | GPT-175B | **38.49** | 60.48 |
| ASSET | GPT-17B | 30.45 | 40.13 |
| | GPT-6.7B | 26.28 | **69.83** |

Table 1: Zero-shot results

**Selection Metric versus Evaluation Metric** Even though this was not one of our main research questions, we observe a strong relation between the metric used for MBL and the metric used for evaluation. For all the model and
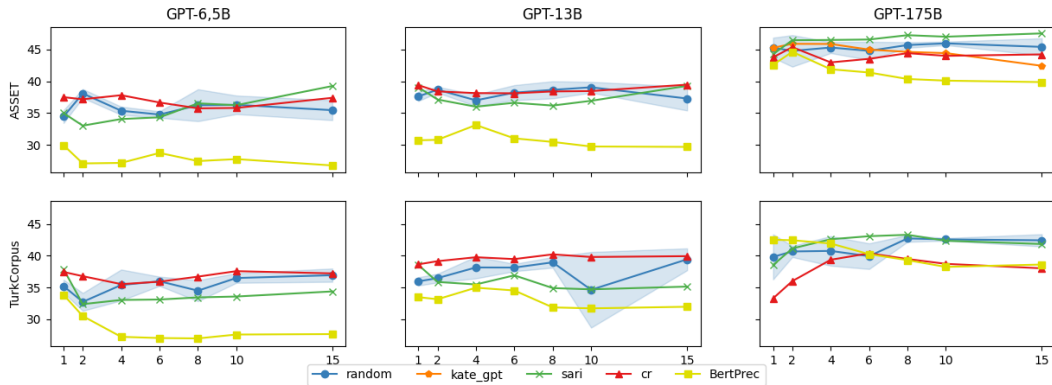
Figure 3: SARI scores for GPT-6.7B, GPT-13B and GPT-175B models on ASSET (top) and TurkCorpus (bottom) datasets. See App. B for BLEU scores.

dataset size settings, we observe that BLEU scores are consistently higher when the examples are selected via BERTScore_Precision. When we evaluate with the SARI score, SARI-metric behaves similarly for the GPT-175B model, however CR-metric performs better for the smaller models. More evidence for the relation between BLEU and BERTScore_Precision can be found in Appendix B. This suggests that the behaviour of large-enough GPT models can be *implicitly controlled* via MBL, which paves the way to a new research direction and needs further investigation.

## 5.2 RQ2: Effect of Order

Previous research (Lu et al., 2022) has shown that the order of the examples may have a significant impact on ICL performance. Commonly used orderings (Lu et al., 2022) include sorting from highest to lowest quality example, vice versa, and random selection. Inspired by these findings, we investigate the robustness of our selection metrics across sample orders. To do so, for each metric we perform three different order arrangements, namely as highest → lowest, lowest → highest, and random ordering for each metric. To have enough variation, we only experiment with $k = 6, 8, 10, 15$. As the baseline, we randomly pick samples and arrange them in 3 different randomized orders.

In Fig 4, we show how the performance of GPT-175B varies on ASSET when the samples that are i) picked randomly, ii) by SARI-based selection and iii) by BERTPrec-based selection are reordered following the above setup. As can be seen, the best-performing metrics, are also the most robust compared to others. To elaborate, SARI-based selection that provided the highest SARI scores has

the least variation, i.e., most robust to order; while BERTPrec-based selection provides the most stable BLEU scores along with the highest.

## 5.3 RQ3: Comparison to State-of-the-art

Finally, we compare our best and average model settings to state-of-the-art fine-tuned models [7]. The results are given in Table 2. Here, Random and SARI averages are calculated from §5.1 results, averaged over all $k$, with random selection being additionally averaged over all three random selections. These averages are reported for GPT-175B results, because it is generally the best model when considering averages across both datasets. As can be seen, the GPT-175B model with SARI-based selection outperforms existing results on all datasets, followed by random best and SARI-averaged. The exact settings (number of examples, model, and ordering) for SARI and Random Best can be found Appendix A.

## 5.4 RQ4: Task Transfer

In order to evaluate the suitability of our approach for unseen tasks and datasets, we experiment with choosing samples from a tune set and testing the performance on an unseen set. Here, we use ASSET and TurkCorpus as the tune set and evaluate on all three datasets: ASSET, TurkCorpus, and FestAbility. To investigate different metrics and language models, we perform experiments with GPT-175B with SARI-based selection and GPT-13B with CR-based selection as both metric selection techniques are generally best on those respective models. We

---

[7]The models which do not provide SOTA (e.g., KATE) are not included in the Table. The statistical significance cannot be provided since there is only one setting for the few-shot setting.
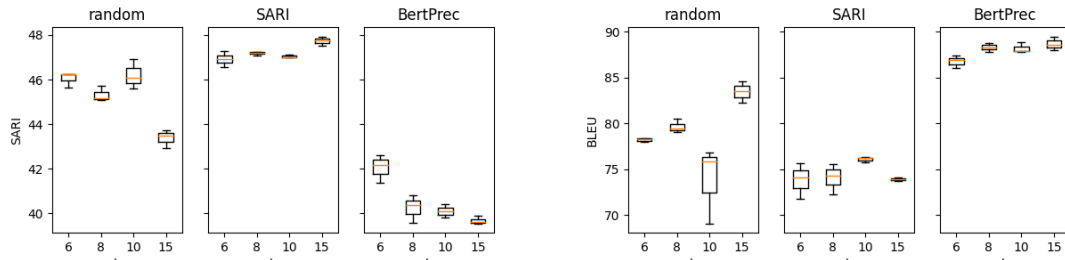
Figure 4: Boxplots for GPT-175B model performance on ASSET with sample (re)ordering via random, SARI-based and BERTPrec-based selections. Performance shown in SARI (left), and BLEU (right)

| Model | ASSET | TurkCorpus | FestAbility |
|---|---|---|---|
| Finetuned T5 | 45.04 | 43.31 | N/A |
| MUSS (BART + ACCESS) | 43.63 | 42.62 | N/A |
| BART-Large+Classifier | 38.76 | N/A | 27.13 |
| (Ours) Random-Best | 46.93 | 43.14 | N/A |
| (Ours) Random-Average | 45.33 | 40.32 | N/A |
| (Ours) MBL-Best | **47.94** | **43.46** | **44.86** |
| (Ours) MBL-Average | 46.63 | 41.78 | 43.55 |

Table 2: Comparison to TS state-of-the-art models. Random- and MBL-best examples are selected from the top examples in all experiments run. The best results are shown in bold. For more information on the exact settings for MBL and Random Best, see Appendix A.

use the best experimental settings from Table 2 in out-of-domain settings, comparing them with their in-domain counterparts. For example, the best setting for TurkCorpus is k=6 with high to low ordering (see Appendix A for more details on optimal experiment settings), so we compare the results of the model when given this setting on both the TurkCorpus and ASSET datasets. State-of-the-art results in this table refer to the best setting for in-domain experiments (i.e. ASSET evaluated on ASSET or TurkCorpus evaluated on TurkCorpus).

Our results are given in Table 3. For easy comparison, the Table also includes in-domain selection results as well as the current state-of-the-art scores taken from Table 2. In the first row, we observe that samples selected from TurkCorpus and tested on ASSET achieve significantly lower SARI scores than their in-domain variant for the GPT-175B setting, whereas the gap is lower for the GPT-13B. On the other hand, for the TurkCorpus test setting (second row), we see that GPT-175B model prompted with the best ASSET examples achieves even better results than the in-domain setting, suggesting a highly successful transfer. This ability cannot be observed for the GPT-13B model with CR-based selection. The final row shows the transfer results to another related but different task. It is apparent

that both models prompted with ASSET examples achieve marginally higher scores than the TurkCorpus ones.

Taking a look at the BLEU scores, we see that out-of-domain configurations on the TurkCorpus and ASSET datasets generally tend to match or even exceed their in-domain counterparts, suggesting a successful transfer. However, on the FestAbility dataset, we observe notably low BLEU scores, which are in-part due to the nature of FestAbility, in which sentences are often simplified in unconventional ways. Additionally, FestAbility sentences are extremely short, with only 1452 unique tokens in the original sentences and 996 unique tokens in the simplified sentences (Chamovitz and Abend, 2022), leading to unconventional results.

| Test Set | Model Setting | Tune Set | SARI | BLEU |
|---|---|---|---|---|
| ASSET | GPT-175B, SARI, high to low, k=6 | TurkCorpus | 43.46 | 79.83 |
| | GPT-175B, SARI, high to low, k=6 | ASSET | 46.93 | 75.67 |
| | GPT-13B, CR, high to low, k=15 | TurkCorpus | 41.73 | 74.57 |
| | GPT-13B, CR, high to low, k=15 | ASSET | 41.9 | 76.49 |
| | *State-of-the-art (MBL-Best)* | | *47.94* | *73.92* |
| | *Zero-shot (GPT-175B)* | | *38.49* | *60.48* |
| TurkCorpus | GPT-175B, SARI, random, k=15 | ASSET | 42.37 | 64.52 |
| | GPT-175B, SARI, random, k=15 | TurkCorpus | 41.48 | 85.89 |
| | GPT-13B, CR, high to low, k=15 | ASSET | 39.44 | 71.15 |
| | GPT-13B, CR, high to low, k=15 | TurkCorpus | 40.37 | 73.83 |
| | *State-of-the-art (MBL-Best)* | | *43.46* | *79.83* |
| | *Zero-Shot (GPT-175B)* | | *32.17* | *42.34* |
| FestAbility | GPT-175B, SARI, random, k=6 | TurkCorpus | 42.24 | 20.76 |
| | GPT-175B, SARI, random, k=15 | ASSET | 44.86 | 17.08 |
| | GPT-13B, CR, high to low, k=15 | TurkCorpus | 25.46 | 23.37 |
| | GPT-13B, CR, high to low, k=15 | ASSET | 36.63 | 12.01 |
| | *State-of-the-art (MBL-Best)* | | *44.86* | *N/A* |
| | *Zero-shot (GPT-175B)* | | *40.77* | *6.9* |

Table 3: ICL out-of-domain results for GPT-175B, SARI-based selection and GPT-13B, CR-based selection. Examples are picked from the *Tune Set* and tested on the *Test Set*. Zero-shot results are from GPT-175$B$ and given in the final rows for each dataset.

## 6 Qualitative Analysis

In this section, we perform a qualitative analysis of different model generated simplifications and

| Metric | Top 2 Examples |
|---|---|
| Compression Ratio | **Complex Sentence** They manifest with either neurological complications or with skin problems (or occasionally both). <br> **Simple Sentence** They show either brain or skin problems (or both). |
| | **Complex Sentence** The psychological state of sympathy is closely linked with that of compassion, empathy and empathic concern. <br> **Simple Sentence** Sympathy is closely linked with compassion and empathy. |
| BertScore Precision | **Complex Sentence** Sthenurine forelimbs were long with two extra-long fingers and claws compared with the relatively small, stiff arms of modern macropods. <br> **Simple Sentence** Sthenurine forelimbs were long with two extra-long fingers and claws compared with the small, stiff arms of modern macropods. |
| | **Complex Sentence** In 1828, Coenraad Johannes van Houten developed the first cocoa powder producing machine in the Netherlands. <br> **Simple Sentence** In 1828, Coenraad Johannes van Houten created the first cocoa powder producing machine in the Netherlands. |
| SARI | **Complex Sentence** The organic matter in soil derives from plants and animals. <br> **Simple Sentence** The organic matter in soil comes from plants and animals. |
| | **Complex Sentence** Dennis Lee Hopper (born May 17, 1936) is an American actor, filmmaker and artist. <br> **Simple Sentence** Dennis Lee Hopper was born on May 17, 1936. He is an American actor, filmmaker and artist. |

Table 4: Top 2 examples from each applicable selection metric (random and KATE-GPT selection were not applicable). All samples taken from the ASSET Validation dataset. We color rephrases first in blue and then in yellow, mark significant deletions in red, and underline sentence splits.

metric-based prompting examples in order to better understand how different settings affect model outputs.

## 6.1 Explaining Performance as $k$ Increases

We aim to understand why certain metrics (BERT-Prec and KATE-GPT) tend to perform worse as $k$ increases, while other metrics (SARI) tend to perform better as $k$ increases when evaluated on SARI scores (as seen in Fig. 3). In fact, this result is commonly seen in other papers (Zhao et al., 2021; Zhang et al., 2022), where they describe that adding more training examples can sometimes hurt accuracy. By analyzing output of metric-based selection on a fixed dataset and model (ASSET, GPT-175B) seen in Appendix D.3, we aim to understand the performance of different metrics as the number of examples, or $k$, increases. Our analysis focuses on three different metrics (KATE-GPT, BERTPrec, and SARI) and a particularly difficult example due to its unconventional subject nature, multiple abbreviations, and unknown words, and objectively confusing sentence structure. In general, we see from earlier trends that KATE-GPT and BERTPrec-selected examples tend to get worse (w.r.t SARI) as $k$ increases (see Figure 1-top). We also observe this qualitatively, as $k$ increases, KATE-GPT and

BERTPrec examples become closer to the original sentence, with BERTPrec generations even matching the original sentence at $k = 15$. However, as the value of k increases, SARI-selected examples show an improvement in quality. We observe that examples selected using the SARI score metric tend to: i) split sentences more frequently, and ii) decode potentially confusing abbreviations, such as "OEL".

**Sentence Splitting:** SARI-selected examples are more prone to splitting sentences (see $k$=2,15), which may be in-part due to the style of the top SARI examples, which include sentence splitting; while this is not present in any of the other metrics. See Appendix D.3 for examples. Sentence splitting is correlated with increased human comprehension of TS outputs (Williams et al., 2003). This is particularly interesting because it leads us to infer that models can potentially learn the "style" of the reference sentences.

**Abbreviations:** In all three cases, ($k = 2, 8, 15$) examples selected by SARI score remove the potentially confusing abbreviation "OEL" and instead replace it with either "original English-language" or "English-language", while KATE-GPT and BERT-Prec selected examples only exhibit this behavior for $k = 2$ (see Appendix D.3).

## 6.2 Analyzing Model Size

Model size plays a significant role in output sentences, with smaller models (especially GPT-6.7B) tending to change very little structurally from the original sentence, regardless of the metric used to select examples. See Appendix D for a complete list of model outputs on all metrics for the original sentence "OEL manga series Graystripe's Trilogy There is a three volume original English-language manga series following Graystripe, between the time that he was taken by Twolegs in Dawn until he returned to ThunderClan in The Sight". From these results, we conclude that GPT-6.7B tends to hardly change sentences at all, with both Random and BERTPrec-selected examples having no change from the original sentence. SARI-selected examples adds a comma, but CR-selection prompts the model to rephrase key parts of the sentence. GPT-13B performs considerably better when looking at a qualitative analysis, as all examples have removed "OEL manga series Graystripe's Trilogy" and restructured the sentence to be more concise, and SARI-selected going as far to remove an ambiguous abbreviation "OEL". These qualitative observations are consistent with our results from Figure 1.

## 6.3 Analyzing Top Metric-Selected Examples

In this section, we analyze the top metric-selected examples for compression ratio, BERTPrec, and SARI. In Table 4 we include the top 2 examples for each metric from the ASSET validation dataset, and in Appendix C we include the remaining top 8 examples for SARI and BERTPrec selection.

Looking at the style of both BERTPrec and SARI score, both metrics' top examples barely change from the original sentences, often only changing one or two words (i.e., movie → film) but leaving the rest unchanged, primarily using deletion or rewriting operations. However, in CR top examples, we see extreme deletions from the original sentences and several rewriting operations done (which is consistent with our understanding of the compression ratio). Notably, we also see that top SARI examples are the only metric that use sentence splitting (see the 2nd example under SARI from Table 4).

## 7 Conclusion and Future Work

In conclusion, we propose a novel and robust method for selecting examples in the TS domain, evaluating its effectiveness on multiple well-known TS datasets and even on downstream tasks like cognitive simplification. Our experiments demonstrate state-of-the-art results in the field of TS and CS, reaching scores of 44.86 on FestAbility, 47.94 on ASSET and 43.46 on TurkCorpus. We hope that future work will generalize our findings in other text generation tasks and other domains.

## Limitations

Our approach is computationally and financially intensive, especially on the GPT-175B model, which limits its scalability to smaller, open-source models. While our approach has shown strong results in the TS domain, we are not yet sure whether using domain-specific selection methods is widely applicable. Our approach also is not applicable in true few-shot settings in which a large validation set is not available to select examples from. Also, our approaches' scalability to other downstream TS tasks outside of cognitive simplification is yet to be tested, especially in different domains. We tested on two well-known TS datasets (ASSET and TurkCorpus), and we did not test on another known TS dataset, Newsela, due to its restrictive licensing. Additionally, we have tested our approach on example numbers up to 15 due to financial constraints, and testing on higher numbers of examples may show additional insights and we leave this for future researchers.

## Ethics Statement

We acknowledge that while our approach reaches high scores on datasets aimed for individuals with disabilities, further research and evaluation from humans with specific disabilities listed in this paper is crucial to determine the true effectiveness of our approach in these scenarios.

# References

Fernando Alva-Manchego, Louis Martin, Antoine Bordes, Carolina Scarton, Benoît Sagot, and Lucia Specia. 2020. ASSET: A dataset for tuning and evaluation of sentence simplification models with multiple rewriting transformations. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4668–4679, Online. Association for Computational Linguistics.

Fernando Alva-Manchego, Louis Martin, Carolina Scarton, and Lucia Specia. 2019. EASSE: Easier automatic sentence simplification evaluation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP): System Demonstrations*, pages 49–54, Hong Kong, China. Association for Computational Linguistics.

Fernando Alva-Manchego, Carolina Scarton, and Lucia Specia. 2021. The (un)suitability of automatic evaluation metrics for text simplification. *Computational Linguistics*, 47(4):861–889.

Eduard Barbu, Maria Martín-Valdivia, Eugenio Martínez-Cámara, and L. López. 2015. Language technologies applied to document simplification for helping autistic people. *Expert Systems with Applications*, 42:5076–5086.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020a. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020b. Language models are few-shot learners.

Eytan Chamovitz and Omri Abend. 2022. Cognitive simplification operations improve text simplification. In *Proceedings of the 26th Conference on Computational Natural Language Learning (CoNLL)*, pages 241–265.

Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, Lei Li, and Zhifang Sui. 2023. A survey for in-context learning.

Hila Gonen, Srini Iyer, Terra Blevins, Noah A. Smith, and Luke Zettlemoyer. 2022. Demystifying prompts in language models via perplexity estimation. *CoRR*, abs/2212.04037.

J Peter Kincaid, Robert P Fishburne Jr, Richard L Rogers, and Brad S Chissom. 1975. Derivation of new readability formulas (automated readability index, fog count and flesch reading ease formula) for navy enlisted personnel. Technical report, Naval Technical Training Command Millington TN Research Branch.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 7871–7880. Association for Computational Linguistics.

Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2022. What makes good in-context examples for gpt-3? In *Proceedings of Deep Learning Inside Out: The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures, DeeLIO@ACL 2022, Dublin, Ireland and Online, May 27, 2022*, pages 100–114. Association for Computational Linguistics.

Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8086–8098, Dublin, Ireland. Association for Computational Linguistics.

Louis Martin, Angela Fan, Éric de la Clergerie, Antoine Bordes, and Benoît Sagot. 2022. MUSS: Multilingual unsupervised sentence simplification by mining paraphrases. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 1651–1664, Marseille, France. European Language Resources Association.

Islam Nassar, Michelle Ananda-Rajah, and Gholamreza Haffari. 2019. Neural versus non-neural text simplification: A case study. In *Proceedings of the The 17th Annual Workshop of the Australasian Language Technology Association, ALTA 2019, Sydney, Australia, December 4-6, 2019*, pages 172–177. Australasian Language Technology Association.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Jipeng Qiang, Yun Li, Zhu Yi, Yunhao Yuan, and Xindong Wu. 2020. Lexical simplification with pre-trained encoders. *Thirty-Fourth AAAI Conference on Artificial Intelligence*, page 8649–8656.

Luz Rello, Ricardo Baeza-Yates, Stefan Bott, and Horacio Saggion. 2013. Simplify or help? text simplification strategies for people with dyslexia. In *Proceedings of the 10th International Cross-Disciplinary Conference on Web Accessibility*, W4A '13, New York, NY, USA. Association for Computing Machinery.

Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2022. Learning to retrieve prompts for in-context learning. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2655–2671, Seattle, United States. Association for Computational Linguistics.

Kim Cheng Sheang and Horacio Saggion. 2021. Controllable sentence simplification with a unified text-to-text transfer transformer. In *Proceedings of the 14th International Conference on Natural Language Generation*, pages 341–352, Aberdeen, Scotland, UK. Association for Computational Linguistics.

Taylor Sorensen, Joshua Robinson, Christopher Michael Rytting, Alexander Glenn Shaw, Kyle Jeffrey Rogers, Alexia Pauline Delorey, Mahmoud Khalil, Nancy Fulda, and David Wingate. 2022. An information-theoretic approach to prompt engineering without ground truth labels. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 819–862. Association for Computational Linguistics.

Sanja Stajner. 2021. Automatic text simplification for social good: Progress and challenges. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 2637–2652, Online. Association for Computational Linguistics.

Elior Sulem, Omri Abend, and Ari Rappoport. 2018. BLEU is not suitable for the evaluation of text simplification. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 738–744, Brussels, Belgium. Association for Computational Linguistics.

Teerapaun Tanprasert and David Kauchak. 2021. Flesch-kincaid is not a text simplification evaluation metric. In *Proceedings of the 1st Workshop on Natural Language Generation, Evaluation, and Metrics (GEM 2021)*, pages 1–14, Online. Association for Computational Linguistics.

Sandra Williams, Ehud Reiter, and Liesl Osman. 2003. Experiments with discourse-level choices and readability. In *Proceedings of the 9th European Workshop on Natural Language Generation (ENLG-2003) at EACL 2003*, Budapest, Hungary. Association for Computational Linguistics.

Wei Xu, Courtney Napoles, Ellie Pavlick, Quanze Chen, and Chris Callison-Burch. 2016. Optimizing statistical machine translation for text simplification. *Transactions of the Association for Computational Linguistics*, 4:401–415.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with BERT. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

Yiming Zhang, Shi Feng, and Chenhao Tan. 2022. Active example selection for in-context learning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pages 9134–9148. Association for Computational Linguistics.

Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 12697–12706. PMLR.

## A    Optimal Settings

In this section, we provide full optimal settings in Table 5 for the "Random Best" and "MBL Best" models. These settings are all in-domain (i.e., AS-SET Validation, ASSET Tune; MTurk Validation, MTurk Tune) and include model size, $k$ (number of examples), and ordering (high/low, low/high, and random).

| Metric and Dataset | Model | Metric | k | Ordering |
|---|---|---|---|---|
| MBL Best TurkCorpus | GPT-175B | SARI | 6 | High/Low |
| Random Best TurkCorpus | GPT-175B | SARI | 8 | Low/High |
| MBL Best ASSET | GPT-175B | SARI | 15 | Random |
| Random Best ASSET | GPT-175B | SARI | 10 | Random |

Table 5: MBL-Best and Random-Best settings for results

## B    BLEU Results

In this section, we include full BLEU results in Figure 5 including all model sizes (175B, 13B, 6.7B), datasets (TurkCorpus and ASSET) and selection techniques.

## C    Top ASSET Examples

In this section, we include extended results from 4, with the top 3-10 results from the ASSET Validation dataset based on top SARI (Section C.1), CR, and BERTPrec (Section C.2) scores.

### C.1    SARI

In this section, we include the top 3 to 10 examples based on SARI score selected from the ASSET Validation dataset in Table 6.

### C.2    BERTPrec

In this section, we include the top 3 to 10 examples based on SARI score selected from the ASSET Validation dataset in Table 7.

## D    Selected Model Generated Outputs

In this section, we analyze select model generated outputs on 4 (5 for GPT-175B) different example-selection methods (Random, SARI, CR, BERTPrec and optionally KATE-GPT) on different models for in-domain configurations of the ASSET dataset. The original sentence in all of these is "OEL manga series Graystripe's Trilogy There is a three volume original English-language manga series following Graystripe, between the time that he was taken by Twolegs in Dawn until he returned to Thunder-Clan in The Sight." This sentence was specifically picked from the ASSET/TurkCorpus test dataset based on three reasons: 1) complexity (potentially confusing abbreviations and unconventional sentence structure) 2) length 3) unfamiliar/domain-specific terms from "Warrior Cats" (e.g. "Thunder-Clan" and "Twolegs"). §D.1 includes generations on GPT-13B, §D.2 includes generatons on GPT-6.7B, and §D.3 includes generations on GPT-175B.

### D.1    Selected GPT-13B Generations

In this section, we include generations from the original sentence mentioned in D on GPT-13B on the ASSET Test set.

### D.2    Selected GPT-6.7B Generations

In this section, we include generations from the original sentence mentioned in D on GPT-6.7B on the ASSET Test set.

### D.3    Selected GPT-175B Generations

In this section, we include generations from the original sentence mentioned in D on GPT-175B on the ASSET Test set. Text in red indicates text that has been successfully been changed from the abbreviation "OEL" to an interpretable phrase (either "English-language" or "original English-language").

| k | Example |
|---|---------|
| 3 | **Complex Sentence** It is adjacent to Lord Wandsworth College.<br>**Simple Sentence** it is next to Lord Wandsworth College. |
| 4 | **Complex Sentence** He took the post of chief conductor of the Netherlands Radio Philharmonic in 1957.<br>**Simple Sentence** He became the chief conductor of the Netherlands Radio Philharmonic in 1957. |
| 5 | **Complex Sentence** It was discovered on February 27, 1995.<br>**Simple Sentence** It was found on February 27, 1995. |
| 6 | **Complex Sentence** Surnames Aaron Schock, member of the U. S. House of Representatives representing the 18th district of Illinois.<br>**Simple Sentence** Surnames Aaron Schock is a member of the U. S. House of Representatives. He represents the 18th district of Illinois. |
| 7 | **Complex Sentence** Mork holds a Professorship at the Norwegian Academy of Music, Oslo.<br>**Simple Sentence** Mork is a Professor at the Norwegian Academy of Music. |
| 8 | **Complex Sentence** The Hubble Space Telescope observed Fortuna in 1993.<br>**Simple Sentence** The Hubble Space Telescope saw Fortuna in 1993. |
| 9 | **Complex Sentence** The lithosphere is underlain by the asthenosphere, the weaker, hotter, and deeper part of the upper mantle.<br>**Simple Sentence** The lithosphere is supported by the asthenosphere, the weaker, hotter, and deeper part of the upper mantle. |
| 10 | **Complex Sentence** The Beatles famously included his face on the cover of Sgt. Pepper's Lonely Hearts Club Band (Guy and Llewelyn-Jones 2004, 111).<br>**Simple Sentence** The Beatles put his face on the cover of Sgt. Pepper's Lonely Hearts Club Band. |

Table 6: Top 3-10 Examples from SARI, ASSET Validation dataset.

| k | Example |
|---|---------|
| 3 | **Complex Sentence** The Convent has been the official residence of the Governor of Gibraltar since 1728. <br> **Simple Sentence** The Convent has been the residence of the Governor of Gibraltar since 1728. |
| 4 | **Complex Sentence** Scholarships, Academic Awards, Flying Eagle Awards and Improvement Awards are given to students with outstanding academic achievements. <br> **Simple Sentence** Scholarships, Academic Awards, Flying Eagle Awards and Improvement Awards are given to students with academic achievements. |
| 5 | **Complex Sentence** The blood vessels in the human body include arteries, veins and capillaries. <br> **Simple Sentence** The blood vessels in the human body are called arteries, veins and capillaries. |
| 6 | **Complex Sentence** Frederick had a summer residence built there for Sophie Charlotte by the architect Johann Arnold Nering between 1695 and 1699. <br> **Simple Sentence** Frederick had a summer residence built for Sophie Charlotte by the architect Johann Arnold Nering between 1695 and 1699. |
| 7 | **Complex Sentence** The film stars Al Pacino, John Cazale, Chris Sarandon, James Broderick, and Charles Durning. <br> **Simple Sentence** The movie stars Al Pacino, John Cazale, Chris Sarandon, James Broderick, and Charles Durning. |
| 8 | **Complex Sentence** According to an interview in the UK newspaper The Sun, Heyman wrote the brand's weekly scripts and submitted them to writers for possible changes, and then Vince McMahon for final approval. <br> **Simple Sentence** According to an interview in the UK newspaper The Sun, Heyman wrote the brand's weekly scripts and sent them to writers for possible changes, and then Vince McMahon for final approval. |
| 9 | **Complex Sentence** In March 2001, the World Wrestling Federation purchased World Championship Wrestling. <br> **Simple Sentence** In March 2001, the World Wrestling Federation bought World Championship Wrestling. |
| 10 | **Complex Sentence** Becker defeated Jim Courier in straight sets to win the 1992 year-end ATP Tour World Championships in Frankfurt. <br> **Simple Sentence** Becker defeated Jim Courier in straight sets to win the 1992 year-end ATP Tour World Championships. |

Table 7: Top 3-10 Examples from BERTPrec, ASSET Validation dataset.

| Metric | Curie Simplifications |
|--------|----------------------|
| Random | there is a three volume manga series following graystripe between the time that he was taken by twolegs in dawn and he returned to thunderclan in the sight . |
| SARI | There is a three-volume original English-language manga series following graystripe, between the time he was taken by twolegs in dawn and he returned to thunderclan in the sight. |
| CR | There is a three volume manga series following graystripe between the time he was taken by twolegs in dawn and he returned to thunderclan in the sight. |
| BERTPrec | There is a three volume original English-language manga series following graystripe , between the time that he was taken by twolegs in dawn and he returned to thunderclan in the sight . |

Table 8: Curie Simplifications for a selected complex sentence, ran on ASSET.
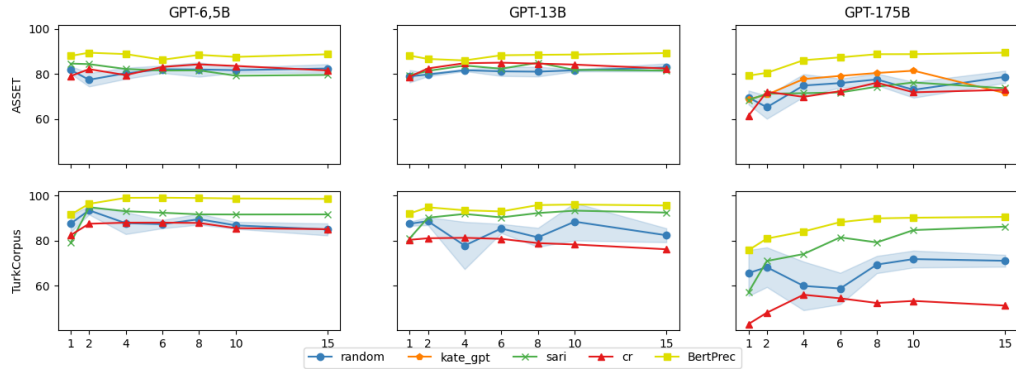
Figure 5: BLEU scores for GPT-6.5, GPT-13B and GPT-175B models on ASSET (top) and TurkCorpus (bottom) datasets.

| Metric | Curie Simplifications |
|---|---|
| Random | oel manga series graystripe 's trilogy there is a three volume original english-language manga series following graystripe between the time that he was taken by twolegs in dawn until he returned to thunderclan in the sight . |
| SARI | oel manga series graystripe 's trilogy there is a three volume original english-language manga series following graystripe , between the time that he was taken by twolegs in dawn until he returned to thunderclan in the sight . |
| CR | OEL manga series, Graystripe 's Trilogy, has a three volume original English-language manga series following Graystripe, between the time that he was taken by twolegs in dawn until he returned to thunderclan in the sight. |
| BERTPrec | oel manga series graystripe 's trilogy there is a three volume original english-language manga series following graystripe , between the time that he was taken by twolegs in dawn until he returned to thunderclan in the sight . |

Table 9: Babbage Simplifications for a selected complex sentence, ran on ASSET.

| Original | OEL manga series Graystripe's Trilogy There is a three volume original English-language manga series following Graystripe, between the time that he was taken by Twolegs in Dawn until he returned to ThunderClan in The Sight. |
|---|---|
| **Reference Simplification** | There is a three volume original English-language manga following Graystripe, between the time that he was taken by Twolegs in Dawn until he returned to ThunderClan in The Sight. |

**k=2**

| KATE-GPT | Graystripe's Trilogy is a three volume original English-language manga series following Graystripe from when he was taken by Twolegs to his return to ThunderClan. |
|---|---|
| **BERTPrec** | Graystripe's trilogy is a three-volume English-language manga series about his time between being taken by twolegs in Dawn and returning to ThunderClan in The Sight. |
| **SARI** | Graystripe's trilogy is a three-volume original English-language manga series. It follows Graystripe between the time he was taken by Twolegs in Dawn and when he returned to ThunderClan in The Sight. |

**k=8**

| KATE-GPT | The Graystripe's Trilogy is an OEL manga series that follows Graystripe's journey between the time he was taken by Twolegs in Dawn and his return to ThunderClan in The Sight. |
|---|---|
| **BERTPrec** | Oel manga series Graystripe's trilogy follows Graystripe between being taken by twolegs in Dawn and returning to ThunderClan in The Sight. |
| **SARI** | Graystripe's trilogy is a three volume English-language manga series about his time away from ThunderClan in Dawn. |

**k=15**

| KATE-GPT | The OEL manga series Graystripe's Trilogy follows Graystripe between the time he was taken by Twolegs in Dawn and when he returned to ThunderClan in The Sight. |
|---|---|
| **BERTPrec** | Oel manga series graystripe's trilogy is a three-volume English-language manga series following Graystripe between his capture by twolegs in Dawn and his return to ThunderClan in The Sight. |
| **SARI** | Graystripe's trilogy is a three-volume English-language manga series. It covers the time Graystripe was taken by twolegs in Dawn until he returned to ThunderClan in The Sight. |

Table 10: Model-generated simplifications on ASSET, GPT-175B.

# Exploring the Naturalness of Cognitive Status Informed Referring Form Selection Models

**Gabriel Del Castillo**\* and  **Grace Clark**\* and  **Zhao Han**\* and  **Tom Williams**

MIRRORLab

Department of Computer Science

Colorado School of Mines

gdelcastillo@mines.edu, geclark@mines.edu, zhaohan@mines.edu, twilliams@mines.edu

## Abstract

Language-capable robots must be able to efficiently and naturally communicate about objects in the environment. A key part of communication is *Referring Form Selection* (RFS): the process of selecting a form like *it*, *that*, or *the N* to use when referring to an object. Recent cognitive status-informed computational RFS models have been evaluated in terms of goodness-of-fit to human data. But it is as yet unclear whether these models actually select referring forms that are any more *natural* than baseline alternatives, regardless of goodness-of-fit. Through a human subject study designed to assess this question, we show that even though cognitive status-informed referring selection models achieve good fit to human data, they do not (yet) produce concrete benefits in terms of naturalness. On the other hand, our results show that human utterances also had high variability in perceived naturality, demonstrating the challenges of evaluating RFS naturality.

**Keywords:** Referring form selection (RFS), computational models, naturalness, Givenness Hierarchy, cognitive status

## 1 Introduction

Referring is a critical part of human communication, especially in situated, task-based interactions. Humans use a variety of referring forms during reference production, including both definite descriptions (e.g., *The red box on the table*) and concise referring forms (e.g., *this box*, *that*, or *it*). While more concise referring forms are less information-rich, they allow speakers to express their intentions more quickly, and allow their listeners in turn to more quickly and effectively infer those intentions (Gundel et al., 1993). The process of choosing what type of referring form to use, known as *Referring Form Selection*, is an important first step in the production of referring language (Krahmer and Van Deemter, 2012).



Figure 1: To investigate the naturalness of referring forms, we conducted a study where participants watched videos of human-human instruction tasks. After each video, participants were shown a referring utterance that could have followed as the next line in the dialogue, and were asked to assess the naturalness of that utterance.

Yet despite the wide variety of referring forms observed in human-human interaction, and the critical role of Referring Form Selection in language production, most research on generating referring language has focused solely on definite descriptions (Van Deemter, 2016; Krahmer and Van Deemter, 2012). While generating effective definite descriptions is a critical task, a speaker solely relying on this referring form would be an inefficient, unnatural, and annoying speaker. This discrepancy is critical not just for the psycholinguistics community, who seek to understand the cognitive dynamics of language production, but also for the Artificial Intelligence and Human-Robot Interaction communities, who seek to enable efficient, natural, and humanlike communication in task-based, situated domains (Tellex et al., 2013; Jackson and Williams, 2022; Cakmak and Thomaz, 2012; Williams et al., 2015; Gervits et al., 2021). As such, we argue that more attention to the problem of Referring Form Selection is needed across multiple areas of cognitive science.

---

\*The first three authors contributed equally to this work.

Recently, a number of researchers have begun to give the Referring Form Selection problem additional consideration, using a variety of experimental and machine learning research techniques (Same and van Deemter, 2020; Pal et al., 2021; Chen et al., 2021; Han et al., 2022; Spevak et al., 2022). For example, Chen et al. (2021) examined the features learned by deep learning models of Referring Form Selection; Han et al. (2022) addressed the ecological validity of the contexts in which Referring Form Selection is studied; Spevak et al. (2022) studied how document planning of task instructions could be designed for optimal referring form selection; and Pal et al. (2021) and Han et al. (2022) studied how models of cognitive status could be used to enable cognitively informed models of Referring Form Selection.

This last set of work is of particular interest: Pal et al. (2021) and Han et al. (2022) leveraged the well-validated (Gundel et al., 2010) Givenness Hierarchy theory (Gundel et al., 1993), a linguistic theory that captures the relation between different referring forms and the *cognitive status* of referents in listeners' minds. For example, per this theory, when a speaker uses *this*, one can infer they assume their target referent to be *activated* in their listener's mind; when a speaker uses *it*, one can infer they assume their target referent to be *in focus*.

Previous work on *cognitive status*-informed models of Referring Form Selection have largely been evaluated in terms of fit to human data using objective metrics like accuracy and notably human evaluations of these computational model in live human-robot interactions by Han and Williams (2023). That is, previous researchers have only assessed whether the referring forms predicted by their models *match* the referring forms that people actually use in human-human interactions.

While assessing fit to human data supports these models *as cognitive models*, it obfuscates a key dimension of Referring Form Selection: when a human selects a Referring Form during Referring Form Selection, there is no one "correct" form for them to select. In many contexts, for example, *the N'* and *that N'* may be relatively equally appropriate. Even when an object is truly in focus, warranting the use of the extremely concise *it*, the use of *the-N* is not *wrong*; and in fact, in some such cases, the use of *the-N* may be advantageous as it is simply more natural sounding.

As such, while cognitive status-informed models

of Referring Form Selection have been shown to achieve good fit to human data, (1) high goodness of fit may be an unnecessarily aggressive benchmark, and (2) it is unclear whether the referring forms selected by these models are actually any more natural than those that would be produced if simpler baseline models were used.

In this work, we thus compared the naturalness of referring forms selected by cognitive status-informed referring form selection models (specifically, that presented by Han et al. (2022)) to those that would be selected by a variety of baselines, including a random baseline, and a definite description baseline (in which a definite noun phrase *the N* is always used).

To do so, we conducted an experiment in which we modified a dataset of task-based referring expressions, systematically varied the referring forms shown to participants, allowing us to collect naturalness ratings for all possible referring forms that could have been used in those referring expressions. We then consider, for each of those referring expressions, what referring form each of the compared models would have predicted, and thus what the perceived naturalness would have been. Averaging these naturalness predictions for each model, we are able to compare the overall naturalness of the considered models.

As we will show, our results suggest that even though cognitive status-informed referring selection models achieve good fit to human data, they do not (yet) produce concrete benefits in terms of naturality. But our results also demonstrate the challenges of performing this type of evaluation, as even the utterances produced by humans had high variability in perceived naturality.

## 2   Related Work

Arnold and Zerkle (2019) argues that linguistic Referring Form Selection models generally fall into two categories: *rational* and *pragmatic*. *Rational* models (e.g. Aylett and Turk (2004); Frank and Goodman (2012)) could explain the use of pronouns from an egocentric perspective, i.e., in terms of their ease of use in conversations. *Pragmatic* models, on the other hand, could explain the use of pronouns from an allocentric perspective, i.e., in terms of the assumptions about interlocutors that lead to their use. These allocentric accounts are typically grounded in theoretical constructs like cognitive status (Grosz et al., 1995). Although

these pragmatic models vary in terms of the constructs they use to explain referring form choice (e.g., givenness (Gundel et al., 1993), and focus (Grosz et al., 1995; Brennan et al., 1987; Grosz and Sidner, 1986), these models are all centered around the assumption that referring form selection is based on the status a referent has in a conversation or in the mind of conversational participants.

While both of these models make important contributions to the literature, neither performs at exceptional levels when it comes to predicting which specific referring forms to use. As Arnold and Zerkle (2019) pointed out, rational models suggest using reduced forms vastly more often than seen in practice, and fail to predict referring forms that are equally short. Furthermore, Arnold and Zerkle (2019) and Grüning and Kibrik (2005) note that both kinds of models focus on individual events or factors, such as recency in conversation (Mann et al., 1989), instead of developing a fully comprehensive model for all of what reference production entails.

Artificial Intelligence (AI) researchers developing Referring Form Selection Models have the same problems (Ge et al., 1998; McCoy and Strube, 1999; Callaway and Lester, 2002; Poesio et al., 2004; Kibble and Power, 2004; Kibrik, 2011; Kibrik et al., 2016). AI RFS models can be broadly categorized as *multi-factorial process modeling*, where the prediction of referring forms is approached as a problem of classification based on various linguistic and contextual features. Much like previous models, those discussed by Kibrik (2011), Van Deemter et al. (2012) and Gatt et al. (2014) opt for predicting pronoun use in general, as opposed to predicting the use of specific referring forms. Additionally, models like those listed above tend to be trained using purely textual domains (e.g., Krasavina and Chiarcos (2007)'s) that are very different from situated domains. Situated domains are highly ambiguous, with large numbers of nearly identical objects, and require speakers to make run-time decisions based on linguistic features (like prosody) and non-linguistic features (like physical distance) that may be assessed using noisy sensors.

To fix these problems, Pal et al. (2020) presented dynamic models of cognitive status based on the Givenness Hierarchy (Gundel et al., 1993), which they used to produce cognitive status-informed RFS models. In addition to cognitive status, these models included situated features like physical distance, leading to promising results (Pal et al., 2021). Han et al. (2022) solved a number of external ecological issues in the task environment that referring form data was collected from, e.g., including repeated and non-present objects. This led to a wider variety of referring forms collected. However, both Pal et al. (2021) and Han et al. (2022) only evaluated their work in terms of model goodness-of-fit to human data[1]; that is, how well the model predicts the referring forms used by others. While goodness-of-fit is a valuable metric, it obscures the fact that when choosing a referring form, there are often multiple referring forms that might be appropriate. As such, focusing on goodness-of-fit simultaneously risks underselling the performance of a model and overselling the benefits of the model. Because of the ecological validity and the wide range of referring forms from Han et al. (2022)'s model, we used it in this work.

## 3 Hypothesis

Due to the incorporation of cognitive statuses of objects in listeners' mind, we believe that referring forms predicted by cognitive status-informed models will have higher perceived naturalness, compared to random selection and compared to use of definite descriptions alone.

## 4 Method

In order to evaluate the naturalness of different referring forms, we used a novel experimental design in which we (1) collected naturalness ratings for each of a large set of referring forms across a large set of referring contexts, and (2) used these ratings to determine the overall naturalness of several competing models, by seeing what the naturalness ratings *would have been* under the referring forms selected by those models across those referring contexts.

To collect naturalness ratings, we performed an online experiment in which participants watched a series of videos from a dataset of human-human task based interactions. After each video, the participant was shown, in text, the next utterance spoken

---

[1]In research performed in parallel, completed after, but formally published before this work, we also conducted in-person human evaluations of the naturalness of our referring form selection model (Han and Williams, 2023). While that work is beyond the scope of the present paper, readers may want to consult that paper, which reinforces and adds nuance to the result of this paper.
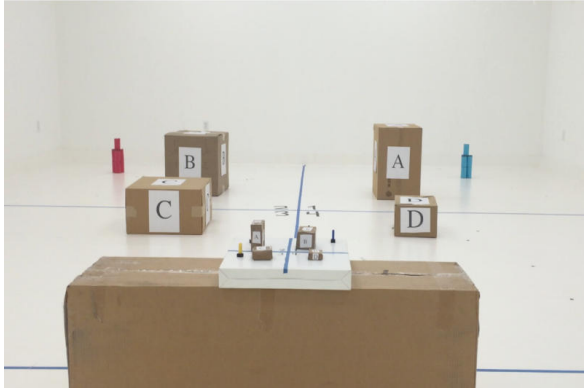
Figure 2: The task environment by Bennett et al. (2017) where a person instructs another person to re-configure objects to the layout of smaller models on the near table.

in the video, and asked to evaluate its naturalness.

We will now step through (1) the specific stimuli used in this experiment; (2) the experimental design that determined how these videos were shown to participants; (3) how we systematically varied the referring forms shown to participants and assessed the naturalness of those referring forms; (4) the overall experimental procedure; and (5) our participant demographics.

## 4.1 Stimuli

The videos we showed to participants were those collected by Bennett et al. (2017). This dataset contains videos from an experiment involving dyadic interactions in which one participant instructs another participant in how to rearrange a set of boxes and cans in order to match a desired configuration. The task environment is shown in Figure 2.

We selected ten videos from Bennett et al. (2017)'s dataset, and divided each into ten sub-videos, each of which ended immediately before the $n^{th}$ referring expression where $n \in \{1 \dots 10\}$. That is, for each video, we constructed ten excerpts, the first of which started at the beginning of the task and ended immediately before the first referring expression, the second started at the beginning of the task and ended immediately before the second referring expression, and so forth. All videos were subtitled for clarity. Figure 3 left shows a video.

We selected videos that contained a wide range of referring forms. As shown in Table 2, the distribution of referring forms in the original dataset is extremely skewed, with *the* $\langle N \rangle$ and $\langle N \rangle$ taking 85%. In contrast, the first three referring forms in the chosen videos approximately take 30% each.

## 4.2 Experimental Design

Each participant watched ten videos, each of which was an excerpt from a different one of the ten videos (i.e., video 0, 1, 2, 3, 4, 5, 6, 7, 8, or 9), and each of which ended at a different cutpoint, (i.e., immediately before referring expression 0, 1, 2, 3, 4, 5, 6, 7, 8, or 9 in that video). The average length of the videos was 35.6 seconds. The sequence of videos watched by each participant was selected using a Graeco-Latin square design (Grant, 1948) to ensure that each participant saw ten different videos of ten different interactions of ten different lengths while controlling for ordering effects.

## 4.3 Manipulations and Measures

As mentioned above, at the end of each video, participants were shown, in text, the utterance that immediately followed where the video cut off. These utterances were manipulated to vary the referring form used in the expression, with the actual referring form from the video replaced by one of the following: {it, this, this-N', that, that-N', the-N', N'}. For example, if the original utterance was "Now push **box D** to the left", participants were shown "Now push **it** to the left", "Now push **this box D** to the left" and so forth. These referring forms were selected at random for each video according to a pre-determined schedule. After being shown this "next utterance", participants were asked to rate its naturalness on a 5-point Likert item, with 1 being very unnatural and 5 being very natural.

## 4.4 Procedure

Participants first completed an informed consent form, read their task instructions, and answered demographic questions. Next, to ensure participants could hear what was said in the videos and avoid bots automatically filling out the questionnaires, participants performed an audio/video check. If participants passed this check, they proceeded to watch ten videos based on their randomly assigned Graeco-Latin square row, after each of which they answered the naturalness questions. Finally, participants answered an attention check question (asking the color of the walls in the task environment). This experiment's design and procedure were approved by the authors' institution.
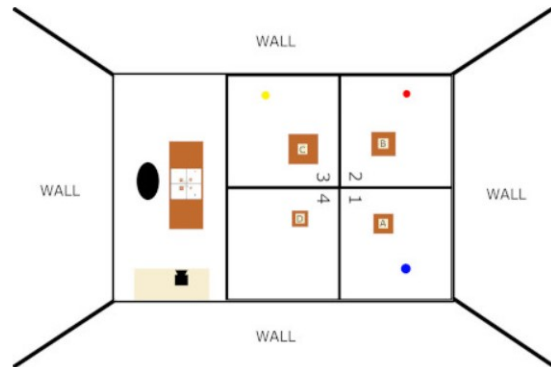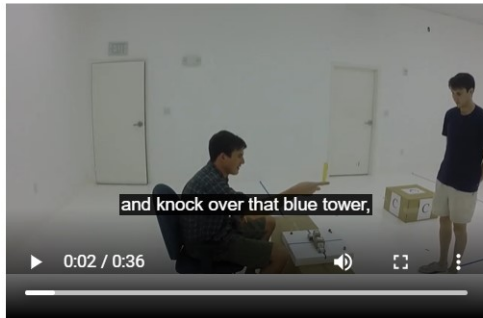
Figure 3: **Left**: The first video of ten videos that was playing and ended immediately before the first referring form. **Right**: The map used for participants to select the object being referred to (here, Box C). After the selection, the correct object would be similarly highlighted on the map.

## 4.5 Participants

361 participants were recruited globally from Prolific (Palan and Schitter, 2018) with a requirement that they must be fluent in English. 14 were excluded for failing the attention check, yielding 347 data points. The average age of the participants was 25.8 years (SD=7.06). The minimum age was 18, and the maximum was 65. 164 (47.3%) participants identified as female, and 175 (50.4%) identified as male. Four (1.2%) participants identified as non-binary and one (0.3%) identified as genderqueer. Three (0.9%) participants declined to identify their gender. Participants were asked to self-identify in terms of race and ethnicity. The categories with more than 5 participants are White/Caucasian (192, 55.3%), Black/African (72, 20.7%), and Latinx/Hispanic (31, 8.9%). All other racial or ethnic identities comprised less than 1.7% (6 participants). Each participant was paid USD $4.00 for their time.

## 5 Analysis

## 5.1 Data Matrix

All naturalness scores were aggregated into a $100 \times 7$ matrix, where each of the 100 rows represented a different video excerpt shown to participants, and each of the 7 columns represented a different referring form. For example, the first cell in the matrix contained the average naturalness scores for the utterance following cutpoint 0 in video 0, when the referring expression in that utterance was replaced with "it". The remainder of the first row contained the average naturalness scores for the other possible referring forms used

following video 0 cutpoint 0. This data matrix was then used to evaluate a set of Referring Form Selection models, as described in the next section.

## 5.2 Models

To test our hypothesis, we compared five models using this data matrix. For each model, we considered each row in the data matrix, and identified which referring form the model would have predicted in the referring context encoded by that row. We then extracted the naturalness score from the column associated with that prediction. This produced a set of 100 naturalness scores for each model.

The five models we compared were (1) a *Random* baseline; (2) a *Definite Description* baseline; (3) a *Human* baseline; (4) Han et al. (2022)'s cognitive status-informed model where utterance-level temporal distance is used instead of object-level temporal distance; (5) A modified model trained with the physical distance being the furthest. This is in line with the data that the cognitive status model was trained on. Ideally, reference-level temporal distance would be utilized for both the CS model and the RF model, as it is more accurate on an object-mention-per-object-mention basis, but this would create conflict between the two models.

### 5.2.1 Random Model

The Random Model served as our first baseline. Under this model, a referring form was selected at random: for each row in our data matrix, the naturalness score from a random column was used.
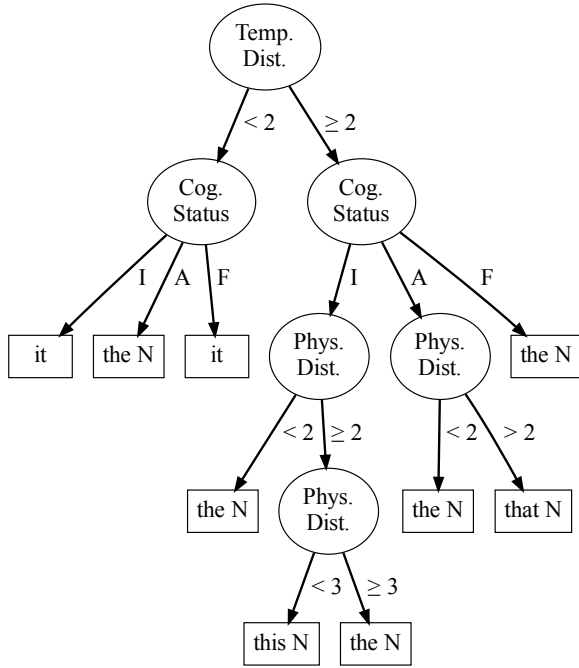
Figure 4: The visualization of the decision tree model used in the HanRFS condition, adapted from Han et al. (2022) with temporal distance at utterance level.

### 5.2.2 Definite Description Model

The Definite Description Model served as our second baseline. Under this model, *The N'* was always used: for each row in our data matrix, the naturalness score from the *The N'* column was extracted.

### 5.2.3 Human Model

The Human Model serves as our third baseline. Under this model, the referring form actually used by participants in the dataset was used. For example, if for a particular cutpoint in a particular video, the utterance following the cutpoint actually contained "that tower", then the naturalness score from the *That N'* column was extracted.

### 5.2.4 Cognitive Status-Informed Model (HanRFS)

For our fourth baseline, we used the cognitive status-informed model presented by Han et al. (2022) (under MIT licence). This model is a decision tree based machine learning model, shown in Figure 4, and uses three features: cognitive status, temporal distance, and physical distance.

*Cognitive status* was predicted using Pal et al. (2020)'s Bayesian cognitive status engine, which itself makes predictions based on an object's linguistic status and previously inferred cognitive status. That is, for a referring form to be predicted

| Far Left Quad: 6 | On Vertical Line between Far Quads: 5 | Far Right Quad: 6 |
| On Left Horizontal Line: 4 | At Line Intersection: 3 | On Right Horizontal Line: 4 |
| Near Left Quad: 2 | On Vertical Line between Near Quads: 1 | Near Right Quad: 2 |

Table 1: Codes for physical distance.

at cutpoint $t$, we fed each referring form from cutpoints $0 \ldots t-1$ (if any) to Pal's cognitive status engine. This produced a distribution over cognitive status that the target referent should have at time $t$. We then used the most likely cognitive status from this distribution as the feature passed to Han et al. (2022)'s decision tree.

*Temporal distance* was calculated as recency of mention: a target referent's temporal distance was calculated as the number of utterances since the utterance where the object was mentioned.

*Physical distance* was calculated in terms of qualitative distance-to-object. Han et al. (2022)'s original model was trained in a tabletop environment, and as such, they operationalized physical distance by assigning a set of distance scores 1-6 to each area in a $3 \times 3$ grid on the tabletop. We elected to do the same, breaking the task environment shown in the video into a $3 \times 3$ grid, and assigning a distance score 1-6 to each quadrant as shown in Table 1.

### 5.2.5 Modified Cognitive Status-Informed Model (HanRFS-RD)

Since the location of objects could have an impact on the choice of referring form and their perceived naturalness, we decided to include another baseline, HanRFS-RD (Remapped Distances), in which all physical distances were set to the furthest possible value (6), since all objects in Bennett et al. (2017)'s environment were further than any object used to train the decision tree model.

### 5.3 Data Analysis

To compare the predictions made by each of our five models, we used the Bayesian statistical framework (Wagenmakers et al., 2018), given its capability to quantify evidence both for and against a hypothesis, compared to the Frequentist approach. Specifically, we used JASP 0.17.1 (JASP Team, 2022) to run Bayesian statistical tests.

One important concept to understand in the

Table 2: Referring form distribution across conditions

| | *it* | *the* $\langle N' \rangle$ | $\langle N' \rangle$ | *this* | *that* | *this* $\langle N' \rangle$ | *that* $\langle N' \rangle$ |
|---|---|---|---|---|---|---|---|
| Original | 0.12 | 0.54 | 0.31 | 0.00 | 0.00 | 0.01 | 0.02 |
| Random | 0.13 | 0.16 | 0.08 | 0.14 | 0.13 | 0.16 | 0.21 |
| The $\langle N' \rangle$ | 0.00 | 1.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| Human | 0.28 | 0.29 | 0.29 | 0.01 | 0.02 | 0.04 | 0.07 |
| HanRFS | 0.30 | 0.32[†] | 0.00[†] | 0.00 | 0.00 | 0.27 | 0.10 |
| HanRFS-RD | 0.30 | 0.60[†] | 0.00[†] | 0.00 | 0.00 | 0.00 | 0.10 |

[†]The two models by Han et al. (2022) merged $\langle N \rangle$ with *the* $\langle N \rangle$
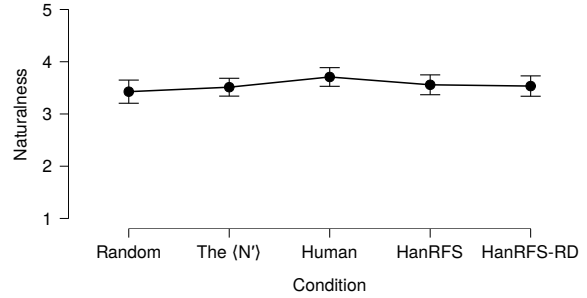


Figure 5: Mean naturalness ratings. Error bars show 95% credible intervals. Results favor no difference across conditions.
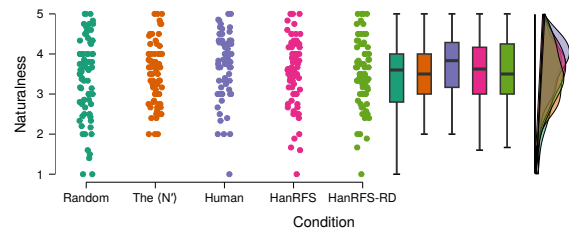


Figure 6: A raincloud plot (Allen et al., 2019) for naturalness ratings, combining a cloud of points with a box plot and a one-sided violin plot.

Bayesian approach is the Bayes factor (BF), defined as the ratio of the likelihood of given data being observed under each of two competing hypotheses, $\mathcal{H}_1$ and $\mathcal{H}_0$. For example, a Bayes Factor of $BF_{10}=5$ indicates a favor of $\mathcal{H}_1$ that the data are five times more likely under $\mathcal{H}_1$ than under $\mathcal{H}_0$.

To help the decision-making process, we used the widely-accepted classification scheme (Lee and Wagenmakers, 2014). For evidence favoring $\mathcal{H}_1$, Bayes factor values are categorized as anecdotal ($BF_{10} \in (1,3]$; inconclusive), weak ($BF_{10} \in (3,10]$), moderate ($BF_{10} \in (10,30]$), strong ($BF_{10} \in (30,100]$), extreme ($BF_{10} \in (100,\infty)$). When evidence favors $\mathcal{H}_0$, these thresholds are inverted, and we can use $BF_{01}$ for easier interpretation (Note the subscript is 01 rather than 10). For example, $BF_{10} = 1/5 = 0.2$ can be expressed as $BF_{01} = 5$.

## 6 Results

### 6.1 Model comparisons

Table 2 shows the distribution of the referring forms across the five conditions. The distribution of random referring forms, by definition, roughly follows a uniform distribution, with deviance due only to sampling noise. The $\langle N \rangle$ condition contains only itself. For the model condition, it predicts 30% of *it*, 31% of *the* $\langle N \rangle$, 27% of *this*, and 10% of *that* $\langle N \rangle$. On the other hand, for the fixed physical distance model, 30% of *it*, 60% of *the* $\langle N \rangle$, and 10% of *that* $\langle N \rangle$ was predicted, in line with the expected changes by increasing physical distance to its maximum value. Note that Han et al. (2022) took a descriptivist view (Frege, 1892; Russell, 2001; Nelson, 2002) and merged bare nouns ($\langle N \rangle$) with definite nouns (*the* $\langle N \rangle$).

### 6.2 Naturalness in Referring Form Selection

As seen from Figure 5, the mean naturalness scores are approximately the same in all five conditions and, surprisingly, the actual referring forms were only rated slightly higher: Random (M=3.427,

SD=0.977), The $\langle N' \rangle$ (M=3.513, SD=0.752), HanRFS (M=3.559, SD=0.837), HanRFS-RD (M=3.534, SD=0.859), and Human (M=3.621, SD=0.814). Out of all of these, the Human baseline performed marginally better (M=3.621), followed by HanRFS (M=3.559), HanRFS-RD (M=3.534), The $\langle N' \rangle$ (M=3.513) and Random (M=3.427); however, due to the SD value being greater than 0.75 for all models, there is no statistically significant difference between them. Figure 6 shows the raw data points with boxplots and distribution estimation.

To see whether there is a difference, we conducted a Bayesian one-way analysis of variance (ANOVA) (Rouder et al., 2012) on the naturalness data. This analysis revealed strong evidence against the effect of different referring form selections and the actual referring forms ($BF_{01} = 28.358$), i.e., favoring $\mathcal{H}_0$. This means that the data are around 28.358 times more likely under models that did not include an effect than under those that did. Thus, the hypothesis is not supported: the referring forms predicted by the cognitive status-informed model were not perceived as more natural.

## 7 Discussion

### 7.1 Naturalness is in the Eye of the Beholder

We hypothesized that the cognitive status-informed RFS model's predictions would be more natural. However, results showed that all RFS models are equally natural, with extremely high variability in perceived naturality even for human-generated referring expressions. There are a variety of possible explanations for these observations.

First, humans may regularly generate unnatural sounding referring expressions. If so, human-level naturality may merely be a "low bar" that NLG research should seek to surpass.

Second, humans may vary dramatically in their perceptions of what is "natural". If so, human judgment may be a poor way to assess referring form naturality.

Third, our experimental paradigm may have been unsuccessful in measuring the naturality of referring forms on their own. Referring forms are always used in the context of a larger utterance, which itself may be viewed as natural or unnatural. To mitigate this concern, we specifically asked participants about the naturalness of the referring forms used in the utterance they were shown. However, it is possible that participants either did not follow these instructions, or were simply unable to adjudicate the naturalness of these forms without considering the broader context of their use. For example, in the utterance "And while you're there can you knock over **the blue tower**", participants rated the usage of **blue tower** ($\langle N \rangle$) as very natural. The separation of the naturalness of the referring form from its context of use is remarkably challenging because, without context, using concise referring forms becomes no longer useful.

Finally, our results may be due to the global population reflected in our sample. While participants were required to be fluent in English, most participants indicated that English was their second language. This may have led to significant variation in our naturalness ratings.

**Takeaway 1**: Future work needs to better separate the perceived naturalness of a referring form from its context, such as dialog. To confirm this, one may need to measure the naturalness of the context as a controlling factor, or may need to be particularly aggressive about reminding participants that they are rating only the referring form itself.

### 7.2 How Far is Far?

We included the modified HanRFS-RD model with remapped distances because of the differences between what is considered "close" and "far" in our analyzed dataset versus the dataset on which Han et al. (2022)'s model was trained on. This raises a larger question, however, of how to model referent distances in a task-agnostic way. Physical distance is clearly an important factor, and is known to play a role in differentiating referring forms like "this" vs "that", as well as differentiating the use of abstract versus precise deictic gestures (Stogsdill et al., 2021). Yet what is considered near versus far is highly task dependent, depending not only on the overall size of the space, but also on the physical affordances and explorability of the space. For example, Han et al. (2022)'s model was created in a space that was smaller than Bennett et al. (2017). But moreover, while in Bennett et al. (2017)'s experiment objects were out of immediate reach of the participants, in Han et al. (2022)'s experiment, objects were reachable without walking around, i.e., no farther than 60cm (2 feet) away, and in fact were touched and manipulated by participants. In other task environments, other features may also become relevant. In large-scale open-world environments, for example, many referents are non-visible (or may not even be known to exist) when they are referred to.

**Takeaway 2**: Future work needs to understand how referring form selection models can encode physical distance features in a way that is agnostic of, or relative to, the size of a task environment; should consider inclusion of a suite of distance features sensitive to different types of task environments; and should consider features related to but distinct from distance, like reachability, manipulability, and visibility.

## 8 Conclusions

To go beyond the focus on goodness-of-fit in cognitive status-informed computational referring form selection model evaluation, we conducted a human-subjects study to explore the naturalness ratings of the predictions. Surprisingly, results did not reveal an improvement in naturalness over random baselines, and in fact suggest that human perceptions of even *human-generated* referring forms are incredibly varied and not significantly different from those random baselines. Our results suggest several directions for future work, and new technical and

methodological considerations that must be made.

## Acknowledgements

## References

Micah Allen, Davide Poggiali, Kirstie Whitaker, Tom Rhys Marshall, and Rogier A Kievit. 2019. Raincloud plots: a multi-platform tool for robust data visualization. *Wellcome open research*, 4.

Jennifer E Arnold and Sandra A Zerkle. 2019. Why do people produce pronouns? pragmatic selection vs. rational models. *Language, Cognition and Neuroscience*, 34(9):1152–1175.

Matthew Aylett and Alice Turk. 2004. The smooth signal redundancy hypothesis: A functional explanation for relationships between redundancy, prosodic prominence, and duration in spontaneous speech. *Language and speech*, 47(1):31–56.

Maxwell Bennett, Tom Williams, Daria Thames, and Matthias Scheutz. 2017. Differences in interaction patterns and perception for teleoperated and autonomous humanoid robots. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 6589–6594. IEEE.

Susan E Brennan, Marilyn W Friedman, and Carl Pollard. 1987. A centering approach to pronouns. In *25th Annual Meeting of the Association for Computational Linguistics*, pages 155–162.

Maya Cakmak and Andrea L Thomaz. 2012. Designing robot learners that ask good questions. In *2012 7th ACM/IEEE International Conference on Human-Robot Interaction*, pages 17–24. IEEE.

Charles B Callaway and James C Lester. 2002. Pronominalization in generated discourse and dialogue. In *ACL*.

Guanyi Chen, Fahime Same, and Kees van Deemter. 2021. What can neural referential form selectors learn? In *Proceedings of the 14th International Conference on Natural Language Generation*, pages 154–166.

Michael C Frank and Noah D Goodman. 2012. Predicting pragmatic reasoning in language games. *Science*, 336(6084):998–998.

Gottlob Frege. 1892. Über sinn und bedeutung. *Zeitschrift für Philosophie und philosophische Kritik*, 100:25–50.

Albert Gatt, Emiel Krahmer, Kees Van Deemter, and Roger Van Gompel. 2014. Models and empirical data for the production of referring expressions. *Lang., Cognition and Neuroscience*, 29(8):899–911.

Niyu Ge, John Hale, and Eugene Charniak. 1998. A statistical approach to anaphora resolution. In *Workshop on Very Large Corpora*.

Felix Gervits, Gordon Briggs, Antonio Roque, Genki A Kadomatsu, Dean Thurston, Matthias Scheutz, and Matthew Marge. 2021. Decision-theoretic question generation for situated reference resolution: An empirical study and computational model. In *Proceedings of the 2021 International Conference on Multimodal Interaction*, pages 150–158.

David A Grant. 1948. The latin square principle in the design and analysis of psychological experiments. *Psychological bulletin*, 45(5):427.

Barbara J Grosz, Aravind K Joshi, and Scott Weinstein. 1995. Centering: A framework for modelling the local coherence of discourse. *Computational Linguistics*.

Barbara J Grosz and Candace L Sidner. 1986. Attention, intentions, and the structure of discourse. *Computational linguistics*, 12(3):175–204.

André Grüning and Andrej A Kibrik. 2005. Modelling referential choice in discourse: A cognitive calculative approach and a neural network approach. *Anaphora processing: Linguistic, cognitive and computational modelling*, 263:163.

Jeanette K Gundel, Mamadou Bassene, Bryan Gordon, Linda Humnick, and Amel Khalfaoui. 2010. Testing predictions of the givenness hierarchy framework: A crosslinguistic investigation. *Journal of Pragmatics*, 42(7):1770–1785.

Jeanette K Gundel, Nancy Hedberg, and Ron Zacharski. 1993. Cognitive status and the form of referring expressions in discourse. *Language*, pages 274–307.

Zhao Han, Polina Rygina, and Thomas Williams. 2022. Evaluating referring form selection models in partially-known environments. In *Proceedings of the 15th International Conference on Natural Language Generation*, pages 1–14.

Zhao Han and Tom Williams. 2023. Evaluating cognitive status-informed referring form selection for human-robot interactions. In *2023 Annual Meeting of the Cognitive Science Society (CogSci)*.

Ryan Blake Jackson and Tom Williams. 2022. Enabling morally sensitive robotic clarification requests. *ACM Transactions on Human-Robot Interaction (THRI)*, 11(2):1–18.

JASP Team. 2022. JASP (Version 0.16.4)[Computer software].

Rodger Kibble and Richard Power. 2004. Optimizing referential coherence in text generation. *Comp. Ling.*

Andrej A Kibrik. 2011. *Reference in discourse*. OUP.

Andrej A Kibrik, Mariya V Khudyakova, Grigory B Dobrov, Anastasia Linnik, and Dmitrij A Zalmanov. 2016. Referential choice: Predictability and its limits. *Frontiers in psychology*, 7:1429.

Emiel Krahmer and Kees Van Deemter. 2012. Computational generation of referring expressions: A survey. *Computational Linguistics*, 38(1):173–218.

Olga Krasavina and Christian Chiarcos. 2007. Pocos-potsdam coreference scheme. In *Linguistic Annotation Workshop*.

Michael D Lee and Eric-Jan Wagenmakers. 2014. *Bayesian cognitive modeling: A practical course*. Cambridge university press.

William C Mann, Christian MIM Matthiessen, and Sandra A Thompson. 1989. Rhetorical structure theory and text analysis. *NASA STI/Recon Technical Report N*, 90:26733.

Kathleen F McCoy and Michael Strube. 1999. Generating anaphoric expressions: pronoun or definite description? In *The Relation of Discourse/Dialogue Structure and Reference*.

Michael Nelson. 2002. Descriptivism defended. *Noûs*, 36(3):408–435.

Poulomi Pal, Grace Clark, and Tom Williams. 2021. Givenness hierarchy theoretic referential choice in situated contexts. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 43.

Poulomi Pal, Lixiao Zhu, Andrea Golden-Lasher, Akshay Swaminathan, and Tom Williams. 2020. Givenness hierarchy theoretic cognitive status filtering. In *Annual Meeting of the Cognitive Science Society*.

Stefan Palan and Christian Schitter. 2018. Prolific. ac—a subject pool for online experiments. *Journal of Behavioral and Experimental Finance*, 17:22–27.

Massimo Poesio, Rosemary Stevenson, Barbara Di Eugenio, and Janet Hitzeman. 2004. Centering: A parametric theory and its instantiations. *Computational linguistics*, 30(3):309–363.

Jeffrey N Rouder, Richard D Morey, Paul L Speckman, and Jordan M Province. 2012. Default bayes factors for anova designs. *Journal of mathematical psychology*, 56(5):356–374.

Bertrand Russell. 2001. *The problems of philosophy*. OUP Oxford.

Fahime Same and Kees van Deemter. 2020. A linguistic perspective on reference: Choosing a feature set for generating referring expressions in context. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4575–4586.

Kevin Spevak, Zhao Han, Tom Williams, and Neil T Dantam. 2022. Givenness hierarchy informed optimal document planning for situated human-robot interaction. In *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*.

Adam Stogsdill, Grace Clark, Aly Ranucci, Thao Phung, and Tom Williams. 2021. Is it pointless? modeling and evaluation of category transitions of spatial gestures. In *Companion of the 2021 ACM/IEEE International Conference on Human-Robot Interaction*, pages 392–396.

Stefanie Tellex, Pratiksha Thakerll, Robin Deitsl, Dimitar Simeonovl, Thomas Kollar, and Nicholas Royl. 2013. Toward information theoretic human-robot dialog. *Robotics*, page 409.

Kees Van Deemter. 2016. *Computational models of referring: a study in cognitive science*. MIT Press.

Kees Van Deemter, Albert Gatt, Roger PG Van Gompel, and Emiel Krahmer. 2012. Toward a computational psycholinguistics of reference production. *Topics in cognitive science*, 4(2):166–183.

Eric-Jan Wagenmakers, Maarten Marsman, Tahira Jamil, Alexander Ly, Josine Verhagen, Jonathon Love, Ravi Selker, Quentin F Gronau, Martin Šmíra, Sacha Epskamp, et al. 2018. Bayesian inference for psychology. Part I: Theoretical advantages and practical ramifications. *Psychonomic bulletin & review*, 25(1):35–57.

Tom Williams, Gordon Briggs, Bradley Oosterveld, and Matthias Scheutz. 2015. Going beyond literal command-based instructions: Extending robotic natural language interaction capabilities. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*.

# System-Initiated Transitions from Chit-Chat to Task-Oriented Dialogues with Transition Info Extractor and Transition Sentence Generator

**Ye Liu**[1,3]**, Stefan Ultes**[2]**, Wolfgang Minker**[3] and **Wolfgang Maier**[1]

[1]Mercedes-Benz AG, Sindelfingen, Germany

`{ye.y.liu, wolfgang.mw.maier}@mercedes-benz.com`

[2]University of Bamberg, Bamberg, Germany

`stefan.ultes@uni-bamberg.de`

[3]Ulm University, Ulm, Germany

`{ye.liu, wolfgang.minker}@uni-ulm.de`

## Abstract

In this work, we study dialogue scenarios that start from chit-chat but eventually switch to task-related services, and investigate how a unified dialogue model, which can engage in both chit-chat and task-oriented dialogues, takes the initiative during the dialogue mode transition from chit-chat to task-oriented in a coherent and cooperative manner. We firstly build a *transition info extractor* (TIE) that keeps track of the preceding chit-chat interaction and detects the potential user intention to switch to a task-oriented service. Meanwhile, in the unified model, a *transition sentence generator* (TSG) is extended through efficient Adapter tuning and transition prompt learning. When the TIE successfully finds task-related information from the preceding chit-chat, such as a transition domain ("train" in Figure 1), then the TSG is activated automatically in the unified model to initiate this transition by generating a transition sentence under the guidance of transition information extracted by TIE. The experimental results show promising performance regarding the proactive transitions. We achieve an additional large improvement on TIE model by utilizing Conditional Random Fields (CRF). The TSG can flexibly generate transition sentences while maintaining the unified capabilities of normal chit-chat and task-oriented response generation.

## 1 Introduction

Spoken dialogue systems (SDSs) have usually been developed targeting only one out of two different categories, task-oriented or chit-chat (aka open-domain). The former focuses on achieving functional goals and the latter aims at creating engaging social conversations without special goals. In recent years, several previous works (Lin et al., 2021; Zhao et al., 2021; Young et al., 2022) have studied unified conversational models that can engage in both chit-chat and task-oriented dialogue. However, the system-initiated transitions that emerge during

switchover between these two dialogue modes have rarely been explored. Especially when a user chats casually with the dialogue system, but implicitly expresses a need for a specific task-related service, it is desired that the dialogue system is able to capture this hidden information and proactively ask the user if they require such a task-oriented service (like booking a train ticket in Figure 1). It has been proven to be beneficial for commercial SDSs to proactively offer or sell their task-related services (Chiu et al., 2022; Liu et al., 2023). Furthermore, these transitions smoothly initiated by the dialogue system are regarded as a proactive feature (Nothdurft et al., 2015) and can greatly improve the user interaction experience (Liu et al., 2022).

The goal of this work is to develop the initiative capabilities of a unified conversational model that is capable of detecting the implicit user intention of using some task-related services, even if they are talking casually, and to proactively bridge the connection from chit-chat to task-oriented dialogue through generating a transition sentence (red in Figure 1). As the dialogue example in Figure 1 shows, the original response at the transition turn is only "I see". If the agent can anticipate in advance that the user wants to visit the "London Kings Cross" through the preceding chit-chat, it can then proactively establish a connection with the task-oriented "train" service that the user needs by saying "If you want, I can look for a train to London Kings Cross for you.".

To enable the initiative capabilities in a unified model, the main contributions of this paper are as follows:

1. To detect the hidden task-related transition domain/slot/value entities from the preceding chit-chat, we propose the transition info extractor (TIE) to keep track of preceding chit-chat dialogue through leveraging natural language understanding (NLU) technology (Chen et al., 2019).
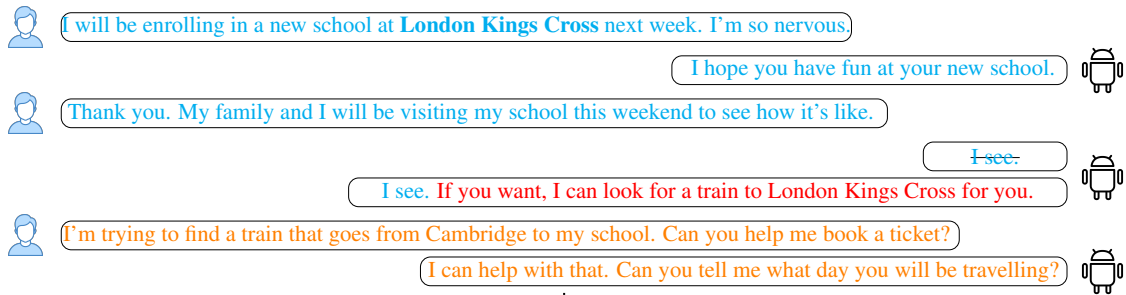
279

Figure 1: A Prepended FusedChat dialogue with an augmented transition sentence (red) for the proactive transition from chit-chat to task-oriented. The blue and orange represents chit-chat and task-oriented interaction respectively. Compared with the original chit-chat (crossed out) response at the transition turn, the transition sentence (red) can enable the dialogue system to proactively switch to task-oriented services.

2. We artificially augment 215 dialogues with a domain guided transition sentence and a domain-slot-value guided transition sentence respectively. We then collect transition sentence templates for different domains and different domain-slot pairs from these human augmented dialogues. The transition sentence templates are further utilized to annotate the remaining unannotated dialogues.

3. We leverage transition prompt learning (Li et al., 2022) and Adapter tuning (Lin et al., 2021) to efficiently extend the transition sentence generation (TSG) in a unified NLG model with the augmented dialogues.

The overall architecture flow of this work is shown in Figure 2. When the TIE successfully extracts the transition information from the preceding chit-chat, the TSG in the unified NLG is activated to generate a transition sentence besides the normal response to proactively guide this switch. The combined flow is highlighted in red. Otherwise, the TIE continually tracks the chit-chat, and unified NLG works as usual to generate a normal chit-chat or task-oriented response without (w/o) a transition sentence.

The remainder of this paper is structured as follows: Section 2 shows related work of our research. Section 3 presents the transition sentence augmentation and templates for the TSG training. Section 4 introduces the proposed TIE model for detecting the task-related transition information from the preceding chit-chat interaction. Section 5 presents the unified NLG extended with TSG through transition prompt and Adapter tuning. Section 6 elaborates on the performance evaluation of this work. Section 7 concludes this work and outlines future research.

## 2   Related Works

NLU is generally a crucial component in task-oriented SDSs and responsible for parsing an utterance into a semantic frame to identify the user's intention (De Mori et al., 2008). With the development of deep learning methods, RNN, CNN, as well as their variations or combinations have been widely for the NLU task (Yao et al., 2013; Mesnil et al., 2014; Yao et al., 2014; Liu and Lane, 2016). Wang et al. (2018) proposed a attention-based encoder-decoder, CNN-BLSTM, for joint intent detection and slot filling. Goo et al. (2018) proposed a slot gate that focused on capturing the relationship between slot and intent. Kenton and Toutanova (2019) and Xu et al. (2020) both used the pre-trained BERT for the joint intent classification and slot filling. The proposed TIE is inspired by NLU modeling.

Beyond that, Xu and Sarikaya (2013) and Ma and Hovy (2016) both utilized the traditional approach, Conditional Random Fields (CRF) (Sha and Pereira, 2003), for sequence labelling with the combination of LSTM and CNN. We also leverage the CRF technology to further improve the performance of the TIE model.

Shuster et al. (2020) introduced the dodecaDialogue task, to assemble important aspects of an engaging conversational agent into a single collection by leveraging 12 tasks. Adapter-Bot (Lin et al., 2021) utilized multiple adapter layers with the pre-trained DialoGPT model to activate new response skills and styles. Zhao et al. (2021) proposed a dialogue model for training chit-chat and task-oriented in a unified data schema, which both include belief states, representation of dataset results, and system acts. However, these models simply fuse chit-chat dialogue and task-oriented dialogue into
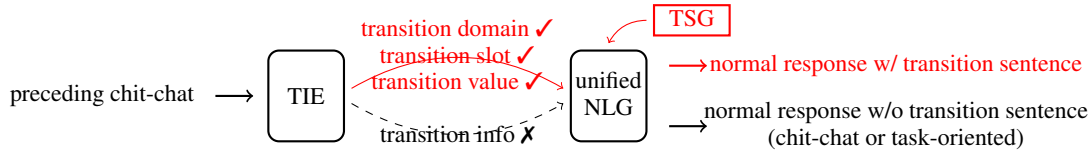
Figure 2: The overall Architecture flow for system-initiated transitions from chit-chat to task-oriented.

one model and do not consider the dependency between different types of dialogues in the multi-turn setting. In contrast, all dialogues in FusedChat released in Young et al. (2022) include both chit-chat and task-oriented turns, and treat them as parallel dialogue modes of equal importance. Chiu et al. (2022) proposed SalesBot and introduced the dialogue transitions from chit-chat to task-oriented. Liu et al. (2020) introduced the proactive transitions in conversational recommendation over multi-type dialogues. Liu et al. (2022) elaborated on three types of system-initiated transitions in a unified dialogue model and discussed the potential challenges respectively. Liu et al. (2023) proposed the system-initiated transitions between chit-chat and task-oriented dialogues, where the transitions from chit-chat to task-oriented and from task-oriented to chit-chat were treated equally. However, we mainly investigate the system-initiated transitions from chit-chat to task-oriented with the Prepended FusedChat dataset for this work.

## 3 Transition Sentence Augmentation and Templates

This section introduces the details of human augmentation of transition sentences and template collection for unannotated dialogues.

We mainly utilize the **Prepended** FusedChat (Young et al., 2022) dataset for initiative transitions from chit-chat to task-oriented in this work. FusedChat is a public available dataset, where human augmented open-domain dialogues are prepended and appended to the dialogues of the popular task-oriented dataset MultiWOZ (Budzianowski et al., 2018; Ye et al., 2021). In the Prepended FusedChat, each dialogue starts with chit-chat interaction and eventually switch to task-oriented requests. Table 1 shows the statistics of the Prepended FusedChat[1] used in this work. As a prepended FusedChat example shown in Figure 1, the user controls the switch

| data type | train | valid | test |
|---|---|---|---|
| dialogue size | 3255 | 474 | 331 |

Table 1: Statistics of Prepended FusedChat.

| domain | train | restaurant | attraction | taxi |
|---|---|---|---|---|
| number of templates | 95 | 56 | 45 | 17 |

Table 2: Statistics of transition sentence templates for different domains.

to task-oriented services.[2] However, our goal is to build a proactive dialogue system that can establish a smooth transition from chit-chat to task-oriented by itself.

To achieve this, we hire one master student with computational linguistics background to augment a domain guided transition sentence and a domain-slot-value guided transition sentence (red sentence in Figure 1) for 215 Prepended FusedChat dialogues respectively. The domain guided transition sentence must explicitly include the domain information. The domain-slot-value guided transition sentence must contain the specific value extracted from the preceding chit-chat dialogue aside from the domain, as the transition sentence in Figure 1, "If you want, I can look for a train to London Kings Cross for you." with "train" domain and "London Kings Cross" value.

After the human augmentation, we collect the templates for transition sentences in different domains and different domain-slot pairs from the augmented 215 dialogues respectively. For the domain-slot-value guided transition sentences, we use "[VALUE]" to replace the specific value to collect the domain-slot templates. Table 2 and Table 3 show template statistics for different domains and domain-slot pairs, respectively. Table 8 and Table 9 in the Appendix show some template examples of transition sentences in different domains and domain-slot pairs respectively. These templates are further used to randomly annotate the remaining

---

| domain | train | | | restaurant | | attraction | | taxi | |
|---|---|---|---|---|---|---|---|---|---|
| slot | day | destination | departure | food | name | type | name | destination | departure |
| number of templates | 22 | 40 | 35 | 45 | 11 | 30 | 15 | 9 | 8 |

Table 3: Statistics of transition sentence templates for different domain-slot pairs.

unannotated Prepended FusedChat dialogues. Then all Prepended FusedChat with augmented transition sentences can be used for training the extended TSG in the unified NLG.

## 4 Transition Info Extractor (TIE)

This section presents our TIE model that can detect potential user intention to switch to task-oriented services. As shown in Figure 3, TIE is adapted from pre-trained RoBERTa (Liu et al., 2019) and has three components, a transition domain classifier, a transition slot classifier and a slot filling layer. When the interaction starts from chit-chat, the TIE keeps track of the preceding chit-chat to predict the potential transition domain and slot, while extracting the specific value through slot filling. For instance, the transition domain-slot-value extracted in Figure 3, is "restaurant-food-Korean restaurant".

### 4.1 Joint RoBERTa for domain/slot classification and slot filling

We utilize the pre-trained RoBERTa (Liu et al., 2019) as the backbone TIE model for jointly predicting transition domain and corresponding slot, also extracting the specific value from the preceding chit-chat dialogue through slot filling task, as shown in Figure 3.

Given the preceding dialogue history until to the current user turn $\mathbf{x} = (x_1, x_2, ...x_n)$, the [CLS] token is inserted into the first place and [SEP] is inserted to split user utterances and system responses. The corresponding slot filling label is $\mathbf{y}^{sf} = (y_1, y_2, ...y_n)$ along with [CLS] and [SEP] tokens. The input $\mathbf{x}$ and slot filling label $\mathbf{y}$ are both padded to maximal length $N$ of the batch data. In addition, $y^d$ and $y^s$ are transition domain and slot label respectively. Let $\mathcal{D} = \{(\mathbf{x}, y^d, y^s, \mathbf{y}^{sf})\}_{m=1}^{M}$ be the dataset of size $M$ for joint RoBERTa training.

Adapted from the pre-trained RoBERTa, the final hidden states of the input are

$$h_{[\text{CLS}]}, h_{x_1}, h_{x_2}, ..., h_{x_n} = \text{RoBERTa}(\mathbf{x}) \quad (1)$$

Two classifier layers in Equation 2 are separately added on the output of [CLS] token, $h_{[\text{CLS}]}$, to pre-

dict transition domain and slot.

$$\hat{y}^d = softmax(\mathbf{W}^d Dropout(h_{[\text{CLS}]}) + \mathbf{b}^d)$$
$$\hat{y}^s = softmax(\mathbf{W}^s Dropout(h_{[\text{CLS}]}) + \mathbf{b}^s) \quad (2)$$

For the domain classifier, four different transition domains,[3] train, restaurant, attraction and taxi, are collected in the Prepended FusedChat. When no explicit user intention is detected in the preceding chit-chat, the domain classifier should recognise it as "UNK" to indicate that the current dialogue turn is not a good moment to switch to task-oriented. Hence, the domain classifier is a 5 classification task.

For the slot classifier, there are six slots,[4] namely day, destination, departure, food, name and type. Also along with "UNK", the slot classifier is a 7 classification task. Some slots are shared in different domains, e.g., "name" in restaurant and attraction domains (see Table 3).

For the slot filling task, the final hidden states in Equation 1 are fed into the slot filling (sf) layer in Equation 3 to classify over slot filling labels.

$$\hat{\mathbf{y}}_n^{sf} = softmax(\mathbf{W}^{sf} Dropout(\mathbf{h}_{\mathbf{x_n}}) + \mathbf{b}^{sf}); n \in 0...N \quad (3)$$

We use the IOB (In/Out/Begin) labelling format (Ramshaw and Marcus, 1999) for the slot filling labels. The dictionary of those labels is as follows and includes 22 tokens:

- 3 special tokens, "[PAD]", "[CLS]", "[SEP]", which are aligned with RoBERTa tokenizer.

- 9 domain-slot combinations in Table 3, but every domain-slot pair is extend with prefix "B-" and "I-". E.g. "B-restaurant-food" and "I-restaurant-food" in the Figure 3. When the specific value has more than one word, the first one is labelled with prefix "B-", the remaining with prefix "I-".

---

[3]"hotel" also exists in Prepended FusedChat as transition domain, but only in two dialogues. We delete those two dialogues to prevent the severe imbalance between different domains.

[4]Two dialogues have "pricerange" as transition slot under restaurant domain and one dialogue has "area' as transition slot under attraction domain. We also remove these dialogues in case of the imbalanced slots.
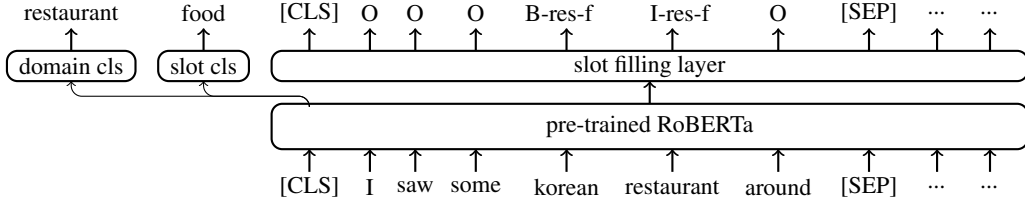
Figure 3: Architecture of the proposed TIE model that includes transition domain/slot classifier and slot filling task. The "B-res-f" and "I-res-f" is the abbreviation of the "B-restaurant-food" and "I-restaurant-food" respectively.

- The "O" is assigned to words not belonging to any specific value in sentences.

The $\mathbf{W}$ and $\mathbf{b}$ in Equation 2 and 3 are a trainable weight matrix and a bias vector. RoBERTa is jointly fine-tuned via minimizing the sum of cross-entropy loss of domain, slot classifier and slot filling task, as shown in Equation 4.

$$l_{joint\ RoBERTa} =$$

$$\sum_{M} (||\hat{y}^d - y^d||^2 + ||\hat{y}^s - y^s||^2 + \sum_{n=0}^{N} ||\hat{\mathbf{y}}_n^{sf} - \mathbf{y}^{sf}||^2)$$

(4)

### 4.2 Conditional Random Fields (CRF)

Beyond joint RoBERTa training for the transition domain/slot classification and slot filling tasks, we also use Conditional Random Fields (CRF) (Lafferty et al., 2001), to model the slot filling sequence jointly instead of decoding each slot filling label independently. CRF has been successfully used to exploit the dependencies within sequence labels corresponding to surrounding words and can highly improve the performance of slot filling task (Ma and Hovy, 2016). In this work, the dropout layer is applied before feeding RoBERTa outputs into CRF layer. The Viterbi algorithm is used for decoding.

We only utilize the preceding chit-chat part of Prepended FusedChat for joint RoBERTa training. To better analyse the proposed TIE model, three different TIE models are trained. As shown in Table 4, "RoBERTa w/o slot filling" only includes transition domain and slot classifiers; "joint RoBERTa" is jointly trained with domain, slot classifier and slot filling task together; and finally "joint RoBERTa + CRF" is our proposed final model, where the CRF is used for the slot filling task. All models are trained with two GPUs, the learning rate is $5e-5$ and batch size is 32. The best model of RoBERTa w/o slot filling is saved at epoch 5 with early stopping. The joint RoBERTa is saved at epoch 4 and joint RoBERTa + CRF at epoch 3.

## 5 Unified NLG extended with Transition Sentence Generator (TSG)

This section firstly introduces the unified NLG model that can reply to both chit-chat and task-oriented requests. Then we mainly elaborate on the TSG integrated in unified NLG through efficient Adapter tuning and transition prompt technologies. The extended NLG with TSG can generate a transition sentence given the transition information extracted by TIE to enable the system-initiated transition. The details of unified NLG and the extension with TSG are shown in Figure 4.

### 5.1 Unified NLG

We briefly presents the unified NLG model. By leveraging the entire FusedChat dataset (Young et al., 2022), where every dialogue includes interdependent chit-chat and task-oriented interaction, we tackle the unified generation problem through fine-tuning conditional GPT-2 (Radford et al., 2019). Given the FusedChat dataset $\mathcal{D}' = \{(u_g, d_g, r_g)_{g=1}^G, (u_l, r_l)_{l=1}^L\}$ with $G$ task-oriented samples (orange in Figure 4) and $L$ chit-chat samples (blue in Figure 4), the goal is to build a unified model parameterized by $\theta$ to be able to respond to both chit-chat and task-oriented requests,

$$p_\theta(r) = \begin{cases} \prod_{t=1}^T p_\theta(r_t|r_{<t}, u, d) & \text{if task-oriented} \\ \prod_{t=1}^T p_\theta(r_t|r_{<t}, u) & \text{if chit-chat} \end{cases}$$

(5)

where $r_{<t}$ indicates all tokens before $t$. The $u$ represents the dialogue context; $d$ means the dialogue actions only exist in task-oriented data and $r$ is the system response which includes $(r_1, ...r_t, ...)$ tokens with length $T$.

During the unified GPT-2 fine-tuning, we add [USER] and [SYSTEM] to the GPT-2 tokenizer to distinguish user utterances from system responses. At most three preceding dialogue turns are used as the dialogue context for response generation because of memory constraints. During training,
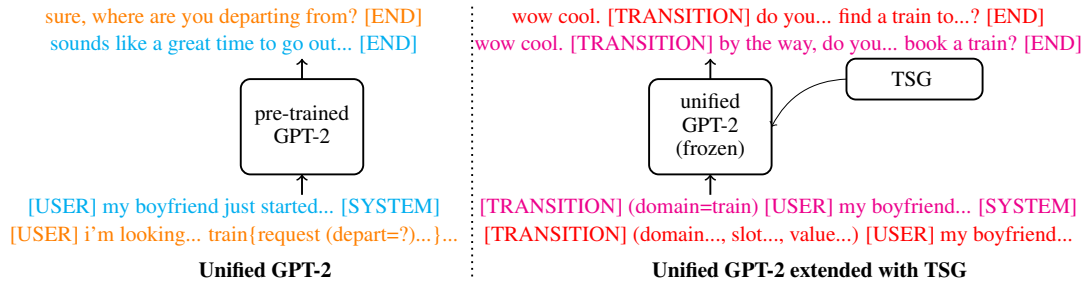
Figure 4: Architecture of unified GPT-2 and extended version with integrated TSG via Adapter tuning and transition prompt. In the unified GPT-2, the orange and blue represents the task-oriented and chit-chat example respectively. Two transition scenarios for each dialogue are used as training data for the TS Adapter tuning, one is only transition domain (magenta) as prompt, the other is transition domain-slot-value (red) as prompt.

the learning rate is $5e-5$, batch size is 20. The best model is saved at epoch 6 with early stopping. We mix top-K sampling and top-p (nucleus) sampling (Holtzman et al., 2019) for decoding. We apply top-K of 5 and top-p of 0.9 for chit-chat response generation and top-K of 10 and top-p of 0.5 for task-oriented response generation respectively.

## 5.2 Transition Prompt and Adapter Tuning

To enable the proactive capabilities, we integrate the efficient Adapter layers (Houlsby et al., 2019; Pfeiffer et al., 2021) into the unified GPT-2. Adapter tuning freezes the parameters of a pre-trained model and injects lightweight modules between layers (Le et al., 2021) to enable a new capability. Hence, the original performance of unified NLG for generating normal responses is retained without any loss. Meanwhile, the capability of generating transition sentences is extended through activating the newly added Adapter layers. To further explicitly control the transition sentence generation, the prompt learning (Liu et al., 2021; Li et al., 2022) is used. More precisely, when the TIE model successfully detects the user intention requiring a task-related service, the integrated Adapter layers are activated meanwhile the transition information extracted via TIE is converted into prompt input to generate a transition sentence to proactively establish the transition from chit-chat to task-oriented.

### 5.2.1 Transition Prompt

Prompt learning can efficiently adapt a given task to pre-trained models without modifying the structure of models (Lester et al., 2021). In this work, we only convert the task-related transition information extracted by TIE to the transition prompt, which is a part of the input for the generation model that explicitly guides the transition sentence generation.

We add a special token [TRANSITION] into the GPT-2 tokenizer and insert this token into the first place of the task-related transition prompt. Two different types of transition prompt are as follows:

1. When only the transition domain information is available, the prompt is like "[TRANSITION] ( domain = train ) ", where "train" is the extracted transition domain (magenta input in Figure 4).

2. When transition domain, slot and value are all extracted via TIE model, then the prompt is like "[TRANSITION] ( domain = train, slot = destination, value = Norwich ) ", where the transition domain is "train", slot is "destination" along with the value "Norwich" (red input in Figure 4).

The dialogue context is prepended with the transition prompt to be the input of the generation model. In addition, [TRANSITION] is also used to separate the transition sentence from normal response at transition turn (responses of magenta and red examples in Figure 4). Hence, the [TRANSITION] in prompt inputs is a signal for the generation model that it is a good moment to guide the transition to task-oriented service because TIE extracts task-related information, while the [TRANSITION] in generated responses is a signal to demonstrate that the NLG model is able to generate a transition sentence for proactive transition.

### 5.2.2 TSG through Adapter Tuning

We utilize the AdapterHub (Pfeiffer et al., 2020), which is a framework that can easily integrate Adapters into pre-trained Transformer-based models (Vaswani et al., 2017). The Houlsby Adapter (Houlsby et al., 2019) includes two bottleneck adapters in each transformer layer, one after the

284

multi-head attention sub-layer and the other after the feed-forward sub-layer. The Pfeiffer Adapter (Pfeiffer et al., 2021) only includes the adapter after the feed-forward sub-layer. Only 1% (Pfeiffer) and 2% (Houlsby) parameters are updated during Adapter tuning with frozen unified GPT-2. Hence, we can efficiently integrate the transition sentence generation into the unified GPT-2, while keeping the original capabilities of generating normal chit-chat and task-oriented responses by deactivating the Adapter layers.

Only the generation at the transition turn is utilized for the training of TSG. Every dialogue has two transition cases: One only consists of transition domain as prompt (magenta input in Figure 4) and the other consists of transition domain-slot-value as prompt (red input in Figure 4). We prepend the transition prompt before the preceding chit-chat context as input. The response includes a normal chit-chat response as well as a transition sentence separated with [TRANSITION] (the response of red and magenta examples in Figure 4). For the TSG, the Houlsby and Pfeiffer Adapters are both trained with the learning rate $5e-5$, batch size 20. The best models are both saved at epoch 16 (early stopping). We apply top-K of 5 and top-p of 0.9 for the response generation at the transition turn.

# 6 Results Comparison

This section evaluates this work and provides detailed performance comparison from different perspectives. We firstly evaluate different TIE models and different generation models separately with test Prepended FusedChat. Then we further evaluate the combined performance of the best TIE model and generation model only at transition turns.

## 6.1 TIE models

Table 4 shows the performance comparison of different TIE models. We use classification accuracy and weighted F1 score to evaluate the performance of transition domain and slot classifiers. Slot filling F1 (sf_f1) score is widely used to evaluate the slot filling task (Chen et al., 2019). In addition, we also use sentence-level slot filling accuracy (sen_sf_acc), which is the ratio of the number of dialogues correctly labelled slot filling to the total number of dialogues. The overall performance of the TIE model is evaluated using sentence-level semantic accuracy (semantic_acc) (Yu et al., 2010; Weld et al., 2021) which measures the proportion of

the correctly predicted triples of transition domain, slot, and extracted slot filling values (including "O" labels).

The performance comparison in Table 4 demonstrates that joint RoBERTa with CRF as the TIE model achieves the best performance over transition domain classifier, slot classifier and slot filling task. It is surprising that not only the slot filling task benefits from the CRF. The performance of transition domain and slot classifiers is improved in the multi-task learning as well.

## 6.2 Generation models

To evaluate generated chit-chat responses, Distinct-1 (Dis-1) and Distinct-2 (Dis-2) (Li et al., 2016) are used to measure the proportion of the distinct unigrams and bigrams in all the generated results to indicate diversity. To evaluate generated task-oriented responses, two *N*-gram matching metrics, BLEU (Papineni et al., 2002) and Meteor (Banerjee and Lavie, 2005) are used to evaluate the overall quality of task-oriented generations. In addition, a machine learned automatic metric, BERTScore (Zhang et al., 2019) is also utilized to evaluate task-oriented and transition sentence generations.

Beyond that, we propose several automatic metrics to evaluate transition sentence generations. *Transition accuracy* detects whether the generated response at transition turn includes the [TRANSITION] special token. With [TRANSITION], we can split the transition sentence from the normal response. This metric can measure high-level capability if the model can generate a transition sentence to proactively switch to a task-oriented service. *d accuracy* detects if the domain guided transition sentence includes the specific domain keyword. *d-v accuracy* detects if the transition domain-slot-value guided transition sentence includes specific domain and value keywords both. *d accuracy* and *d-v accuracy* can evaluate the capability of the proposed transition prompt for explicitly controlling transition sentence generation to a large extent.

We found that almost all generated transition sentences by TSG with TP are of high quality and include the extracted transition information (several cases are shown in Table 7 in Appendix), instead of generic transition responses like "Do you need anything else?" or "Do you need some help?".

To better understand the performance of our models, we also retrain the unified GPT-2 without

| | domain classifier | | slot classifier | | slot filling | | |
|---|---|---|---|---|---|---|---|
| | accuracy | weighted f1 | accuracy | weighted f1 | sen_sf_acc | sf_f1 | semantic_acc |
| RoBERTa w/o slot filling | 78.57% | 79.57% | 66.52% | 66.84% | – | – | – |
| joint RoBERTa | 82.41% | 82.92% | 71.86% | 73.84% | 68.02% | 48.64% | 61.94% |
| joint RoBERTa + CRF | **93.71%** | **94.15%** | **82.41%** | **82.30%** | **80.28%** | **61.82%** | **73.67%** |

Table 4: Performance of transition domain/slot classification and slot filling task in different TIE models.

| | | Chit-Chat | | Task-Oriented | | | domain TS | | | domain-slot-value TS | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Dis-1 | Dis-2 | BLUE | Meteor | BERTScore (F1) | BERTScore (F1) | *transition accuracy* | *d accuracy* | BERTScore (F1) | *transition accuracy* | *d-v accuracy* |
| unified GPT-2 | | **1.74%** | **12.70%** | **34.77%** | **55.65%** | **93.20%** | – | – | – | – | – | – |
| retrain | w/o TP | 1.67% | 11.41% | 32.86% | 53.52% | 92.91% | 88.82% | 98.25% | 58.19% | 89.29% | 98.97% | 30.15% |
| | w/ TP | 1.60% | 11.18% | 32.58% | 53.33% | 92.94% | 90.19% | 98.43% | 99.21% | 91.70% | 98.79% | 92.63% |
| TSG (Houlsby) | w/o TP | **1.74%** | **12.70%** | **34.77%** | **55.65%** | **93.20%** | 89.04% | 98.67% | 62.48% | 89.40% | **99.34%** | 27.19% |
| | w/ TP | **1.74%** | **12.70%** | **34.77%** | **55.65%** | **93.20%** | 90.28% | **99.40%** | 99.15% | **91.84%** | 99.21% | **96.80%** |
| TSG (Pfeiffer) | w/o TP | **1.74%** | **12.70%** | **34.77%** | **55.65%** | **93.20%** | 88.90% | 97.82% | 59.52% | 89.33% | 98.25% | 25.98% |
| | w/ TP | **1.74%** | **12.70%** | **34.77%** | **55.65%** | **93.20%** | 90.34% | 98.13% | **99.70%** | 91.83% | 98.43% | 96.62% |

Table 5: Performance of different NLG models, including unified GPT-2 and retrained without Adapter, extended with Houslby and Pfeiffer TSG separately, and all with transition prompt (w/ TP) and w/o TP respectively.

Adapter to enable its transition sentence generation (without TSG). From the comparison between the retrained model and unified GPT-2 in Table 5, we can see that the performance on chit-chat and task-oriented response generations has a loss, even though the retrained GPT-2 is still able to generate transition sentences. In contrast, our TSG extended in unified GPT-2 through Adapter tuning can retain the original capability for chit-chat and task-oriented generations, while maintaining a better performance on transition sentence generation.In addition, the retraining is not memory-efficient, while TSG only updates the Adapter parameters with frozen GPT-2.

To better assess the effects of our proposed transition prompt method, we retrain the model and extend TSG both along with the transition prompt (w/ TP) and without the transition prompt (w/o TP) respectively. Through the comparison between w/o TP and w/ TP in different models (highlighted in gray background in Table 5), the *d accuracy* and *d-v accuracy* metrics are highly improved with transition prompt guidance. This demonstrates that transition prompt can explicitly control the transition sentence generation. The performance comparison between Pfeiffer and Houlsby Adapter tuning has no big difference, however, the Pfeiffer Adapter uses only half of the trainable parameters, and is therefore the more effective choice for this work.

### 6.3 Combined TIE and generation model

To better reflect the overall performance of this work, we evaluate the combined TIE and generation models at transition turns, i.e., given the preceding chit-chat, the TIE model predicts transition domain/slot and extracts values, then this generated transition information by TIE is used as the transition prompt to guide transition sentence generation at the transition turn. Table 6 shows the combined performance of TIE and unified GPT-2 with Houlsby and Pfeiffer TSG, respectively.

Given the higher domain accuracy compared to slot accuracy, it is sensible to only use domain prediction as transition information to guide transition sentence generation when generated transition slot or extracted values are not reliable. This also validates our initial idea to propose two kinds of transition prompts. Regarding the lower slot accuracy, we found that the TIE model tends to confuse "destination" and "departure" under the "train" domain; over 60% of slot misjudged dialogues are in these cases. This would further affect the overall performance of the TIE model, which is shown by the semantic_acc metric.

Each Prepended FusedChat dialogue has only one turn for the transition from chit-chat to task-oriented. We directly define this turn as the transition turn, where the initiative dialogue model proactively switches to a task-oriented service through generating a transition sentence. Also, dialogue interactions could be more sophisticated in real life and it is difficult to accurately define the most appropriate moment to initiate a proactive transition. Furthermore, it gets more complicated if there are multiple transitions in one dialogue. A further, deeper investigation of appropriate moments for

| | | TIE (joint RoBERTa + CRF) | | | | Extended GPT-2 with TSG | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | domain cls | slot cls | slot filling | | domain TS | | | domain-slot-value TS | | |
| | | accuracy | accuracy | sf_f1 | semantic_acc | BERTScore (F1) | *transition accuracy* | *d accuracy* | BERTScore (F1) | *transition accuracy* | *d-v accuracy* |
| TSG (Houlsby) | w/ TP | | | | | 90.10% | 99.40% | 92.87% | 91.10% | 99.21% | 82.78% |
| TSG (Pfeiffer) | w/ TP | 93.35% | 65.56% | 64.71% | 50.15% | 90.08% | 98.49% | 93.53% | 91.25% | 98.37% | 83.02% |

Table 6: Overall performance of combined TIE and extended GPT-2 with TSG at transition turns.

a dialogue mode transition will be done in future work.

## 7 Conclusion

This work investigates the dialogue transition from chit-chat to task-oriented initiated by a dialogue agent. We build a TIE model adapted from pre-trained RoBERTa to keep track of the preceding chit-chat and predict transition domain, slot, while extracting the specific value from the chit-chat history via slot filling task. A unified generation model adapted from the pre-trained GPT-2 is built and extended its proactive capability for transition sentence generation through efficient Adapter tuning and transition prompt learning. Our proposed work shows promising performance both on transition information extraction and transition sentence generation. We will continue working on system-initiated transitions in other dialogue scenarios in the future.

## 8 Ethics Statement

This work develops proactive transitions from chit-chat to task-oriented dialogue in a unified dialogue system. Proactivity is always desired during the development of voice assistants. It can improve user interactive experience and serve users more efficiently. The dataset used in this work is public available and manually collected. Furthermore, our research is limited to a specific case, i.e, the user starts casual chat and eventually switches to a task-oriented service. However, more hidden challenges and ethics issues should be discussed further in the real scenarios. Would users prefer to be proactively served if the dialogue system successfully detects the user intention? Will they feel their privacy is violated if the dialogue system proactively provides task-related services? Such potential issues could be addressed by asking for user consent before providing the proactive interaction, which raises the additional question how many users would turn on such a feature from the start.

## References

Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72.

Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, Iñigo Casanueva, Stefan Ultes, Osman Ramadan, and Milica Gasic. 2018. Multiwoz-a large-scale multi-domain wizard-of-oz dataset for task-oriented dialogue modelling. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 5016–5026.

Qian Chen, Zhu Zhuo, and Wen Wang. 2019. Bert for joint intent classification and slot filling. *arXiv preprint arXiv:1902.10909*.

Ssu Chiu, Maolin Li, Yen-Ting Lin, and Yun-Nung Chen. 2022. Salesbot: Transitioning from chit-chat to task-oriented dialogues. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6143–6158.

Renato De Mori, Frédéric Bechet, Dilek Hakkani-Tur, Michael McTear, Giuseppe Riccardi, and Gokhan Tur. 2008. Spoken language understanding. *IEEE Signal Processing Magazine*, 25(3):50–58.

Chih-Wen Goo, Guang Gao, Yun-Kai Hsu, Chih-Li Huo, Tsung-Chieh Chen, Keng-Wei Hsu, and Yun-Nung Chen. 2018. Slot-gated modeling for joint slot filling and intent prediction. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 753–757.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration. In *International Conference on Learning Representations*.

Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*, pages 2790–2799. PMLR.

Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pages 4171–4186.

John Lafferty, Andrew McCallum, and Fernando CN Pereira. 2001. Conditional random fields: Probabilistic models for segmenting and labeling sequence data.

Hang Le, Juan Pino, Changhan Wang, Jiatao Gu, Didier Schwab, and Laurent Besacier. 2021. Lightweight adapter tuning for multilingual speech translation. In *The Joint Conference of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (ACL-IJCNLP 2021)*.

Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059.

Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and William B Dolan. 2016. A diversity-promoting objective function for neural conversation models. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 110–119.

Lei Li, Yongfeng Zhang, and Li Chen. 2022. Personalized prompt learning for explainable recommendation. *arXiv preprint arXiv:2202.07371*.

Zhaojiang Lin, Andrea Madotto, Yejin Bang, and Pascale Fung. 2021. The adapter-bot: All-in-one controllable conversational model. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 16081–16083.

Bing Liu and Ian Lane. 2016. Attention-based recurrent neural network models for joint intent detection and slot filling. *Interspeech 2016*, pages 685–689.

Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021. Pretrain, prompt, and predict: A systematic survey of prompting methods in natural language processing. *arXiv preprint arXiv:2107.13586*.

Ye Liu, Stefan Ultes, Wolfgang Minker, and Wolfgang Maier. 2023. Unified conversational models with system-initiated transitions between chitchat and task-oriented dialogues. *arXiv preprint arXiv:2307.01664*.

Ye Liu, Yung-Ching Yang, Wolfgang Maier, Wolfgang Minker, and Stefan Ultes. 2022. On system-initiated transitions in a unified natural language generation model for dialogue systems. In *Proceedings of the 26th Workshop on the Semantics and Pragmatics of Dialogue - Poster Abstracts*, Dublin, Ireland. SEM-DIAL.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.

Zeming Liu, Haifeng Wang, Zheng-Yu Niu, Hua Wu, Wanxiang Che, and Ting Liu. 2020. Towards conversational recommendation over multi-type dialogs. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1036–1049.

Xuezhe Ma and Eduard Hovy. 2016. End-to-end sequence labeling via bi-directional lstm-cnns-crf. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1064–1074.

Grégoire Mesnil, Yann Dauphin, Kaisheng Yao, Yoshua Bengio, Li Deng, Dilek Hakkani-Tur, Xiaodong He, Larry Heck, Gokhan Tur, Dong Yu, et al. 2014. Using recurrent neural networks for slot filling in spoken language understanding. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 23(3):530–539.

Florian Nothdurft, Stefan Ultes, and Wolfgang Minker. 2015. Finding appropriate interaction strategies for proactive dialogue systems—an open quest. In *Proc. of the 2nd European and the 5th Nordic Symposium on Multimodal Communication 2014*, pages 73–80. LiU Electronic Press.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. 2021. Adapterfusion: Non-destructive task composition for transfer learning. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 487–503.

Jonas Pfeiffer, Andreas Rücklé, Clifton Poth, Aishwarya Kamath, Ivan Vulić, Sebastian Ruder, Kyunghyun Cho, and Iryna Gurevych. 2020. Adapterhub: A framework for adapting transformers. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 46–54.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Lance A Ramshaw and Mitchell P Marcus. 1999. Text chunking using transformation-based learning. In *Natural language processing using very large corpora*, pages 157–176. Springer.

288

Fei Sha and Fernando Pereira. 2003. Shallow parsing with conditional random fields. In *Proceedings of the 2003 Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics*, pages 213–220.

Kurt Shuster, Da Ju, Stephen Roller, Emily Dinan, Y-Lan Boureau, and Jason Weston. 2020. The dialogue dodecathlon: Open-domain knowledge and image grounded conversational agents. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2453–2470.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NIPS*.

Yufan Wang, Li Tang, and Tingting He. 2018. Attention-based cnn-blstm networks for joint intent detection and slot filling. In *Chinese Computational Linguistics and Natural Language Processing Based on Naturally Annotated Big Data*, pages 250–261. Springer.

Henry Weld, Xiaoqi Huang, Siqu Long, Josiah Poon, and Soyeon Caren Han. 2021. A survey of joint intent detection and slot filling models in natural language understanding. *ACM Computing Surveys (CSUR)*.

Puyang Xu and Ruhi Sarikaya. 2013. Convolutional neural network based triangular crf for joint intent detection and slot filling. In *2013 ieee workshop on automatic speech recognition and understanding*, pages 78–83. IEEE.

Weijia Xu, Batool Haider, and Saab Mansour. 2020. End-to-end slot alignment and recognition for cross-lingual nlu. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5052–5063.

Kaisheng Yao, Baolin Peng, Yu Zhang, Dong Yu, Geoffrey Zweig, and Yangyang Shi. 2014. Spoken language understanding using long short-term memory neural networks. In *2014 IEEE Spoken Language Technology Workshop (SLT)*, pages 189–194. IEEE.

Kaisheng Yao, Geoffrey Zweig, Mei-Yuh Hwang, Yangyang Shi, and Dong Yu. 2013. Recurrent neural networks for language understanding. In *Interspeech*, pages 2524–2528.

Fanghua Ye, Jarana Manotumruksa, and Emine Yilmaz. 2021. Multiwoz 2.4: A multi-domain task-oriented dialogue dataset with essential annotation corrections to improve state tracking evaluation. *CoRR*, abs/2104.00773.

Tom Young, Frank Xing, Vlad Pandelea, Jinjie Ni, and Erik Cambria. 2022. Fusing task-oriented and open-domain dialogues in conversational agents.

Dong Yu, Shizhen Wang, and Li Deng. 2010. Sequential labeling using deep-structured conditional random fields. *IEEE Journal of Selected Topics in Signal Processing*, 4(6):965–973.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

Xinyan Zhao, Bin He, Yasheng Wang, Yitong Li, Fei Mi, Yajiao Liu, Xin Jiang, Qun Liu, and Huanhuan Chen. 2021. Unids: A unified dialogue system for chit-chat and task-oriented dialogues. *arXiv e-prints*, pages arXiv–2110.

# 9   Appendix

Table 7 shows several dialogue examples with generated transition sentences in this work. Table 8 and Table 9 show transition sentence templates for different domains and domain-slot pairs, respectively.

| | |
|---|---|
| dialogue context | **User:** A group of my friends from another university is coming to visit my campus!<br>**System:** What is the occasion?<br>**User:** My friends are actually working together to film a short movie and they would like shoot some footage here.<br>**System:** It seems your friends are working on quite an interesting project.<br>**User:** I am going to meetup with them this Saturday before bringing them to my university. |
| response at transition turn | *[d]:* That is a good opportunity to catch up with your friends. *By the way, I can help you with the **train**.*<br>*[d-v]:* That sounds like a great opportunity to discuss the film with your friends. *By the way, i am happy to look for a **train** on **Saturday** for you.* |
| dialogue context | **User:** My sister likes to watch musical plays.<br>**System:** Nice!<br>**User:** I am thinking of bringing her to a theatre to watch one soon since it will be her birthday next week. |
| response at transition turn | *[d]:* Sounds fun! *If you are looking for an **attraction**, I am happy to help.*<br>*[d-v]:* Sounds great! *If you are finding any **theatres**, our **attraction** service can help.* |
| dialogue context | **User:** I think I have had enough time in Cambridge, its a small town and you really get tired of seeing the same things over and over again .<br>**System:** How long have you been here for?<br>**User:** I have lived here for about 3 years now, I think I want a change of fresh air. |
| response at transition turn | *[d]:* That sounds great. *If you want, I can look for a **train** for you.*<br>*[d-v]:* I think change is always good. *if you are planning to travel, I am happy to look for a **train** leaving from **Cambridge** for you.* |
| dialogue context | **User:** Without the additional workload, I should being accompany my family on a visit now. Now i have to let my family visit Hobsons house first and meet my family there.<br>**System:** Sorry to hear that, hope you can finish the work early.<br>**User:** I have finished the work and left the company. I will try to meet my family by bus or taxi, whichever is faster. |
| response at transition turn | *[d]:* I am sure your family will have an enjoyable time there. *Shall I get a **taxi** for you getting there?*<br>*[d-v]:* I am sure your family will understand. , *By the way, if you want to book a taxi to **hobsons house**, feel free to use our **taxi** service.* |

Table 7: Several dialogue examples with transition sentence (highlighted in red) generated by the extended NLG with TSG. The *[d]* means only the transition domain as transition prompt and *[d-v]* means the transition domain-slot-value as transition prompt to guide the transition sentence generation. Transition domains and values present in transition sentences are highlighted in **bold**.

| domain | templates of the transition sentence |
|---|---|
| restaurant | I am happy to give recommendation on restaurants.<br>I can recommend some restaurants if you want.<br>Do you want my recommendation on the restaurants?<br>I can also provide you more information on this restaurant.<br>Maybe you would like to use our restaurant service to know more about it.<br>... |
| attraction | By the way, you can reach to our attraction service to know more about this place.<br>Besides, our attraction service provides various information.<br>I can recommend some attractions to you.<br>By the way, have you checked out our attraction service to know more about this place?<br>If you are finding any attraction, I am always happy to help.<br>... |
| train | Additionally I could help with looking for train tickets for you.<br>By the way, I can help you to find thee trains to get there.<br>Let me arrange the train for you.<br>Please refer to our train service if you need any help with the booking.<br>I am glad to give you more information on the train.<br>... |
| taxi | Do you need help with booking a taxi to get there?<br>Do you want me to look for a taxi for you?<br>Do you need a taxi afterwards?<br>Maybe you would like my help with the taxi?<br>If you need to get there soon, I can help you book a taxi.<br>... |

Table 8: Transition sentences templates for different domains.

291

| domain-slot | templates of the transition sentence |
| --- | --- |
| restaurant-food | I am happy to give recommendation on [VALUE] restaurants.<br>I can recommend some [VALUE] restaurants if you want.<br>You can find more information on [VALUE] restaurants in our restaurant service.<br>It's my pleasure to recommend some [VALUE] restaurants if you want.<br>... |
| restaurant-name | I can also provide you more information on this restaurant named [VALUE].<br>Maybe you would like to use our restaurant service to know more about [VALUE].<br>I will be more than pleasant to help with booking a table at the restaurant called [VALUE].<br>Feel free to ask for more information about this restaurant named [VALUE].<br>... |
| attraction-name | By the way, you can reach to our attraction service to know more about [VALUE].<br>Do you want to plan your trip to [VALUE] using our attraction service?<br>By the way, I can provide more attraction information on [VALUE].<br>Talking about attractions, do you need more information about [VALUE].<br>... |
| attraction-type | Besides, our attraction service provides various information on [VALUE].<br>If you are looking for attraction that has [VALUE] activities, i am happy to help you.<br>In our attraction service, you can find more information on visiting [VALUE]s.<br>... |
| train-day | Additionally I could help with looking for train on [VALUE] for you.<br>Let me arrange the train for [VALUE] for you.<br>If you want, you can use our service to book the train for [VALUE].<br>I would love to help you with the train tickets for [VALUE].<br>... |
| train-destination | By the way, I can help you to find the trains to [VALUE].<br>If you want, I can look for a train to [VALUE] for you.<br>Additionally, you can use our service to book a train to [VALUE].<br>... |
| train-departure | I think our service might be helpful in booking the train leaving from [VALUE].<br>I am happy to look for a train leaving from [VALUE] for you.<br>Shall I find you some train tickets departing from [VALUE].<br>... |
| taxi-departure | By the way, do you need help with booking a taxi departing from [VALUE]?<br>Do you want me to look for a taxi depart from [VALUE] for you?.<br>Will you need my help with the taxi leaving from [VALUE].<br>... |
| taxi-destination | Shall I get a taxi for you to get to [VALUE]?<br>By the way, if you need a taxi to [VALUE], please feel free to use our taxi service.<br>If you need a taxi to get to [VALUE], feel free to use our taxi service.<br>... |

Table 9: Transition sentences templates for different domain-slot pairs. The specific values in the human augmented transition sentences are replaced by the special [VALUE] token to collect the templates.

# HL Dataset: Visually-grounded Description of Scenes, Actions and Rationales

**Michele Cafagna**[1]    **Kees van Deemter**[2]    **Albert Gatt**[1,2]

[1]University of Malta, Institute of Linguistics and Language Technology
[2]Universiteit Utrecht, Information and Computing Sciences
michele.cafagna@um.edu.mt
{a.gatt, c.j.vandeemter}@uu.nl

## Abstract

Current captioning datasets focus on object-centric captions, describing the visible objects in the image, e.g. "people eating food in a park". Although these datasets are useful to evaluate the ability of Vision & Language models to recognize and describe visual content, they do not support controlled experiments involving model testing or fine-tuning, with more high-level captions, which humans find easy and natural to produce. For example, people often describe images based on the type of scene they depict ('people at a holiday resort') and the actions they perform ('people having a picnic'). Such descriptions draw on personal experience and commonsense assumptions. We present the High-Level Dataset [1]; a dataset extending 14997 images from the COCO dataset, aligned with a new set of 134,973 human annotated (high-level) captions collected along three axes: *scenes*, *actions* and *rationales*. We further extend this dataset with confidence scores collected from an independent set of readers, as well as a set of narrative captions generated synthetically, by combining each of the three axes. We describe this dataset and analyse it extensively. We also present baseline results for the High-Level Captioning task.

## 1 Introduction

Conceptual grounding broadly refers to the idea that symbols (e.g. language) are grounded in perception (Barsalou et al., 2008). Perceptually grounded communication is made possible by the fact that perceptual experiences are largely shared. However, individual experience can also license subjective inferences which inform not just what we express through language, but also what we choose to assume and leave unexpressed (Bisk et al., 2020).

Among the many modalities available in the perceptual spectrum, visual grounding has always been of primary interest as it provides a relatively straightforward way to link linguistic expressions to physical objects. Consistent with this claim, a glance at many widely used datasets and models in image captioning reveals a bias towards 'object-centric' descriptions, whereby models are trained on image-text pairs where the text consists of explicit mentions of objects visible in the scene. However, experience and perception also motivate other, non-object-centric ways of talking about the world, for example, when we talk about scenes, or when we describe actions or their underlying rationales. While such 'high-level' descriptions are also perceptually grounded, they incorporate world knowledge and subjective experience.

For example, the object-centric description in Table 1 certainly describes the visual content, though it is based mainly on the recognition of objects in the scene. By contrast, the three high-level captions (*scene, action, rationale*, from the HL-Dataset described below), provide three different perspectives of the scene among the many possible ones, which are triggered by expectations and assumptions based on subjective experience and world knowledge.

In this work, we tackle the issue of grounding high-level linguistic descriptions in the visual modality, proposing the High-Level (HL) Dataset: a resource for Vision and Language (V&L) modeling which aligns existing object-centric captions with human-collected high-level descriptions of images along three different axes: *scenes, actions* and *rationales*. The high-level captions capture the human interpretation of the scene which are complementary to object-centric captions used in current V&L datasets, e.g. in COCO (Lin et al., 2014). We take a step further, and we collect *confidence scores* from independent annotators, which serve to shed

---

| Image | Axis | Caption |
|---|---|---|
|  | scene | the picture is shot in a ski resort |
| | action | they are just relaxing after a round of skiing |
| | rationale | they want to have a good time together |
| | object-centric (COCO) | a woman and a boy sitting in the snow outside of a cabin. |

Table 1: Example of High-Level captions. It is shown one of the three captions available for the three axes collected: *scene, action, rationale*, combined with the object-centric captions from COCO.

light on the extent to which the high-level captions in the dataset correspond to widely-shared assumptions, or to idiosyncratic interpretations. Finally, we consider the task of generating captions that incorporate these different axes, yielding a more narrative-like description of images. Our contributions are:

- We present and release the HL Dataset, a new V&L resource, grounding high-level captions in images along three different axes and aligned with existing object-centric captions;

- We describe the collection protocol and provide an in-depth analysis of the data;

- We present baselines for the High-Level Captioning task and describe further potential uses for our data.

## 2 Related work

Hodosh et al. (2013), in their influential work, argue that image captioning is mostly interested in 'conceptual descriptions', which focus on what is actually in the image and differ from the so-called non-visual descriptions, which provide additional background information. This line of thought has been broadly followed in the field, resulting in datasets emphasizing object-centric content in V&L tasks involving text generation, like image captioning (Lin et al., 2014; Sharma et al., 2018; Agrawal et al., 2019) and visual question answering (Antol et al., 2015; Zhu et al., 2016).

For instance, in the instructions used to collect COCO (Lin et al., 2014), the annotators are explicitly asked to mention entities visible in the image. This is beneficial to enhance cross-modal interactions: Zhang et al. (2021) show that improving the visual backbone on object recognition tasks, improves the performance of visio-linguistic models in downstream tasks. Li et al. (2020) show that

using object labels to bridge the two modalities improves grounding capabilities of V&L models. Object-centricity is also a feature of widely-used web-scraped datasets: in the Conceptual Captions dataset for instance, Sharma et al. (2018) filtered out all captions which did not overlap with object labels automatically identified by a computer vision model in the corresponding image.

Some efforts have been made to understand how low-level concepts improve generalization capabilities and connect to high-level concepts. Object-centric captions help to improve the generalization over unseen objects (Hu et al., 2021) and play a role in the model understanding of abstract concepts (Cafagna et al., 2022; Wang et al., 2022b). In our work, we are interested in the relations between what Hodosh et al. (2013) refer to as 'conceptual' and 'non-visual' descriptions, which we re-frame as a distinction between low-level (object-centric) and high-level descriptions in multimodal learning. We release a novel dataset to foster research in this direction.

Motivation for the present work is also provided by recent research exploring the visual correlates of inferences, temporal and causal relationships (e.g., Park et al., 2020), which also have implications for generation. In visual storytelling, for instance, a model has to understand actions and interactions among the visually depicted entities (Huang et al., 2016; Hu et al., 2020; Lukin et al., 2018; Hong et al., 2023). Identifying actions is a prerequisite for predicting their motivations or rationales as well as explaining automatically generated descriptions of images (Hendricks et al., 2018). Actions and intention are paramount to performing commonsense and temporal reasoning on visual inputs. Along these lines, Park et al. (2020) creates dynamic stories on top of static images, where the task is to predict priors and subsequent actions and rationales. Our work is similar in spirit, as we align

high-level descriptions of *actions* and *rationales* with low-level descriptions of static images.

Some work has also been done to test multimodal model grounding capabilities from a more linguistic perspective. Parcalabescu et al. (2022) build a benchmark to test models on a variety of linguistic phenomena, like spatial relations, counting, existence, etc. Pezzelle et al. (2020) assess the integration of complementary information of V&L models across modalities, while Thrush et al. (2022) test multimodal models on compositional reasoning. In this context, the HL Dataset proposed here can offer another benchmark for V&L models' understanding of high-level descriptions of images. Such descriptions are licensed by the entities depicted in the visual modality and the relationships between them but they do not mention them explicitly.

## 3 Data

In this section, we describe the protocol used to collect annotations for *scenes, actions* and *rationales* and the subsequent collection of confidence scores through crowdsourcing. Differently from previous works, such as COCO, where human annotators are instructed to be objective and to mention only the objects clearly visible in the picture, we elicit high-level concepts in the form of captions by encouraging the annotators to rely on their subjective interpretation of the image.

### 3.1 Data collection

The task of collecting high-level descriptions is by nature hard to define and requires a clear and careful formulation, therefore we run a pilot study with the double goal of collecting feedback and fine-tuning the task instructions. Full details of the pilot are reported in Appendix D.

**Procedure** The participants are shown an image containing at least one human subject and three questions regarding three aspects or axes: *scene, actions* and *rationales* i,e. *Where is the picture taken?*; *What is the subject doing?*; and *Why is the subject doing it?* We explicitly ask the participants to rely on their personal interpretation of the scene and add examples and suggestions in the instructions to further guide the annotators. Moreover, differently from other VQA datasets like (Antol et al., 2015) and (Zhu et al., 2016), where each question can refer to different entities in the image, we systematically ask the same three questions about the

same subject for each image. See Appendix D for the full instructions and Appendix C for details regarding the annotations costs.

**Images** As mentioned in Section 1 the COCO dataset has a very explicit object-centric orientation, therefore it provides a good starting point to select images, such that we can couple object-centric and high-level captions in a resource-lean approach. Moreover, the alignment of object-centric and high-level captions permits an investigation of the relationship between them.

We randomly select 14,997 images from the COCO 2014 train-val split. In order to answer questions related to *actions* and *rationales* we need to ensure the presence of a (human) subject in the image. Therefore, we leverage the entity annotation provided in COCO to select images containing at least one person.

The whole annotation is conducted on Amazon Mechanical Turk (AMT). We split the workload into batches in order to ease the monitoring of the quality of the data collected. Each image is annotated by three different annotators, therefore we collect three annotations per axis.

### 3.2 Confidence Scores

The high-level descriptions are collected by asking the participants to interpret the scene leveraging their personal experience. The element of subjectivity leads us to expect some variation in the resulting descriptions, especially where annotators need to infer actions and rationales. In order to distinguish what can confidently be considered widely-shared, or 'commonsense' descriptions, from more idiosyncratic interpretations, we conduct a separate study where we crowd-source *confidence scores* for each high-level caption. We ask independent participants to score the likelihood of a high-level description given the image and the corresponding question on a Likert scale from 1 to 5. For a detailed example of the form see Figure 23 in Appendix D.

**Agreement-based worker selection** The confidence scores are collected following the same protocol used to collect the high-level descriptions. Using the data from our pilot study, which was carried out with participants who had been thoroughly briefed on the task, we ran a preliminary qualification task where we employed an *automatic worker selection method* to hire qualified annotators from the crowd-sourcing platform.

Let's consider the participants of the pilot as gold annotators (as they were trained on the task) and their annotations as reference annotations. The inter-annotator agreement computed on the reference annotations can be considered the gold inter-annotator agreement $\alpha_{gold}$ of the task.

We run the qualification task using the same set of items used in the pilot, then for each worker $w$ we re-compute the inter-annotator agreement (Hayes and Krippendorff, 2007), combining the workers and the reference annotations, obtaining $\alpha_w$. We compute an agreement ratio

$$r = \frac{\alpha_w}{\alpha_{gold}} \qquad (1)$$

Then, we select the worker $w$ if $r > t$, where $t$ is a threshold empirically set to $0.5$. This is equivalent to choosing workers such that their contribution does not negatively affect $\alpha_{gold}$ by a factor greater than $t$. In other words, the workers are selected if they are relatively compliant with the gold annotators.

## 4 Dataset Analysis

In this section, we analyse the captions collected in the High-Level Dataset. To provide insights on the kind of captions collected, we analyse the distribution of the captions across different axes, also comparing them with the object-centric COCO captions[2]. Furthermore, we perform a grammatical error analysis, which we report in Appendix A.1.

### 4.1 High-Level descriptions

We collected 3 annotations per axis over a set of 14,997 images for a total of 134,973 captions. An example of high-level descriptions aligned with the original object-centric caption from COCO is shown in Table 1. We expect to observe shorter texts in the high-level captions as annotators were not giving highly descriptive details typical of object-centric captions. This is visible in Figure 1, which shows that the length of the high-level captions is roughly half of the object-centric COCO captions. Though shorter, they have a comparable number of unique tokens over all the axes (as reported in Table 2); this suggests that the high-level captions are not repetitive and contain a fair amount of lexical variability. A more detailed comparison of the statistics is reported in Table 2.

---

[2]The analysis is performed by using Spacy v.3 pipeline for English using the `en_core_web_md` model to analyse the part of speech of the texts.

| Data | # Tok | Avg Len | # Uniq | # Cap |
|------|-------|---------|--------|-------|
| actions | 271168 | 6.02 | 7326 | 44991 |
| scenes | 233232 | 5.18 | 4157 | 44991 |
| rationales | 306396 | 6.81 | 8301 | 44991 |
| HL (tot) | 810796 | 6.00 | 12296 | 134973 |
| COCO | 857218 | 11.42 | 13300 | 75019 |

Table 2: HL dataset caption statistics compared the COCO captions (object-centric) for the shared set of images. We report the number of tokens (# Tok), average length (Len), number of unique tokens (# Uniq), and number of captions (# Cap).
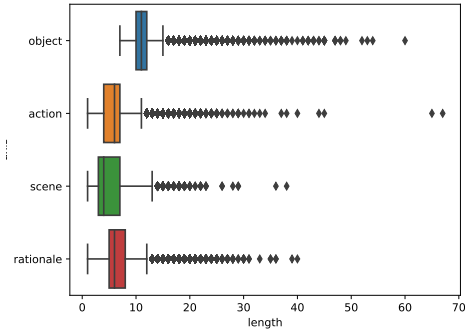


Figure 1: Caption length of the HL captions divided per axis (action, scene, rationale) in comparison to the object-centric COCO captions (object).

Moreover, as already mentioned, the COCO captions are object-centric, that is, these captions are collected to objectively represent the visual content. Although this is convenient in recognition-oriented tasks, they lack the situational knowledge required to contextualize scenes; knowledge that is instead an essential part of the cognitive processes underlying the grounding of language in vision. Indeed, as shown in Figure 2, the most frequent lemmas in the original COCO captions for the images used in the HL Dataset denote mostly objects visible in the picture. The high-level captions represent the same visual content with the addition of situational knowledge coming from the three axes, and this is also visible in different lexico-semantic choices in the texts. For example, Figure 3 shows the most frequent lemmas found in the *scene* axis. Because we align them to the same images, the dataset gives us a clean way to explore the relationship between objects and high-level axes.

**Disentangling the content across the axes** Asking the same three questions about the same subject for each image allows us to consistently compare the content of our captions across three well-defined axes. We analyse the most frequent nouns
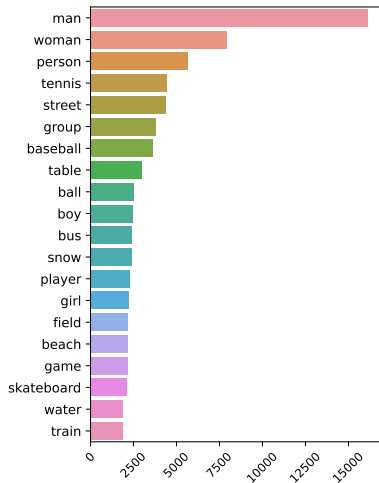
296

Figure 2: The most frequent nouns in the COCO captions of the shared set of images with the HL dataset. The majority of the terms correspond to physical objects visible in the image.
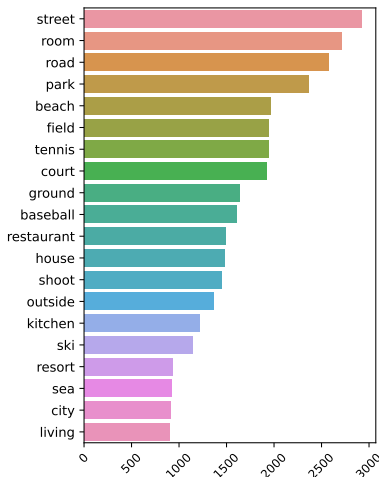


Figure 3: The most frequent lemmas of the captions in the *scene* axis of the HL dataset.



Figure 4: The most frequent verb lemmas of the captions in the *action* axis of the HL dataset.



Figure 5: The most frequent noun and verb lemmas of the captions in the *rationale* axis of the HL dataset.

in the *scene* axis in order to characterize the kind of scenes mentioned in the captions collected. The top most frequent scenes include *street, room* and *road*. These are scene types which can encompass a very broad variety of objects. However, we can also identify scenes for which a narrower range of objects would be diagnostic, for example those related to sport activities like *baseball, tennis, ski, ground* and *court*, or domestic environments like *house, kitchen* and *living* (referring to 'living rooms'). For a more complete view see Figure 3 where we report the top 20 most frequent scenes in the HL dataset.

Similarly, we can characterize also the *action* and the *rationale* axes. We identify the *action* dis-

tribution by analysing the verbs contained in the captions. In Figure 4 we observe that the most frequent actions are related to sports activities, consistently with what was observed in the *scene* axis distribution. The most frequent verbs are *play, ski, surf, skateboard*, but we can also find generic actions like *hold, walk, sit* and *eat*.

In the *rationale* axis we analyse both nouns and verbs. In this axis we expect to observe more subjectivity and content variability, with more lemmas denoting intents, mental states and events, including psych verbs. Our hypothesis is that the annotators leverage their personal experience to infer these answers to a greater extent than they do for scene descriptions.

The majority of the rationales express intentions; in fact, *want* is by far the most frequent term in the

lemmas distribution. As observed with the other two axes, terms related to sports activities are more frequent (*play, game, tennis, practice*), but also related to leisure (*enjoy, fun, vacation, love, family*) along with generic activities (*work, wait, try, eat*). For more details see Figure 5.

The systematic disentanglement of the content along three axes can serve as a filter to identify or analyse sub-samples of the data with specific characteristics. For instance, as observed so far, we can confidently say that sports-related activities are predominant in the dataset.
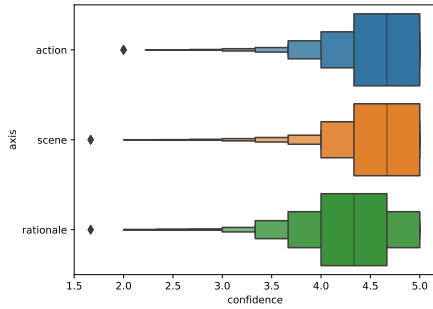
**Connecting high- and low-level concepts** One of the main goals of this resource is to enable the discovery of connections between high- and low-level captions, that are, descriptions of the same images at different levels of abstraction. By construction, the alignment provided by the HL Dataset allows us to identify concrete objects in images which provide 'support' to infer high-level concepts such as scenes, actions and rationales.

We dive deeper into our analysis and study the connection between high-level concepts related to scene, action and rationale, to low-level objects present in the aligned COCO captions. We ask: 'What are the most informative objects for a high-level concept (e.g. *enjoy*) found in a specific axis (e.g *rationale*)?'

We leverage the Point-wise Mutual Information (PMI) (Church and Hanks, 1990) to find the most informative objects linked to a high-level concept. This is helpful to discover connections between concepts across different levels of abstraction but also gives clues on the content distributions within the axes. We filter out object mentions which have a frequency less than 100 in the low-level captions. This leaves 475 object-denoting lemmas. Then, we compute the PMI between content words in the high-level captions and all these lemmas. For example, Figure 6 shows the nouns in the object-centric captions which have the strongest PMI with the verb 'enjoy' in the rationale axis.

We can observe that high-level captions can express different nuances of the same abstract concept. To take another example, *love* (in Figure 7) can refer to the love between an animal and its owner, between two partners (e.g. *wedding*) or the love for sports (e.g. *skate, snowboard*). In the same way, as shown in Figure 6 a general concept like *enjoy* can be characterized by object-level concepts leaning toward a specific nuance of meaning,



Figure 6: Most informative objects for the word *enjoy* in the *rationale* axis. Font size is proportional to PMI.



Figure 7: Most informative objects for the word *love* in the *rationale* axis. Font size is proportional to PMI.

like sports activities (e.g. *kite, snowboarder, skier*) or places (e.g. *sandy shore, ocean, lake*). More examples are provided in Appendix A.2.

## 4.2 Confidence scores analysis

Our confidence scores are similar in spirit to the *self-confidence* scores collected in the VQA dataset (Antol et al., 2015). However, they differ insofar as our scores are not self-reported by the authors of the captions, but collected from independent annotators. The inclusion of an external judgment plays an important role in determining the reliability of interpretation operated by the annotators in the caption collection and therefore, in shedding light on the extent to which an annotator's interpretation of a scene relies on 'shared' or 'commonsense' knowledge, or is entirely idiosyncratic.

We observe an average confidence score of 4.47 on a Likert scale from 1 to 5 (with a standard deviation of 0.78 and a median of 5) over all the axes. This suggests that, overall, according to independent judges, our high-level captions succeeded in capturing shared or 'commonsense' high-level interpretations of the scene.

Furthermore, the confidence scores provide an

Figure 8: Axis-wise confidence score distribution of the high-level captions.



| Idx | Scene caption | Confidence |
|---|---|---|
| 1 | in the restaurant | 1 |
| 2 | in the entrance of the library | 1 |
| 3 | the picture is taken outside a library | 3 |

Figure 9: Example of a 'hard' sample in the HL dataset where the scene captions have low confidence scores.

| Model | Axis | Cider | SBLEU | Rouge-L |
|---|---|---|---|---|
| GIT | action | 110.63 | 15.21 | 30.43 |
| | rationale | 42.58 | 5.90 | 18.57 |
| | scene | 103.00 | 24.67 | 33.92 |
| BLIP | action | 123.07 | 17.16 | 32.16 |
| | rationale | 46.11 | 6.21 | 19.74 |
| | scene | 116.70 | 26.46 | 35.30 |
| ClipCap | action | **176.54** | **27.37** | **39.15** |
| | rationale | **78.04** | **11.71** | **25.76** |
| | scene | **145.93** | **36.73** | **42.83** |

Table 3: Automatic metrics for baselines (GIT, BLIP, and ClipCap) fine-tuned along the three axes (*scene, action*, and *rationales*) of the HL dataset. The results are the average of 5 evaluation runs, by keeping the same decoding strategy and parameters for all the models.

## 5  Baselines and results

In this section, we show how the dataset can be used to finetune models to generate high-level, aspect-specific descriptions, e.g. image-to-scene or image-to-action. Below, in Section 6, we also describe a data augmentation and generation experiment, to merge the three axes into more 'narrative-like' descriptions of images.

We provide baselines for this task by fine-tuning three models, namely GIT (Wang et al., 2022a), BLIP (Li et al., 2022), and ClipCap (Mokady et al., 2021) on each separate axis. All the baselines were trained for a maximum of 10 epochs using a learning rate of $5e-5$, Adam optimizer, and half-precision (fp16).

Table 3 displays automatic evaluation results for the three models, on each axis. The first observation is that ClipCap outperforms by far the other models in each separate axis. Differently from the other models, which are natively multimodal, Clip-Cap leverages a LLM to generate captions, conditioning the text generation on a prefix representing the visual information, which is obtained by a mapping network trained to generate the prefix from CLIP's (Radford et al., 2021) image embeddings.

A second observation, consistent with the analysis presented in earlier sections, is that on all metrics, models fine-tuned to generate rationale-based descriptions receive lower scores. We hypothesise that this is due in part to the greater variability in this axis, and to its inherent difficulty, as reflected in lower confidence scores. Future work could leverage these scores as additional signal in fine-tuning models on captions that require more inference, compared to more descriptive ones.

additional perspective under which our data can be characterized: by performing an axis-wise analysis of the confidence scores distribution (see Figure 8), we observe that the *scene* and *action* captions feature the highest overall confidence, while the *rationale* axis lags behind by a small margin. We expect such differences, since determining the rationale of an action depicted in a static image is challenging, in particular, because annotators can leverage significant visual cues, but have no access either to temporal information or the subject's stated intentions. Therefore, they need to resort to their own priors and expectations which can also lead to idiosyncratic interpretations which independent judges – as in our confidence score analysis – would find relatively unlikely.

One important use of confidence scores is to provide a measure of uncertainty of the data, which can be used, for instance, to identify hard samples; an example is shown in Figure 9. The scene is hard to interpret even for humans and the scene captions display more variability and have low confidence scores. A detailed analysis of lexical and semantic variability in the presence of high-confidence scores is reported in Appendix A.3.

## 6 Data augmentation and narrative generation

We now describe how we extend the dataset to combine the three axes to compose a short 'narrative', which describes the scene, action and rationale in tandem. We call this new dataset HL Narratives. To do this, we leverage the individual axes and synthesise this part of the data using a pre-trained language model. Since scenes, actions, and rationales were elicited individually in a visually grounded and controlled setting, a synthesised version of the three individual captions should also be true of the image to the same extent (modulo the variations in confidence that we observe).

### 6.1 Data generation process

We frame the synthesis of narrative captions as a paraphrasing task. We follow a human-in-the-loop approach consisting of three stages: (i) we manually annotate a small sample of gold data; (ii) we fine-tune a large pre-trained language model (LPLM); (iii) we use the fine-tuned model to generate a sample of data, which is manually corrected and then (iv) added to the gold annotations before fine-tuning again. This procedure allows us to use only a few iterations to annotate quickly a considerable amount of data because the model improves the quality of the generated data, making manual correction progressively easier.

We use a version of T5 (Raffel et al., 2020) already fine-tuned on paraphrase generation[3] as LPLM data generator. We initialise the process with manually paraphrased annotations for 50 images ($3 \times 50 = 150$), fine-tune the model for 2 epochs, and generate 150 captions for another 50 images, which are manually corrected and added to the original 150. The model is then fine-tuned for a further two epochs. In each iteration, we reserve $10\%$ as validation data. After two epochs, we observe that the validation loss does not improve further. Finally, in the last iteration, we use all gold data to fine-tune the model and generate synthetic high-level captions for the whole HL dataset, obtaining 14,997 synthetic captions for training and 1499 for testing. In addition to the T5 paraphrase model, we also experimented with LLaMA (Touvron et al., 2023) in a few-shot setting; however, we find that T5 outperforms LLAMA in this task.

---

[3]Details about the T5 fine-tuned on paraphrase generation are available at `https://huggingface.co/Vamsi/T5_Paraphrase_Paws`.

| Model | SacreBLEU | ROUGE-L | Cider |
|---|---|---|---|
| GIT (PRE) | 1.23 | 11.91 | 18.88 |
| BLIP (PRE) | 3.47 | 15.21 | 24.15 |
| ClipClap (PRE) | 8.72 | 19.45 | 40.47 |
| GIT (FT) | 11.11 | **27.61** | 75.78 |
| BLIP (FT) | **11.70** | 26.17 | **79.39** |
| ClipCap (FT) | 8.15 | 24.53 | 63.91 |

Table 4: Results of the narrative generation task, averaged over 5 runs using the same decoding parameters for all models. PRE: pretrained models; FT: finetuned on the synthetic data.

See Appendix B for full details.

### 6.2 Results

We build three baselines by fine-tuning the same three large pre-trained models used in Section 5: GIT, BLIP, and ClipCap on our synthetic narrative captions. We fine-tune for 3 epochs with batch size 8, learning rate $5e^{-5}$, and Adam optimizer with weight decay (Loshchilov and Hutter, 2017). We test on our gold human-annotated data. As shown in Table 4, where we report results for automatic metrics, overall the models achieve worse results than in the aspect-specific caption generation task (reported in Table 3). This further highlights the difficulty of generating narrative captions of this kind for models trained on object-centric captions.

Notably, the best-performing model in the aspect-specific caption generation task, namely ClipCap, is the worst in the narrative caption generation, though by a small margin (Table 4). This suggests that although a conditioned LLM can greatly adapt to generate high-level descriptions of specific aspects of the scene, it struggles in generating comprehensive high-level descriptions involving multiple high-level aspects of the scene. Ultimately, this suggests that the multimodal representations learned by multimodal models are more robust and effective in generating natural captions than conditioned unimodal models such as ClipCap.

However, the exposure to a small amount of synthetic high-level captions is sufficient to drive the models' generated text toward more narrative-like outputs. See Appendix F for more examples from all models. Further progress can be done in this direction, for example by incorporating confidence scores during finetuning.

## 7 Further uses of the HL Dataset

We envision a wide set of use cases and tasks enabled by the HL Dataset.

GIT (PRE): a group of people on the beach
GIT (FT): people enjoying sunbathing, the picture was taken on the beach and are going to have fun and entertainment

GIT (PRE): two girls looking at their cell phones
GIT (FT): they are reading a text message outside on the street, waiting for their friend.

Figure 10: Comparison between the object-centric captions generated by GIT pre-trained (PRE) and the high-level caption generated by the fine-tuned (FT) model. The generated high-level caption embeds high-level information regarding action, rationale, and scene, depicted in the visual content.

**V&L generative tasks** Our captions support image captioning generation tasks which encompass a broader range of visually grounded linguistic descriptions than the highly object-centric, 'conceptual' descriptions which dominate the captioning literature Hodosh et al. (2013). Moreover, the decomposition along three axes can be exploited to compose narratives of the image, as in image paragraph generation (Wang et al., 2019) and visual storytelling (Huang et al., 2016; Hu et al., 2020). They can be used in combination with the question each axis corresponds to, in order to generate micro-dialog scenarios.

We would also argue that the high-level captions are also more natural and human-like, since they were collected without enforcing any restriction on the content to be described. Given that the images are also aligned with object-centric captions, it is possible to envisage a scenario in which a model is trained to generate high-level captions, which are 'explained' or justified with reference to low-level, object-centric properties (see Hendricks et al., 2016, 2018, for some work in this direction). In this way, the dataset can be leveraged to provide captions and explanations. Furthermore, the confidence scores serve for the identification of hard samples in the data, both for evaluation purposes and to provide additional training signals, as recently shown by Ouyang et al. (2022).

**Multimodal Grounding** HL Dataset is also a useful resource to benchmark the grounding capabilities of large pre-trained V&L models. Along these lines, Cafagna et al. (2021) study the capability of V&L models to understand scene descriptions in zero-shot settings, finding that only large-scale pre-trained V&L models have enough generalization capabilities to handle unseen high-level

scene descriptions. Cafagna et al. (2022) analyse the impact of exposure to high-level scene descriptions on multimodal representations in models pre-trained on object-centric captions. They show that exposure to high-level concepts mainly affects the model's attentional resource allocation over the visual input, even though the low-level concepts learned during pre-training provide enough signal to support and easily adapt to scene descriptions during fine-tuning. This is also supported by Wang et al. (2022b) who find that low-level concepts are needed to learn higher-level concepts, though this does not hold in the other direction.

## 8   Conclusions

In this paper, we introduced the High-Level (HL) Dataset. We extended 14,997 images from the popular COCO dataset with 134,973 human-annotated high-level descriptions systematically collected over three axes: *scene*, *action*, and *rationale*. We aligned high-level captions with object-centric captions and we provided human-collected confidence scores to measure the degree of commonsense expressed in the high-level captions. We also provided baseline results on generating captions for individual axes, as well as synthesised narrative captions by combining these three high-level axes of description.

Differently from current V&L captioning datasets, the high-level captions capture the human interpretation of the scene allowing for inference and expectations. We discussed how they can be used also in combination with low-level captions to improve research in visual commonsense reasoning and multimodal grounding of visual concepts into linguistic expressions and for generative tasks, hoping to foster future research in this direction.

## Ethical Considerations

The data collection received ethical approval from the University of Malta Research Ethics Committee. This data is intended to be used for training, fine-tuning, and performing experimental evaluations of machine learning models. The dataset from which the images were originally sourced is a widely-studied, publicly available resource. As far as we are aware, the data does not contain harmful or offensive content. However, we acknowledge that any biases in the collection of images and/or captions in the original dataset will also be present in the HL Dataset.

## Supplementary Materials Availability Statement:

The HL Dataset is publicly released on GitHub[4] and Huggingface[5]. The syntetic HL Narratives Dataset described in Section 6, is publicly released on Huggingface[6]. All the baselines described in Section 5 and 6 are available on Huggingface[7].

## References

Harsh Agrawal, Karan Desai, Yufei Wang, Xinlei Chen, Rishabh Jain, Mark Johnson, Dhruv Batra, Devi Parikh, Stefan Lee, and Peter Anderson. 2019. Nocaps: Novel object captioning at scale. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8948–8957.

Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433.

Lawrence W Barsalou et al. 2008. Grounded cognition. *Annual review of psychology*, 59(1):617–645.

Yonatan Bisk, Ari Holtzman, Jesse Thomason, Jacob Andreas, Yoshua Bengio, Joyce Chai, Mirella Lapata, Angeliki Lazaridou, Jonathan May, Aleksandr Nisnevich, Nicolas Pinto, and Joseph Turian. 2020. Experience grounds language. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8718–8735, Online. Association for Computational Linguistics.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Michele Cafagna, Kees van Deemter, and Albert Gatt. 2021. What vision-language models 'see' when they see scenes. *arXiv preprint arXiv:2109.07301*.

Michele Cafagna, Kees van Deemter, and Albert Gatt. 2022. Understanding cross-modal interactions in V&L models that generate scene descriptions. In *Proceedings of the Workshop on Unimodal and Multimodal Induction of Linguistic Structures (UM-IoS)*, pages 56–72, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Kenneth Ward Church and Patrick Hanks. 1990. Word association norms, mutual information, and lexicography. *Computational Linguistics*, 16(1):22–29.

Andrew F Hayes and Klaus Krippendorff. 2007. Answering the call for a standard reliability measure for coding data. *Communication methods and measures*, 1(1):77–89.

Lisa Anne Hendricks, Zeynep Akata, Marcus Rohrbach, Jeff Donahue, Bernt Schiele, and Trevor Darrell. 2016. Generating Visual Explanations. In *Proceedings of the 2016 European Conference on Computer Vision (ECCV'16)*, Amsterdam. ArXiv: 1603.08507.

Lisa Anne Hendricks, Ronghang Hu, Trevor Darrell, and Zeynep Akata. 2018. Generating counterfactual explanations with natural language. *arXiv preprint arXiv:1806.09809*.

Micah Hodosh, Peter Young, and Julia Hockenmaier. 2013. Framing image description as a ranking task: Data, models and evaluation metrics. *Journal of Artificial Intelligence Research*, 47:853–899.

Xudong Hong, Asad Sayeed, Khushboo Mehra, Vera Demberg, and Bernt Schiele. 2023. Visual writing prompts: Character-grounded story generation with curated image sequences. *Transactions of the Association for Computational Linguistics*, 11:565–581.

Junjie Hu, Yu Cheng, Zhe Gan, Jingjing Liu, Jianfeng Gao, and Graham Neubig. 2020. What makes a good story? designing composite rewards for visual storytelling. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 7969–7976.

Xiaowei Hu, Xi Yin, Kevin Lin, Lei Zhang, Jianfeng Gao, Lijuan Wang, and Zicheng Liu. 2021. Vivo: Visual vocabulary pre-training for novel object captioning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 1575–1583.

Ting-Hao Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Aishwarya Agrawal, Jacob Devlin, Ross Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, et al. 2016. Visual storytelling. In *Proceedings*

---

of the 2016 conference of the North American chapter of the association for computational linguistics: Human language technologies*, pages 1233–1239.

Klaus Krippendorff. 2018. *Content analysis: An introduction to its methodology*. Sage publications.

V. I. Levenshtein. 1966. Binary Codes Capable of Correcting Deletions, Insertions and Reversals. *Soviet Physics Doklady*, 10:707.

Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. 2022. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR.

Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. 2020. Oscar: Object-semantics aligned pre-training for vision-language tasks. In *European Conference on Computer Vision*, pages 121–137. Springer.

Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.

Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer.

Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.

Stephanie Lukin, Reginald Hobbs, and Clare Voss. 2018. A pipeline for creative visual storytelling. In *Proceedings of the First Workshop on Storytelling*, pages 20–32, New Orleans, Louisiana. Association for Computational Linguistics.

Kiwan Maeng, Alexei Colin, and Brandon Lucia. 2017. Alpaca: Intermittent execution without checkpoints. *Proceedings of the ACM on Programming Languages*, 1(OOPSLA):1–30.

Ron Mokady, Amir Hertz, and Amit H Bermano. 2021. Clipcap: Clip prefix for image captioning. *arXiv preprint arXiv:2111.09734*.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

Letitia Parcalabescu, Michele Cafagna, Lilitta Muradjan, Anette Frank, Iacer Calixto, and Albert Gatt. 2022. VALSE: A task-independent benchmark for vision and language models centered on linguistic phenomena. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8253–8280, Dublin, Ireland. Association for Computational Linguistics.

Jae Sung Park, Chandra Bhagavatula, Roozbeh Mottaghi, Ali Farhadi, and Yejin Choi. 2020. Visualcomet: Reasoning about the dynamic context of a still image. In *European Conference on Computer Vision*, pages 508–524. Springer.

Sandro Pezzelle, Claudio Greco, Greta Gandolfi, Eleonora Gualdoni, and Raffaella Bernardi. 2020. Be different to be better! a benchmark to leverage the complementarity of language and vision. In *Findings of the association for computational linguistics: EMNLP 2020*, pages 2751–2767.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Belgium, Brussels. Association for Computational Linguistics.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551.

Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online. Association for Computational Linguistics.

Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, Melbourne, Australia. Association for Computational Linguistics.

Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Candace Ross. 2022. Winoground: Probing vision and language models for visio-linguistic compositionality. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5238–5248.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4566–4575.

Jianfeng Wang, Zhengyuan Yang, Xiaowei Hu, Linjie Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu, and Lijuan Wang. 2022a. Git: A generative image-to-text transformer for vision and language. *arXiv preprint arXiv:2205.14100*.

Jing Wang, Yingwei Pan, Ting Yao, Jinhui Tang, and Tao Mei. 2019. Convolutional auto-encoding of sentence topics for image paragraph generation. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, IJCAI'19, page 940–946. AAAI Press.

Zhecan Wang, Haoxuan You, Yicheng He, Wenhao Li, Kai-Wei Chang, and Shih-Fu Chang. 2022b. Understanding ME? multimodal evaluation for fine-grained visual commonsense. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9212–9224, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. 2021. Vinvl: Revisiting visual representations in vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5579–5588.

Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. 2018. Texygen: A benchmarking platform for text generation models. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, pages 1097–1100.

Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. 2016. Visual7w: Grounded question answering in images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4995–5004.

## Appendix

## A  Data Analysis Details

### A.1  Quantitying grammatical errors

We ask two postgraduate students experts in linguistics to correct grammatical errors in a sample of 9900 captions, 900 of which are shared between the two experts. They are shown the image-caption pairs and they are asked to edit the caption whenever they identify a grammatical error. The most common errors reported by the annotators are:

- Misuse of prepositions;

- Wrong verb conjugation;

- Pronoun omissions.

In order to quantify the extent to which the corrected captions differ from the original ones, we compute the Levenshtein distance (Levenshtein, 1966) between them.

We observe that 22.5% of the sample have been edited and only 5% with a Levenshtein distance greater than 10. This suggests a reasonable level of grammatical quality overall, with no substantial grammatical issues. This can also be observed from the Levenshtein distance distribution reported in Figure 11. Moreover, the human evaluation is quite reliable as we observe a moderate inter-annotator agreement ($\alpha = 0.507$, (Krippendorff, 2018)) computed over the shared sample.

### A.2  PMI analysis examples

The PMI analysis can provide interesting insight into the connection between object-level and high-level captions on all the three axes available.

On the *scene* axis, for instance, the PMI gives some clues on the extent to which an object can be considered diagnostic for a scene. For instance, two semantically similar scenes like *restaurant* (see Figure 12) and *kitchen* (see Figure 14) share several diagnostic objects, as we would expect. However, we can identify important semantic nuances: the scene *restaurant* contains objects related to the food (i.e. *pizza, cheese, wine, sandwhich*) whereas *kitchen* contains objects related to the preparation of food (i.e. *stove, oven, tray, refrigerator*). Another example is shown in Figure 13, where the most relevant objects for the action *look* encompass a wide variety of contexts, like looking at a screen or a device (e.g. *device, screen, cellphone*) or entertainment (e.g. *zoo, zebra, giraffe*). For more examples see Table 5, where are shown the top most relevant objects for the top three lemmas in the *scene, action* and *rationale* axes.

These semantic differences, while quite easy for humans to interpret, are not usually present in object-centric V&L datasets. They are made explicit and easy to identify in the HL dataset, where captions with different levels of abstraction are aligned with the same image.



Figure 11: Distribution of the Levenshtein distance computed between the original and the corrected high-level captions in a sample of 9900 captions.



Figure 12: Most informative objects for the word *restaurant* in the *scene* axis. Font size is proportional to PMI.

Figure 13: Most informative objects for the word *look* in the *action* axis. Font size is proportional to PMI.



Figure 14: Most informative objects for the word *kitchen* in the *scene* axis. Font size is proportioanl to PMI.

| Axis | Top Lemmas | Top Objects (PMI) |
|------|-----------|-------------------|
| scene | street | intersection, decker, meter |
| | room | living, wii, nintendo |
| | road | traffic, decker, intersection |
| action | play | nintendo, wii, swing |
| | ride | rider, carriage, wave |
| | hold | controller, remote, rain |
| rationale | want | mirror, bathroom, sink |
| | enjoy | wave, kite, ocean |
| | fun | wii, nintendo, controller |

Table 5: Top most informative objects of the top most frequent lemmas in the three axes (*scene, action, rationale*) according to PMI.

## A.3 Quantifying Lexical and Semantic Diversity

In Section 4.2, we showed that in the presence of low confidence, there can be variation or disagreement among high-level captions given by different annotators for the same axis. In such cases, the captions focus on different aspects or refer to different interpretations. Although this phenomenon has been observed for captions with a low confidence score, it is conceivable that it might also happen with high-confidence captions, for example, two captions annotated by different annotators, while differing in the interpretation of an image, could nevertheless be considered highly likely. To quantify this phenomenon, in this section we further expand our analysis by studying the lexical and semantic diversity of our captions.

**Purity score** We leverage the BLEURT score (Sellam et al., 2020), a trainable metric used to evaluate semantic differences in Natural Language Generation, to compute a score measuring the semantic diversity among the high-level captions associated with an image. To do so, we first compute such scores across each axis, and then we combine them to obtain a final score for the item. In this way, we can unpack the semantic diversity item-wise and axis-wise.

Let $C$ be the set of high-level captions of a given axis (e.g. scenes) for a given image. For simplicity, we do not report the index of the image and the axis in the following notation. We compute the BLEURT score of the caption as follows:

$$s_i = BLEURT(c_i, ref) \qquad (2)$$

where $s_i$ is the resulting BLEURT score, $c_i$ is a high-level caption, and $ref$ is the set of reference captions defined as follows:

$$ref := \{c_j \mid c_j \in C \text{ and } j \neq i\} \qquad (3)$$

In other words $ref$ is the set of remaining captions along the axis and therefore, $s_i$ is measuring the semantic diversity of the caption with respect to the other captions along the same axis.

By averaging the caption-wise scores across a single axis and across all the axes we obtain a *purity score* measuring the semantic consistency both axis-wise and item-wise.

**Diversity score** Along the same lines, we propose the *diversity score*, to measure the lexical diversity of the captions. The *diversity score* follows

306

the same logic implemented to compute the *purity score* introduced in the previous paragraph, but the BLEURT score in Eq. 2 is replaced by the BLEU score (Papineni et al., 2002) and then normalized between 0 (similar) and 1 (very different). Our score is similar in spirit to self-BLEU (Zhu et al., 2018) as it measures the similarity of the captions within their own distribution. However, its computation concerns only axis-wise and item-wise captions.

### A.3.1 Results and discussion

As shown in Figure 15 the purity scores obtained are mostly negative, this is due to lexical variations, which the BLEURT score is known to be sensitive to (Sellam et al., 2020). However, BLEURT is not defined in any specific interval thus, it is usually hard to interpret (Sellam et al., 2020) if not considered in relative terms. Based on that, we use it to



Figure 15: Axis-wise purity score distribution.

compare the semantic purity across items and axes within our dataset. As shown in Figure 15, *action* and *scene* share similar purity score distributions whereas the *rationale* is more skewed to the left than the other axes. This shows that the rationales feature a higher semantic diversity (lower overall BLEURT) than the other axes.

The *rationale* axis is also the one featuring the highest lexical diversity, whereas the *scene* and the *action* have similar distributions. This is shown in Figure 16 where the *rationale* density estimate (in green) has a higher peak skewed on the right-hand side than *scene* and *action* density estimate (respectively in orange and blue).

We have similar observations for both *purity* and the *diversity* scores and this confirms what was



Figure 16: Axis-wise diversity score distribution. The scores have been normalized between 0 and 1.



Figure 17: Pearson correlation between confidence, diversity and purity scores.

observed in the confidence score analysis in Section 4.2, namely that the task of determining the rationale of an action from a static image produces more variation and divergent interpretations leading to higher semantic and lexical diversity. Moreover, we find that both the *diversity* and the *purity* scores positively correlate with the confidence scores (See Figure 17).

### A.3.2 Item-based analysis

An item in the HL dataset is an image along with all the high-level captions of all the axes. For instance, Figures 18 and 19 show the item-wise *diversity score* and *purity score* distribution respectively, along with their average value across the whole dataset. An item on the right-hand side of the distribution is systematically more consistent across its axes with respect to the measure considered (*purity* or *diversity*). This information can
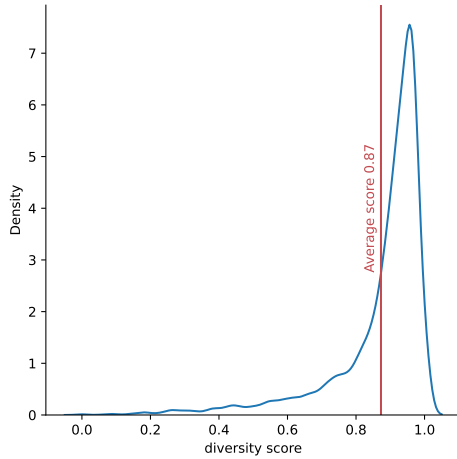
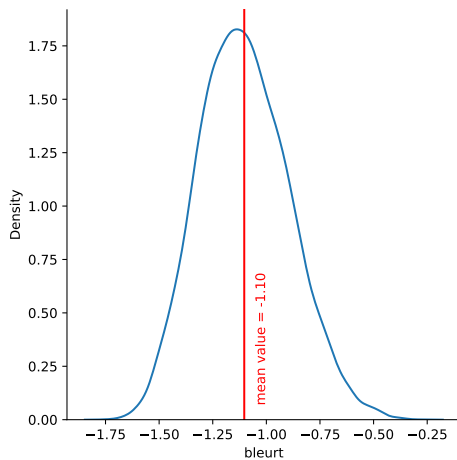307

Figure 18: Item-wise diversity score distribution.



Figure 19: Item-wise purity score distribution.

be combined with confidence scores to perform a more fine-rained sample selection. For example in zero-shot testing, we might want to use a hard sample to test our model with, we can select items with similar lexicons, low-semantic purity, and low confidence scores.

## B  Narative Caption Generation Task Details

### B.1  Few-shots Prompting Data Generation

We test an alternative data generation pipeline by leveraging the in-context learning capabilities featured by the most recent large language models (LLM) (Brown et al., 2020; Maeng et al., 2017; Touvron et al., 2023). This data generation approach has the advantage of not requiring any model fine-tuning.

We design a prompt for our task and we use it to generate data from the recently developed LLaMA

Given three sentences merge them into one sentence, and make sure that the sentence is grammatically correct. Here is an example:'in a beach',' holding an umbrella',' so they won't get a sunburn' <holding an umbrella in the beach so that they won't get a sunburn.>\n The three sentences are: **'scene','action','rationale'** <

Figure 20: Prompt used for the data generation. The parts in bold are replaced with the corresponding high-level descriptions for the given sample.

model (Touvron et al., 2023). The prompt consists of the task description, followed by an example and the inputs of the task written in natural language. The full prompt is shown in Figure 20. The resulting output is then post-processed to extract the generated high-level caption.

**Discussion**    As described in Section 6, we build baseline image captioning models starting from GIT-base and fine-tuning on the LLaMA- and T5-generated synthetic data. The best model is chosen on a combination of qualitative models' output inspections and automatic metrics (SacreBLEU (Post, 2018), ROUGE-L (Lin, 2004) and Cider (Vedantam et al., 2015)) computed over the gold data.

In Table 6 we show the results of the evaluation based on the automatic metrics. First, we observe that the performance of the pre-trained model (PRE) is extremely poor, in the high-level caption generation task, highlighting the substantial difference between captions of this kind with traditional object-centric captioning the pre-trained model is trained on.

Second, focusing on the fine-tuned models, we observe that GIT fine-tuned on T5-generated data performs better than the LLaMa-based counterpart on the automatic metrics. We argue that the model trained on T5-generated synthetic data benefits from the exposure of the data generator to the gold data distribution. However, we point out that the few-shot data generation pipeline remains a valid alternative as it achieves comparable performance without requiring any further fine-tuning.

## C  Annotation Costs

In this section, we report the costs related to the data collection.

**High-level caption collection**    Overall 1033 participants took part in the caption data collection, they were paid $ 0.04 per item corresponding to the hourly minimum rate in the United Kingdom. In total, the data collection cost $ 1938.

| Model | SacreBLEU | ROUGE-L | Cider |
|---|---|---|---|
| GIT(PRE) | 1.23 | 11.91 | 18.88 |
| GIT(T5) | **11.07** | **31.37** | **74.79** |
| GIT(LLaMA) | 10.96 | 24.71 | 65.05 |

Table 6: Automatic metrics computed over the gold annotated high-level captions; the scores are the average results of 5 runs using the same decoding parameters for all models. We compare the pre-trained model (PRE) with the model finetuned on T5-generated (T5) and LLaMA-generated (LLaMA) data.

**Confidence Scores collection**   The qualification task for confidence scores led to the recruitment of 53 annotators. We found that this task was harder than the high-level caption annotation in terms of complexity but not in terms of execution time which was indeed shorter. Therefore, in order to encourage good quality annotations, we pay $ 0.04 per item. Considering the time needed to perform the task, this corresponds to 4 times the hourly rate of the minimum wage in the United Kingdom. The qualification task and the data collection cost respectively $ 93 and $ 1938.

## D   Annotation Details

### D.1   Pilot

We run a pilot study with the double goal of collecting feedback and defining the task instructions. The pilot is run with 9 participants who were trained on the task, with high proficiency in English and a background in computer science and linguistics.

With the results from the pilot we design a beta version of the task and we run a small batch of cases on the crowd-sourcing platform. We manually inspect the results and we further refine the instructions and the formulation of the task before finally proceeding with the annotation in bulk. The final annotation form is shown in Figure 22. It is important to notice that the instructions, shown in Figure 21 are always visible to the workers.

Figure 23 shows the annotation form used for the confidence score collection. Also in this case, the instructions are always visible to the worker and each image is presented along with the original question and the answer.

## E   Additional Data Examples

In Table 7 we show further examples of images and their corresponding captions in the HL Dataset.

**Instructions**:
You are going to see some pictures. Each picture involves one or more people ('the subject'). You will be asked some questions about the picture
Don't think too much, feel free to give your personal interpretation using your knowledge or common sense. Try to answer using full English sentences. **If you're not sure what the answer could be, give your best guess.** Avoid using expressions like "I think" or "I suppose" or "Maybe.
**Do not propose options or possibilities** saying for instance: something "or" something else. **Make your best guess** and state the one you choose.
Write a statement, **don't write a one-word answer**, avoid acronyms or slangs and write a **full sentence**.

1. **Where is the picture taken**: give your best guess about the type of place where the action is happening (for example, "in a ski resort");

2. **What is the subject doing**: Try to describe what the people are doing as concisely as possible.
If there is more than one person, try to choose a description that captures what all of them are doing (for example, "They are skiing")

3. **Why is the subject doing it**: here, write your best guess about why the person or persons are doing the action (for example, "They are on a family holiday")

The **What** question and the **Why** question **cannot have the same** answer.

The answers must be **written correctly in English**, check the spell and most importantly **don't forget the subject of the sentence in your answer** (he, she, it, they)

Figure 21: Final version of the instructions presented to the workers during the collection of the high-level captions. These instructions are always visible to the annotators.

## F   Examples of Narrative Caption generations

In Figure 24 we show examples of narrative caption generations from our fine-tuned baselines.

Figure 22: Annotation form presented to the worker during the high-level captions collection. The instructions (shown in Figure 21), are always visible to the annotators.



Figure 23: The confidence scores annotation form. We show the instructions, the image, the question, and the corresponding answer.

BLIP: they are playing football in a soccer field and are spending time together
Gold (T5): They are playing in a stadium they are in a game.



BLIP (FT): he skates in a snowy field and wants to enjoy the ride.
Gold (T5): He is snowboarding in a ski resort and he is on vacation.



ClipClap (FT): They are waiting for a bus to take them to the bus station
Gold (T5): at the bus stops he needs to be taken to his destination..



ClipClap (FT): He is skating on a skateboard in a skate park.
Gold (T5): He is skateboarding at a skatepark for fun.



GIT (FT): they are riding horses in the beach, they want to go on vacation.
Gold (T5): They are riding in a beach, they are in a trip..



GIT (FT): the cat is watching the dog in the kitchen, it is watching television.
Gold (T5): Two cats are watching tv in a living room and wait to be served food.

Figure 24: Examples of captions generated by the fine-tuned (FT) models and corresponding T5-generated (T5) data on the narrative caption generation task.

| Image | Axis | Caption |
|---|---|---|
|  | scene | the picture is taken in a construction site |
| | action | he is operating machinery |
| | rationale | he is clearing up debris with the machine. |
| | object-centric (COCO) | A blue flatbed truck with a yellow backhoe behind on a residential street. |
|  | scene | The photo is taken in a toilet |
| | action | the subject is sitting on the toilet seat. |
| | rationale | doing it just for fun |
| | object-centric (COCO) | A man in blue shirt sitting on toilet next to sink and mirror. |
|  | scene | the picture is taken at old town street |
| | action | one car is in the picture to turn to old town |
| | rationale | they are coming to old town |
| | object-centric (COCO) | A car driving on a street in the town center |
|  | scene | in the restaurant. |
| | action | they are having their snacks. |
| | rationale | to taste it. |
| | object-centric (COCO) | A dad and his daughter eating a meal at a small table. |
|  | scene | this is inside a garage |
| | action | the bike is just standing alone. |
| | rationale | no one is working on or trying to ride the bike. |
| | object-centric (COCO) | Custom motorcycle has a wooden barrel as a sidecar |

Table 7: Examples of instances of the High-Level Dataset. It is shown one of the three captions available for each of the three axes collected: *scene, action, rationale*, aligned with the object-centric captions from COCO.

# Validating Predictive Models Of Evaluative Language For Controllable Data2Text Generation

**Maurice Langner**
Linguistic Data Science Lab
Ruhr-Universität Bochum, Germany
Maurice.Langner@rub.de

**Ralf Klabunde**
Linguistic Data Science Lab
Ruhr-Universität Bochum, Germany
Ralf.Klabunde@rub.de

## Abstract

In data2text generation, tabular data is transformed into a text that expresses information from that source domain. While some text types, such as instructions, demand objective and neutral language without any expressive and evaluative content, many other text types are expected to provide expressions for these kinds of subjective meanings. In controllable, pipelined neural NLG separate learning models, notably regression models, can be used to predict whether some feature deviates sufficiently strongly from an expected value, so that evaluative language would be appropriate for verbalizing this finding. In this paper, we present an empirical study on the comprehension of evaluative adverbs and adjectival modifiers in car reviews, a text type that is characterized by a mixture of factual information with evaluations expressing positive or negative surprise. We show to what extend regression-based decision boundaries for producing evaluative content in controllable data2text NLG match the reader's expectations that are raised by those evaluative markers. Finally we show that regression values in combination with standard deviation of the technical input data constitute reasonable Boolean thresholds for both positive and negative surprise, which provide the basis for the development of more complex models that also include the scalar base of adverbs and modifiers.

## 1 Introduction

In controllable data-to-text natural language generation (henceforth data2text), tabular data is transformed into surface text that structures, relates and verbalizes the relevant pieces of information as close as possible to the domain-specific characteristic features and structures of the corresponding text types. The goal of adjusting the output to style, tone and structure of typical texts in the respective domain is achieved by either manually encoding template-based systems, or by fine-tuning large language models (LLMs) like BERT or GPT-X on a subset of the corpus.

An application domain that is especially challenging in regard to data2text NLG are car reviews. In these texts, technical information is seamlessly interwoven with subjective impressions of the test driver and, even more important, with the test driver's assessment of the car's features against his experience with cars of comparable quality and class.

Although a trustworthy system should not generate information that is not present in the underlying database of car features, it is possible to approximate the domain knowledge needed to reproduce the expert's assessment of the car's technical features. For example, sentence (b) in Table (1) from a driving report about the Lotus Elise Sport 2015 cannot be generated from purely technical information about the vehicle. It contains additional, subjective information on the driving experience. If a system for controllable data2text NLG shall be capable of generating sentences of this kind, the respective information must be added or the verbalisation be hard-wired into the system, including a mechanism for deciding when it is suitable. Talking about being pressed back into the seat would be inadequate, even quirky, if the car was a compact car with a 60 HP motor. LLMs may encounter such expressions during training and reproduce it when realizing the features for the target text, but controlling style and content for data2text with an LLM is a non-trivial task, which is why most NLG systems in-use do not rely on them (yet).

Accordingly, in sentence (a) in Table (1), where the strong acceleration is contrasted with a comparably low motor power, controllable data2text NLG needs a mechanism for determining when the mismatch between acceleration and power output is large enough to permit the usage of an evaluative adverb like *surprisingly* and respective modifiers like *enormous*. Again, LLMs may learn using

(a) Surprisingly, the sports car has an enormous acceleration from 0 to 60 in under 6.5 seconds, given the comparably low power of 136 hp.

(b) When you hit the gas, the first law of motion kicks in and you feel that well-known feeling rise in your stomach when being pressed back into your seat.

(c) The reason for this is that the Lotus Elise has an exceptionally low kerb weight of 1931 pounds.

Table 1: Three subsequent example sentences from a driving report

evaluative expressions from the input data, but a controllable generation, deciding where such an expression is suitable on the basis of the technical data input, should be preferred. Some examples generated by OpenAI's ChatGPT, which we report below, will clarify this point.

Sentence (c) finally gives the reason for the high acceleration given the low HP: an *exceptionally low kerb weight of only 1931 pounds*. The adverb *exceptionally* signals the impossibility to predict a value of 6.5 seconds for acceleration when only considering the power output, but when considering the joint effect with a tiny weight, the acceleration is, albeit a rare combination, technically possible.

The exemplary excerpt demonstrates the relevance of evaluative expressions on sentence and discourse level. However, although evaluative meaning often implies a subjective estimation why some features or states have a positive or negative stance, these evaluations rely on experiences and expectations that can be inferred from underlying data.

We consider evaluative expressions as a key element to generating humanlike, pragmatically rich text and understanding the underlying mechanisms for producing evaluative language in general. This is of special importance when solely being informative is not the fundamental driving force for communication. Evaluative adverbs and modifiers, which we focus on, are just one way of realizing a mentionable data mismatch that is present at the early stage of document planning. Different NLG systems, either traditional or neural models, may then proceed differently for realizing the respective evaluative content and including it in the generated output.

In this paper, we describe an empirical study on the comprehension of evaluative expressions in driving reports, performed with non-experts in order to identify the relation between their evaluative interpretation and the underlying technical data. The study inverts the typical acceptability rating approach of fixed sentences. The participants

have to define lower and upper thresholds for numerical features that make the respective phrase acceptable, allowing us to empirically approximate intervals and their match with evaluative scales. We then compare the results with regression-based approaches to data analyses to show whether these models are adequate for dynamically determining decision boundaries of evaluative language in text generation.

We use the car domain with its comprehensive technical specifications, but the approach we are presenting should be transferable to any data-oriented generation model and corresponding texts with evaluative content related to these data, such as technical devices and their use, weather reports, or data and reports from sports events, such as Football games, as in the ROTOWIRE corpus (Wiseman et al., 2017).

## 2 Related Research

Evaluative expressions (and expressive constituents in general) have received much attention in Formal Semantics and Pragmatics, revolving around the questions how these items can be integrated into a compositional framework, and how their meanings interact with truth-conditional and inferred content (Potts, 2005; Gutzman, 2015). As meaning analyses that are focusing on conditions for the use of these expressions, they can ignore the question of what world knowledge or data their use is based on, which is an indispensable aspect of data2text NLG however (Reiter and Dale, 2000; Ferreira et al., 2020; Gatt and Krahmer, 2018).

Evaluative adverbs and modifiers, the linguistic topic of this paper, are also attributed to emotion generation or affective language generation (de Rosis and Grasso, 2000), where evaluative items are generated in order to convey information with a specific stance (Elhadad, 1991), for example communicating exam marks (Mahamood et al., 2007) or tailoring texts to specific users (Balloccu et al., 2020). The same intention is inherent to car reviews that include driving reports, where the tone of the text is coloured by the author's opinion. Experimental studies in this area of research are quite complex due to the multitude of influential factors regarding audience, personality, individual preferences and the level of knowledge or common ground. Experimental results are often non-reproducible (Mahamood, 2021), since the target group is a decisive factor. Affective language generation has

been implemented in both, template-based NLG systems (de Rosis and Grasso, 2000), as well as in neural end2end language models (Singh et al., 2020; Santhanam and Shaikh, 2019). Regarding controllability, generating such evaluative markers with Transformer models and LLMs in general poses a challenge due to toxicity and fact hallucination (Ji et al., 2022; Dušek and Kasner, 2020; Shen et al., 2020). First attempts have been made to adapt GPT models to and condition on the intended evaluative stance of the text to be generated (Singh et al., 2020).

In recent approaches, end-to-end language models have been re-pipelined in order to enhance controllability of intermediate representations. Ferreira et al. (2020) propose a re-modularization of neural generation networks, chaining separately trainable and evaluable networks that are specialized for the different tasks of content selection, document planning and surface realisation. They show that these pipelined neural generation models outperform end-to-end networks, especially on unseen data, where the latter tend to produce topic-unrelated, incoherent texts and hallucinations. In pipelined data2text generation it is possible to integrate the necessary calculations for determining the adequacy of evaluative content on a polarized scale (Benamara et al., 2017) with regression and classification (Langner and Klabunde, 2022) and using large language models such as T5 (Xue et al., 2021) for transformation into surface text.

## 2.1 Large language models and evaluative tasks

ChatGPT shows in a clear way that trusting completely on LLMs does not necessarily result in the appropriate use of evaluative expressions. The following excerpt is ChatGPT's answer to the task of evaluating the specific values of the features 'acceleration' and 'motor power' of the car under discussion, viz. the Lotus Elise Sport 2015.[1]

*The acceleration of under 6.5 seconds from 0 to 60 mph with a motor power of 136 HP is pretty impressive. It suggests that the vehicle has a good power-to-weight ratio, which allows it to accelerate quickly.*
*In terms of rating, it depends on the context of the vehicle. If it's a small sports car or a high-performance sedan, then this acceleration is quite good. On the other hand, if it's a larger vehicle such as an SUV or a pickup truck, then it's still decent but not as impressive.*
*Overall, I would rate an acceleration of under 6.5 seconds from 0 to 60 mph with a motor power of 136 HP as a very respectable performance for a car.*

The first paragraph perfectly captures the fact that the Lotus Elise may accelerate so fast with comparably small power due to a "good power-to-weight ratio" (as also explained in the driving report, see sentence (c) in Table 1) and evaluates it as "pretty impressive". The second paragraph reveals the inconsistency in regard to the evaluation, a lack of coherence and the wrong calculation behind it. For a small sports car, the acceleration is now judged as "quite good", where initially it was "impressive", clearly a shift on the acceleration scale to the negative side. Furthermore, it names SUVs and pickup trucks as contrastive vehicles, for which the given acceleration would only be decent. This is in contradiction to the technical facts about motors and vehicle types. In the final paragraph, ChatGPT uses a third evaluative expression of the features, (*a very respectable performance*), and relates acceleration to power output in a generalizing statement.

Prompting ChatGPT with the task to produce a sentence for a driving report with the features above, using an adverb to express its opinion[2], results in the following sentences that were produced in the given order:

1    Impressively, the car with 136 hp achieved an acceleration of 6.5 seconds from 0 to 60 mph, indicating that it is remarkably well-designed for performance.

2    Honestly, the car's performance was underwhelming with a 6.5 second acceleration from 0 to 60 mph despite its 136 hp power output.

Both sentences correctly fulfilled the task, but only the first agrees to expert opinions on the Lotus Elise, and sentence number 2 completely

---

[1]Original prompt: *"How would you rate an acceleration of under 6.5 seconds from 0 to 60 mph with a motor power of 136 HP?"*

[2]Original prompt: *"Please produce a sentence for a road test report of a car with 136 hp and an acceleration of 6.5 seconds from 0 to 60 mph, expressing your opinion with a sentential adverb."*

contradicts the first one in its criticism. In sum, ChatGPT – as a representative example of relying on LLMs – is very well capable of enriching text with evaluations and generates sophisticated formulations, but adequacy and agreement with expert knowledge is in need of improvement.

## 3 Regression models for predicting expectations

The technical data we are using have been extracted from the ADAC database, the database of Europe's largest automobile association. The database contains technical information on vehicles of various types as well as independently written reports on these vehicles from ADAC's vehicle experts. We use this database due to the comprehensive technical details it provides in combination with the vehicle reports. Commercial interests are not present.

We have trained different regression models. The first is a standard linear regression model, which we expect to underfit the data since some of the features, e.g. 'power output' and 'acceleration', have exponential rather than a linear relation. The second is a polynomial regression model, which may underfit for the same reasons of mixed relation types between predictors and response. Furthermore, we implemented a deep neural network with intermediate, non-linear layers for regression. For these networks, we used Adam optimizer and mean absolute error as loss function. The best fitting model is the DNN with an MSE of 2.73 and an R2 score of 0.62 for the examples given in Figure (1).

In general, experts in automotive engineering have expectations about certain features of a car, given its technical details such as a certain acceleration given the power output and kerb weight. The images in Figure (1) show relations between each pair of the features 'acceleration', 'power output' and 'weight' as extracted from the database. The color of each "+" marker represents the feature that is not present on the axes. For example, in sub-figure (c) the color encodes the weight, light-blue indicating small numbers and light-brown indicating large numbers.

The green dot in each sub-figure marks the data point for the Lotus Elise which has been described by the example sentences in Table (1). The position of these data points outside of the data cloud is a strong indicator already for using evaluative expressions. However, we require a general approach for deciding whether information is evaluative at all.

In Figure (1a), there seems to be a linear relationship between power output and weight with strong variation of the power value for instances higher than 2500 kg. Due to large variation, regression models agree ever less with each other with increasing power, the polynomial parabola turning down again, the DNN taking a steeper increase upwards.

Figure (1b) shows a rather broad distribution of the relation between weight and acceleration. There seems to be a linear decrease of seconds between 1000 and 2000 kg, which forms a baseline. But the variation range suggests that weight seems to be of less importance for predicting the acceleration value. Regression models are also less equivalent at the borders above 2800 kg where fewer data is available and variation is largest.

In Figure (1c), there is an exponential decrease of acceleration time with increasing power output. Compared with the other sub-figures, variation is small so that there is a well defined relationship that can be modeled with far better fit than the other ones. According to the distribution of data points, variation seems to grow with the decrease of seconds. The curve finally converges towards a horizontal line. Here, the DNN fits the data perfectly. The polynomial model fits as well, but for power output values higher 400, the curve rises again, which is rather unrealistic and does not fit the data - it is just the nature of a second degree polynomial. The linear regression model overestimates the acceleration for power outputs higher than 400 and due to its linear nature assumes a constant decrease which is not reflected in the data either.

Using these regression models allows us to take features $x_0^c,...,x_{n-1}^c$ of car $c$ in order to predict an expected value for feature $x_n^c$ that shall be verbalised in the text. A deviation of the real feature value in the tabular data from this expected value both qualifies and quantifies the generation of evaluative adverbs or modifiers. Other options of lexicalizing the intended affective tone exist as well, but they are not subject of the empirical study described in this paper.

Applying regression models to the Lotus Elise, we see that these models predict the car to be slower, to weigh more and to have nearly double the power. All regression values therefore would justify the usage of evaluative adjectives with positive stance like *surprisingly* or *incredibly*.

There is one caveat to this setup. How do we quantify the divergence from the real value that triggers evaluative content? Its empirical counterpart is: How do evaluative adverbs influence the reader's expectation about a certain feature? Due to data sparseness, we cannot choose the threshold in such a way that the distribution of evaluative expressions in the empirical data best matches the distribution predicted by the models. For this reason, we decided to perform an empirical study for determining the intervals of values that license the usage of evaluative content on the basis of the reader's expectations.

# 4 Empirical study on expectation values raised by evaluative adverbs

In order to evaluate the adequacy of the regression models, we need to consider how evaluative adverbs and modifiers are interpreted and whether the regression models correctly capture expectations of readers. For this reason, we conducted a study on the interpretation of these evaluative items and their influence on the expectations about numerical features in driving reports. The study is designed as a webserver application, participants from Germany and Austria as well as the USA were acquired via Prolific. Participants were selected w.r.t. their first language, highest education level (at least an academic bachelor degree), possession of a driving licence and ownership of a car. For each study in German and English, 50 participants took part. Additionally, 50 further participants took part in a shorter ablation study on modifiers in German. The participants were paid 13 pounds per hour. 20 minutes were scheduled for each participant, but most were significantly faster (8-13 minutes median time). For the ablation study, 12 minutes were scheduled. For the German studies we analyse in this paper, no outliers had to be excluded. We took into account the features 'acceleration', 'mileage', 'maximum speed', 'power output' and 'displacement', the former four because they are the most well-known and intuitive features, the latter representing features that are less intuitively accessible. We used the German counterparts to the evaluative adverbs *surprisingly*, *disappointingly*, *amazingly* and *unfortunately*, as well as the modifiers *good*, *bad*, *low*, *high*, *slow*, *fast*, *average* and *enormous*.

| polarity | item |
|---|---|
| -2 | Disappointingly, the car goes slowly from 0 to 60 mph in [...] seconds with a power output of 200 hp. |
| 0 | With a power output of 200 hp, the car goes from 0 to 60 mph in [...] seconds. |
| 2 | Amazingly, the car goes from 0 to 60 mph in only [...] seconds with a power output of 200 hp. |

Table 2: Example items and their polarity of group e (estimating acceleration given power output of 200 hp)

## 4.1 Methods

The study comprises two tasks. The first task type is selecting an option from a menu as answer to a question in the car domain. These questions function as distractors, while also being a means for collecting information on the participants experience with cars that can be used for further research in regard to text production. For some items the participants were instructed to select an evaluative adverb that agrees with their judgement of the given features, as in *How would you rate the acceleration of the following car in relation to the power (hp)? "The sports car goes from 0 to 60 mph in 7.5 seconds with an output of 560 hp."*, with answer options *{surprisingly fast, normal, disappointingly slow}*.

The main task type presents a sentence from a road test report (henceforth *item*), where two features are named. These sentences are grouped by their degree of neutrality or polarity towards a positive or negative evaluation. Seven categories are possible, with polarity ratings between -2 and 2, being the most negative and most positive expressions. Sentences of category 0 are neutral, 0.5 and -0.5 contain only modifiers, whereas categories -1 and 1 only contain an evaluative adverb with respective polarity. Categories -2 and 2 contain both an adverb and one or two additional modifiers in the grammatical phrases that contain the features. Examples for three polarities are given in Table (2). In the respective group, participants were asked to estimate acceleration for a car with 200 hp given the differently polar expressions.

These items were collected by automatically extracting sentence adverbs from the ADAC corpus and manually selecting evaluative ones. These sentences containing the adverbs were randomly assigned to a predefined group of polarities. The missing polar items in each group were then manually constructed by modifying the corpus extracted
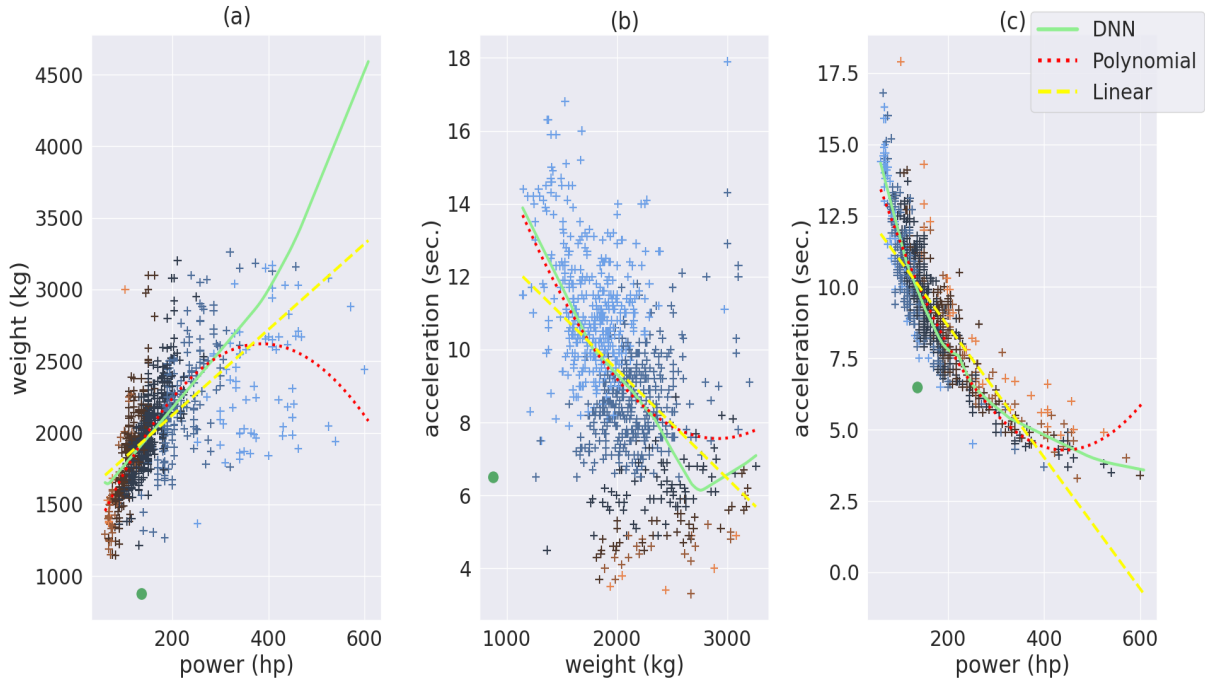
Figure 1: Bivariate plots for weight, acceleration and power output. Each "+" represents a data point. The larger green dot marks the Lotus Elise, which is an outlier in all graphs. The line plots represent three types of regression models and their fit to the data, the linear (yellow dashed), the polynomial (red dotted) and the deep neural network for regression (lightgreen solid).

item, so that the effect of different polarities can be tested within each group.

Instead of letting the participants rate the acceptability of fixed sentences, we decided to remove one of the numerical statements from the items and let the participants determine the lower and upper threshold such that the resulting interval of values agrees to their expectations the expression raises. The participants could either use a slider or two text fields to enter the thresholds. Minimal and maximal values had been determined on the basis of our database on technical features.

Each sentence provides a single session item the participants have to deal with; the order of the 15 items is randomized in order to prevent bias. Towards the end of the study 3 expressions have been shown simultaneously that express the same features, but with polarity categories -2, 0 and 2. These 3 items are not randomized and agree to 3 items all participants have seen before. The task is to adjust the thresholds also in comparison to the choices they make for the other items on the page. We integrated this final page in order to assess whether results vary when the participant's expectations for one item are delineated more clearly by the expectations raised by alternative evalua-

tions.

The study is based on the following hypotheses: We assume that the expectation values differ significantly between neutral sentences and those with evaluative adverbs in agreement with polarity. For example, if the sentence states with positive stance that a sports car has a surprisingly high maximum speed given a certain power output, we assume the participants' expectation interval to be higher than for the statement without evaluative adverb. Likewise, we assume that negative polar expressions vary significantly from positive polar questions in regard to the expected values, since it should differ from the neutral one, but in the opposite direction to the positive expression.

### 4.2 Evaluation

27 out of 32 binary comparisons between two differently polar items and the distributions of estimated values are significant. For each item, participants estimated a lower and an upper threshold that matches their expectation raised by the item's polarity and thus its evaluative stance. Across participants, this results in two normal distributions, one for each threshold. 15 items are grouped by the two contained features, the source feature as orientation and the target feature, whose numeric value

is masked and shall be estimated. Groups consist of 2 to 5 items, which means that not all possible polarities are tested for each feature pair. In each group, we tested all possible pairs of expressions for significance using ANOVA.

For illustration reasons, we concentrate on an example where participants should estimate acceleration on the basis of a given power output (compare Table 2). We have tested a multitude of feature combinations across all studies. Results are generalisations made from all these items.

As shown in Table (2), participants estimated acceleration from 0 to 100 km/h for a car with 200 hp. In Figure (2), there are two graphs containing the distribution plots for the lower threshold (left) and the upper threshold (right) for each of the three items in Table (2). The two maxima of the neutral expression's curve for the lower threshold are located between 3 and 7.5 seconds for acceleration, the single maximum for the upper threshold is positioned at about 11 seconds. We judge these curves as approximations of the lower and upper bound of acceleration values that are normal for a car with 200 hp. Now, the distributions for both, the positive item and the negative one, vary from the neutral item, both in opposite directions of each other but in agreement with their evaluative pole. Distributions for both thresholds of the positive item are shifted to the left towards better acceleration, with maxima of 2.5 seconds and 7.5 seconds respectively. The standard deviation is much smaller for the positive item, which means there was more agreement on the estimated values. On the other side, the distributions for the negative item are shifted to the right towards worse acceleration, the maximum of the lower threshold located around 10 seconds and the maxima of the upper threshold at 13 and 18 seconds. The differences between all curves are strongly significant (compare Table 3).

As Figure (1) shows, acceleration values for cars with 200 hp vary between 6.5 and 10.5 seconds, which agrees surprisingly well with the maxima of the kernel density estimates for the neutral item. Therefore, the interval between 0 and 7.5 may permit the usage of a positive evaluative adverb, whereas values between 10 and 20 seconds may license the usage of the negative evaluation.

As for the three threshold pairs for the simultaneously shown items on the final page of the study, we can compare their distributions with sep-

arately shown items in order to assess the influence of presenting alternative evaluations at the same time. Again, we compare the distributions of the corresponding pairs of items with the same polarity in regard to both thresholds. Only two thresholds are significantly different, namely the neutral upper bound and the positive upper bound. All other thresholds are more or less equivalent. These results indicate that participants do have a good intuition about the intervals that agree with a certain evaluative stance even without presence of alternatives as an orientation point.

One item that is rather special contains the feature *displacement* in a group with a neutral and a strongly positive item with adverb and modifiers. This feature is rather unknown among non-experts of the domain and therefore it is also harder to estimate reasonable values for it. Variance in the data is also really strong, since the relation between *displacement* and *power* depends on many other factors. This is also mirrored in the distributions for lower and upper threshold between the neutral and the positive item. Standard deviation is much larger, there is no significant difference and regression models perform less precise and less confident.

Another hypothesis concerns the effect of modifiers, which we expected to intensify the influential effect of the evaluative adverb. There is no empirical proof of this in our data, however.

The ablation study on the effect of modifiers confirms the findings for evaluative adverbs. Although the bell curves show that for both thresholds, the expected values differ between neutral and modifier-enriched sentence into the same direction as the neutral and adverb-enriched sentences, the effect-size of modifiers seems to be smaller than for evaluative adverbs. One interesting thing to be mentioned here is the influence of dual modifiers (one for each noun phrase of the two opposing car features), when one modifier puts its feature into perspective, such as in one group of items in this ablation study. Participants are supposed to guess acceleration by maximum speed. The modifier-enriched sentences ask for a good acceleration given a "rather average" maximum speed of 200 km/h. The positive sentence with adverb only elicits expectations with best acceleration values, whereas sentences with only modifiers and both adverb and modifiers are nearly indistinguishable and lie in between the neutral and the adverb-enriched
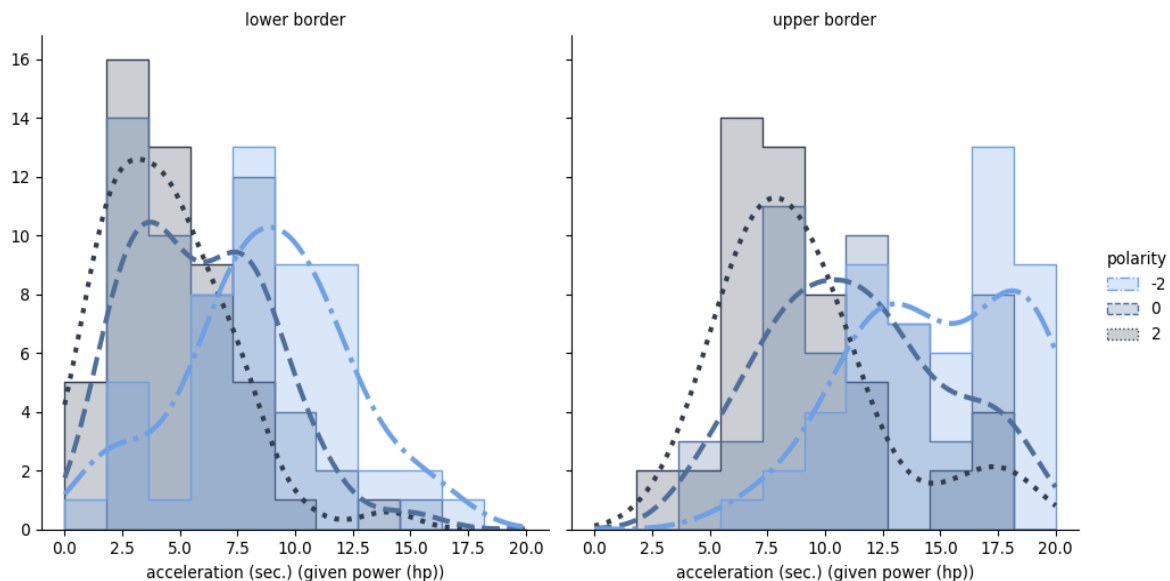
Figure 2: Distribution plots for lower border and upper border estimations on acceleration, given a specific power output value. Lower values on x-axis are better. The dotted line is the positive expression, the dashed line represents the neutral expression and the dash-dot line is the negative polar expression.

| comp | lower border | upper border |
|------|-------------|--------------|
| -2 vs. 0 | 4.2 (p=5.65e-5) | 4.7 (p=7.58e-6) |
| 2 vs: 0 | 2.79 (p=0.006) | 3.11 (p=0.002) |
| 2 vs. -2 | 6.90 (p=4.91e-10) | 7.97 (p=2.68e-12) |

Table 3: Item pairs and their significance scores (residual probability) for curves in Figure 2

sentence. This indicates that modifiers may also delimit the effect size of evaluative adverbs in case they give additional information that narrows down where expectations should be centered.

Once again, for the ablation study, the maxima of the distributions of participants estimations correlate surprisingly well with the real observed data for each item.

Overall, when comparing the maxima of the distribution plots to the subset of the technical database that agrees with the feature that is named in the items (e.g. 200 hp for the items in Table 2), the majority of participants, who are no experts on car technology, have a surprisingly good intuition. Nonetheless data for less known features such as displacement shows that the participants must have enough knowledge about the source feature, otherwise the match between real data and maxima of distribution curves deteriorates.

## 4.3 Bringing empirical data and regression together

The empirical study supports a transparent mapping from regression to the use of evaluative adverbs and modifiers.

Figure (3) depicts a tripartite plot. The upmost part contains the regression model predictions of acceleration given a power output of 200 hp. The light-blue asterisks mark the three model predictions, whereas the dark-blue triangles pointing upwards and downwards to both sides result from adding or subtracting the standard deviation to or from the regression values. The close vicinity of the marks show how close regression models are in prediction in this case. The central part of the plot is a bar plot of cars with 200 hp (+/- 2.5%) and their respective acceleration value. The undermost part contains the kernel density estimation curves for the upper threshold of an item that asks for maximum speed given power output. The left, light-brown curve displays the distribution for the negative item, the central curve represents the neutral item and the rightmost curve displays estimations for the positive item. Non-neutral items contain both adverb and modifiers in this group. As the figure shows, the maximum of the neutral item's curve is neatly aligned with the maximum of the already narrow distribution of real data points as well as the regression values, which visualizes that regression
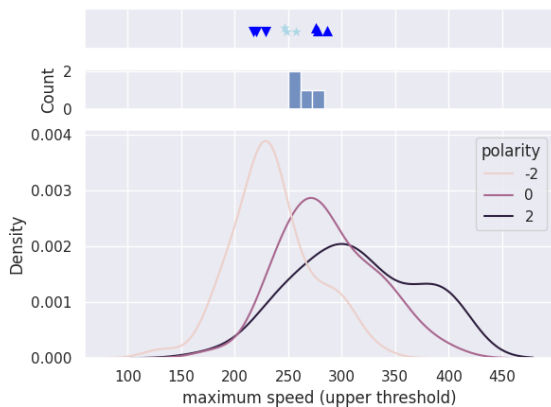
Figure 3: Predicted values (+/- std),
real data and estimations for speed given power output

is very precise and participants have a surprisingly good intuition matching real data. Even more intriguing is the fact that the regression values after modification by either adding or subtracting standard deviation also neatly match the maxima of the curves for the positively and negatively polar items respectively. Across items, one can recognize the pattern that depending on the semantics of the response feature (higher is better or lower is better), non-neutral sentences correspond to the maxima of the respective distributions of estimated values. In regard to thresholds for binary decisions on generating evaluative content, this would mean that if the regression value deviates more than standard deviation from the real value, evaluative content is licensed according to the polarity of the deviation.

This study makes no statement about which specific adverb and/or modifier should be used. It shows that regression values combined with standard deviation are a good starting point for deciding when to become negatively or positively evaluative in a description, but the subtle meaning differences between, e.g., *surprisingly*, *very surprising*, and *totally surprising* were not a topic of this paper. Such a fine-grained lexicalization process, which is a downstreamed module in a pipelined NLG system, requires access to grammatical, semantic and discourse-related constraints, which is outside the scope of the model presented in this paper.

## 5 Conclusion

For controllable generation of evaluative adverbs in data2text NLG, we devised a learning-based, generalizable approach to approximate thresholds for binary decisions on the presence of evaluative lan-

guage and validated our findings with an empirical study on the expectations raised by comprehending evaluative adverbs. The elicited data supports the adequacy of the model and shows a surprisingly good match between regression predictions, real data and human estimations. In a pipelined neural generation system, these learning-based models permit determining generation of evaluative language at an early stage in document planning and therefore improve controllability of evaluative content before applying transformers for surface realisation.

**Supplementary Materials Availability Statement:** Source code of regression models, source code of the web application for the empirical study, empirical data collected during the study as well as source code for analysis of the data including functions for reproducing all graphs and figures in this paper are available online on Github: `https://github.com/MMLangner/evalAdvInData2TextNLG/`. The database containing technical data of cars is proprietary, which is why we are not allowed to distribute it. Please contact the authors for information on how to obtain the technical database in the required format from the ADAC.

## References

Simone Balloccu, Steffen Pauws, and Ehud Reiter. 2020. A NLG framework for user tailoring and profiling in healthcare. In *Proceedings of the First Workshop on Smart Personal Health Interfaces co-located with 25th International Conference on Intelligent User Interfaces, SmartPhil@IUI 2020, Cagliari, Italy, March 17, 2020*, volume 2596 of *CEUR Workshop Proceedings*, pages 13–32. CEUR-WS.org.

Farah Benamara, Maite Taboada, and Yannick Mathieu. 2017. Evaluative language beyond bags of words: Linguistic insights and computational applications. *Computational Linguistics*, 43(1):201–264.

Fiorella de Rosis and Floriana Grasso. 2000. *Affective Natural Language Generation*, pages 204–218. Springer Berlin Heidelberg, Berlin, Heidelberg.

Ondřej Dušek and Zdeněk Kasner. 2020. Evaluating semantic accuracy of data-to-text generation with natural language inference. *arXiv*, pages 131–137.

Michael Elhadad. 1991. Generating adjectives to express the speaker's argumentative intent. In *AAAI Conference on Artificial Intelligence*.

Thiago Castro Ferreira, Chris van der Lee, Emiel van Miltenburg, and Emiel Krahmer. 2020. Neural data-to-text generation: A comparison between pipeline

and end-to-end architectures. *EMNLP-IJCNLP 2019 - 2019 Conference on Empirical Methods in Natural Language Processing and 9th International Joint Conference on Natural Language Processing, Proceedings of the Conference*, pages 552–562.

Albert Gatt and Emiel Krahmer. 2018. Survey of the state of the art in natural language generation: Core tasks, applications and evaluation. *Journal of Artificial Intelligence Research*, 61(c):1–64.

Daniel Gutzman. 2015. *Use-Conditional Meaning: Studies in Multidimensional Semantics*. Oxford Studies in Semantics and Pragmatics. Oxford University Press, Oxford.

Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2022. Survey of hallucination in natural language generation. *CoRR*, abs/2202.03629.

Maurice Langner and Ralf Klabunde. 2022. Realizing a denial of expectation in pipelined neural data-to-text generation. In *Proceedings of the 6th Workshop on Advances in Argumentation in Artificial Intelligence 2022 co-located with the 21st International Conference of the Italian Association for Artificial Intelligence (AIxIA 2022), Udine, Italy, November 28, 2022*, volume 3354 of *CEUR Workshop Proceedings*. CEUR-WS.org.

Saad Mahamood. 2021. Reproducing a comparison of hedged and non-hedged NLG texts. In *Proceedings of the 14th International Conference on Natural Language Generation*, pages 282–285, Aberdeen, Scotland, UK. Association for Computational Linguistics.

Saad Mahamood, Ehud Reiter, and Chris Mellish. 2007. A comparison of hedged and non-hedged nlg texts. In *Proceedings of the Eleventh European Workshop on Natural Language Generation*, ENLG '07, page 155–158, USA. Association for Computational Linguistics.

Christopher Potts. 2005. *The Logic of Conventional Implicatures*. Oxford Studies in Theoretical Linguistics. Oxford University Press, Oxford.

Ehud Reiter and Robert Dale. 2000. *Building Natural Language Generation Systems*. Cambridge University Press, Cambridge.

Sashank Santhanam and Samira Shaikh. 2019. Emotional neural language generation grounded in situational contexts. *CoRR*, abs/1911.11161.

Xiaoyu Shen, Ernie Chang, Hui Su, Jie Zhou, and Dietrich Klakow. 2020. Neural data-To-Text generation via jointly learning the segmentation and correspondence. *arXiv*, (2019):7155–7165.

Ishika Singh, Ahsan Barkati, Tushar Goswamy, and Ashutosh Modi. 2020. Adapting a language model for controlled affective text generation. *CoRR*, abs/2011.04000.

Sam Wiseman, Stuart Shieber, and Alexander Rush. 2017. Challenges in data-to-document generation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2253–2263, Copenhagen, Denmark. Association for Computational Linguistics.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multilingual pre-trained text-to-text transformer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498, Online. Association for Computational Linguistics.

# The Next Chapter: A Study of Large Language Models in Storytelling

**Zhuohan Xie**     **Trevor Cohn**[*]     **Jey Han Lau**
School of Computing and Information Systems
The University of Melbourne
zhuohanx@student.unimelb.edu.au, {t.cohn, laujh}@unimelb.edu.au

## Abstract

To enhance the quality of generated stories, recent story generation models have been investigating the utilization of higher-level attributes like plots or commonsense knowledge. The application of prompt-based learning with large language models (LLMs), exemplified by GPT-3, has exhibited remarkable performance in diverse natural language processing (NLP) tasks. This paper conducts a comprehensive investigation, utilizing both automatic and human evaluation, to compare the story generation capacity of LLMs with recent models across three datasets with variations in style, register, and length of stories. The results demonstrate that LLMs generate stories of significantly higher quality compared to other story generation models. Moreover, they exhibit a level of performance that competes with human authors, albeit with the preliminary observation that they tend to replicate real stories in situations involving world knowledge, resembling a form of plagiarism.

## 1 Introduction

Automatic story generation poses a significant challenge as it requires more than just individually coherent sentences. A good story should exhibit a natural flow, adhere to commonsense logic, and be captivating to the reader. In recent times, the prevailing approach in story generation involves fine-tuning pre-trained language models (PLMs) like GPT-2 (Radford et al., 2019) or BART (Lewis et al., 2020) on specific datasets. These models generally excel in generating fluent sentences, devoid of glaring grammar issues. However, they often struggle to construct a coherent story that adheres to commonsense and fails to create an engaging narrative (See et al., 2019; Guan et al., 2021a). To overcome these challenges, state-of-the-art (SOTA) story generation models integrate higher-level features, such as plots and commonsense knowledge.

Prompt-based learning (Liu et al., 2021) is a recent paradigm specifically tailored for large language models (LLMs) that possess in-context learning capabilities (Brown et al., 2020; Zhao et al., 2023). In contrast to the conventional "pre-train and fine-tune" approach, which necessitates a substantial amount of data for fine-tuning, prompt-based learning enables LLMs to learn a task by providing them with multiple examples as a "prompt", eliminating the need for gradient-based fine-tuning (Liu et al., 2021). Recently, LLMs have demonstrated remarkable performance across various language generation tasks, with notable attention on models such as ChatGPT and GPT-4 (Qin et al., 2023; Liu et al., 2023b; OpenAI, 2023). For instance, a comparative analysis by Qin et al. (2023) highlighted the superior performance of fine-tuned LLMs over smaller pre-trained models in zero-shot scenarios for tasks like dialogue and summarization. It is worthwhile to note though, that story generation was not specifically examined in their experiments.

This paper aims to address this research gap by conducting a comprehensive evaluation of automatic story generation. Specifically, we compare the performance of LLMs, with a particular focus on GPT-3, using prompt-based learning, against SOTA models. We compare generated stories in terms of various automatic evaluation metrics from lexical and semantic matching ones to recently proposed model-based ones. We follow the best practice in literature to conduct rigorous human evaluations including both crowdworkers from Amazon Mechanical Turk and in-house judges, and assess story quality at a fine-grained level, such as coherence and logicality. To summarise, our contributions are:

- We conduct an empirical comparison between GPT-3 and other SOTA techniques for open-ended story generation on three different corpora that differ in style, register, and length.

---

[*]Now at Google DeepMind.

- We test with a wide variety of automatic story evaluation metrics, and find that recent model-based ones work better, consistent with the literature.

- We conducted experiments using two types of annotators: crowdworkers and in-house judges, to assess the quality of stories on various aspects. The results obtained from two groups are consistent. We release this annotated resource as a testbed for developing new automatic metrics in story generation tasks.[1]

- Our experimental findings provide comprehensive evidence that the stories produced by GPT-3 exhibit significant improvement compared to SOTA techniques, and are comparable to stories authored by humans across various aspects.

- We conduct a preliminary study on story plagiarism and find that GPT-3 tends to (soft) "plagiarise" real stories when generating news, even though it does not directly copy the source text, raising further questions as to what extent GPT-3 recycles stories in its memory rather than generating new narratives.

## 2 Related Work

**Story Generation**   See et al. (2019) find that fine-tuned GPT-2 can already generate stories with fluent sentences, but more attentions are needed to incorporate commonsense and higher-level story planning. Most works then use PLMs such as GPT-2 or BART as the backbone and incorporate higher level features to aid the generation process. Specifically, Rashkin et al. (2020); Goldfarb-Tarrant et al. (2020); Tan et al. (2021) construct a storyline to guide the generation process. Guan et al. (2021a); Yu et al. (2021); Xie et al. (2021) incorporate intersentence relationships such as coherence and discourse relationships into the generation process. Guan et al. (2020); Peng et al. (2021) explore using external knowledge such as commonsense for story generation. Xu et al. (2020); Ammanabrolu et al. (2021) combine storyline planning and commonsense reasoning.

There are also studies which explore the use of GPT-3 for story generation. For example, Clark et al. (2021) conducts a Turing test between GPT-3 generated and human-written stories and Lucy

[1]https://github.com/ZhuohanX/TheNextChapter

and Bamman (2021) probe for gender and representation bias in GPT-3 generated stories. These studies, however, do not provide a systematic evaluation that assesses GPT-3 against the SOTA story generation models.

**Story Evaluation**   Automatic story evaluation is admittedly a challenging task, and the lack of standardized evaluation metrics has somewhat impeded progress of story generation (Guan et al., 2021b). Human evaluation is usually considered as the gold standard for story quality evaluation, but it is expensive and time-consuming (Guan and Huang, 2020) and it can not capture diversity (Hashimoto et al., 2017). Subsequently, several automated evaluation metrics are introduced as alternative measures to evaluate the quality (the degree of readability) and diversity (the extent of variation) of the generated stories. For quality, most metrics measure lexical overlap between strings (Papineni et al., 2002; Lin, 2004; Tan et al., 2021) or semantic similarity by comparing embedding of models (Zhao et al., 2019; Zhang et al., 2020) between generated stories and their human references. Recently, learning (Sellam et al., 2020) and generation (Yuan et al., 2021) based methods are explored and they are based on pre-trained language models such as BERT (Devlin et al., 2019) and BART. Nevertheless, these evaluation metrics are limited in that they provide a single score to indicate the overall quality of the story, and few metrics are specifically designed to assess specific aspects such as logicality (the adherence to commonsense) or interestingness (the level of reader engagement) (Chhun et al., 2022).

## 3 Experimental Setup

### 3.1 Story Generation Models

To ensure a comprehensive comparison, we conducted extensive experiments involving GPT-3 and a wide range of SOTA story generation models.

In our experiments, we utilized the largest initial version of GPT-3, namely text-davinci-001, which was initially introduced in June 2020 and comprises 175B parameters. It is perhaps worth noting that this model was considered the most powerful at the time of our experiment (March 2022), although subsequent models like GPT-4 have since been released, boasting even greater capabilities. As such, the results we report here can be interpreted as a "lower bound" of LLM's story generation performance. To adapt GPT-3 to the story domain without

explicit fine-tuning, we employed a prompt-based learning approach. We selected a small number of stories, typically 2 or 3, to serve as exemplars for GPT-3 in the target domain.

For SOTA story generation models, we use 1) knowledge enhanced based models: **KGGPT2** (Guan et al., 2020) and **HINT** (Guan et al., 2021a); 2) storyline planning based model: **PROGEN** (Tan et al., 2021); and 3) **MTCL** (Xu et al., 2020) that combines both storyline planning and commonsense reasoning. We also fine-tune **BART** as an additional baseline. For consistency, all models use nucleus sampling (Holtzman et al., 2020) with $p = 0.95$ as the decoding method. We summarise these models in Table 1, and more details can be found in Appendix A.

### 3.2 Story Datasets

The most popular story dataset is ROCStories (ROC) (Mostafazadeh et al., 2016), which is composed of short commonsense stories and is used by most story generation works. There are also more difficult and longer story datasets, such as WritingPrompts (WP) (Fan et al., 2018) and CNN News (CNN) (Hermann et al., 2015) which are composed of fictional and news stories (two different domains). In our experimental setup, we utilized all three datasets. The ROC dataset was used to evaluate the generation of short stories comprising 5 sentences. The WP dataset was employed to assess medium-length stories, which were trimmed down to 10 sentences. Lastly, the CNN dataset was utilized to evaluate the generation of long stories, and each story have around 20 sentences. For further details about these datasets, please refer to Appendix B.

Whenever possible we evaluate all models on each story dataset. However, this is sometimes infeasible because some models are designed to work on a particular dataset and thus cannot be adapted to other datasets easily. Moreover, we focus on conditional story generation in this work, this means there is some *context* upon which we generate the stories (details below).

**ROC**    We evaluate all models in this dataset. The context we use to generate stories is the first sentence, and so the models are trained to generate the last 4 sentences. Evaluation results are computed over 800 generated stories using randomly sampled leading sentences from the test partition.

**WP**    We assess HINT, PROGEN, GPT-3 and BART on this dataset. The context is a short paragraph ("prompt") that describes the idea of the story. We randomly sample 1000 prompts from the test partition for automatic evaluation.

**CNN**    We only run GPT-3, BART, PROGEN on CNN, as HINT is developed for ROC and WP originally and it does not work well when applied to CNN. Stories of CNN are generated conditioned on the news titles. We randomly sample 600 titles from the test partition for automatic evaluation.

## 4    Automatic Evaluation

### 4.1    Evaluation Metrics

We use two types of automatic evaluation metrics: 1) reference-based metrics, where we compare the generated stories to human reference stories based on the same conditioning context; and 2) reference-free metrics, where we assess the quality of the stories directly.

#### 4.1.1    Reference-based Metrics

Most reference-based metrics measure the lexical or semantic closeness between generated stories and their human references. We experiment with metrics based on string based matching (CBL, MSJ) and embedding based matching (BES) and a learning based metric (BRT), to assess the quality of generated stories. We also use a recall based metric (BBL) to assess the diversity of generated stories. Specifically, **Corpus BLEU (CBL)** computes the average BLEU scores (Papineni et al., 2002) for each generated story against all human references (Caccia et al., 2020; Xie et al., 2021). **MS-Jaccard (MSJ)** measures lexical overlap by computing the n-gram overlap between generated and referenced stories using the Jaccard index (Alihosseini et al., 2019). **BERTScore (BES)** measures the maximum similarity of each token's contextual embedding between generated and referenced stories (Zhang et al., 2020). **BLEURT (BRT)** is trained on synthetic data to predict a similarity score between generated and referenced stories (Sellam et al., 2020). **Backward BLEU (BBL)** computes the coverage of n-grams in the reference stories against the set of generated stories (Shi et al., 2018).[2]

---

[2] We use BLEU4 for CBL and BBL; 4-grams overlap for MSJ; roberta-large model for BES; bert-base-128 for BRT.

| Model | Backbone | Size | Method | Story Datasets |
|-------|----------|------|--------|----------------|
| GPT-3 | text-davinci-001 | 175B | Prompt-based learning with several examples from the story dataset (3 for ROC and WP and 2 for CNN) | ROC, WP, CNN |
| KGGPT2 | GPT-2 small | 124M | Fine-tuned on commonsense data before more fine-tuning with auxiliary classification tasks | ROC |
| PROGEN | BART large | 400M | Three-stage generation where at each stage a fine-tuned BART generates stories based on word importance in the story datasets | ROC, WP, CNN |
| MTCL | GPT-2 small BERT large | 124M 336M | (1) a GPT-2 model to generate keywords; (2) a BERT model to rank retrieved knowledge triples; and (3) a second GPT-2 model that takes top-ranked knowledge triples and context as input for story generation | ROC |
| HINT | BART base | 140M | BART is first fine-tuned on BookCorpus with additional objectives to learn internal structure in a story and then further fine-tuned on the story datasets | ROC, WP |
| BART | BART large | 400M | Baseline model that is fine-tuned on the story datasets using a standard language modelling objective | ROC, WP, CNN |

Table 1: The backbone ("Backbone") of the story generation models and their number of parameters ("Size"). "Story Datasets" indicates which datasets are used to generate stories for a particular model. KGGPT2 and MTCL stories are obtained from the original authors; for PROGEN and HINT we re-run the implementation provided by the authors.

### 4.1.2 Reference-free Metrics

Reference-free metrics evaluate generated stories without comparing them to their human-authored references. We experiment with diversity metrics based on intra-story (D-3, LR-n) and inter-story diversity (SBL). We also compute negative log-likelihood from BART of a story conditioned on the context (BAS) for relatedness, and story length in terms of words (LEN) for complexity.

Specifically, **Lexical Repetition (LR-n)** computes the average percentage of 4-grams appearing at least $n$ times in the generated stories (Shao et al., 2019). **Distinct-3 (D-3)** computes the average ratio of distinct 3-grams to all 3-grams (Li et al., 2016). **Self-BLEU (SBL)** measures inter-story diversity that computes the average BLEU score of each generated story using all generated stories as reference (Zhu et al., 2018). **BARTScore (BAS)** computes generative likelihood of a story conditioned on the context (i.e., leading sentence for ROC, prompt for WP and title for CNN) to measure the extent to which a generated story relates to its condition (Yuan et al., 2021).[3] **Length (LEN)** measures the average length of the generated stories, which is used as a rough indicator of generation complexity.

---

[3]We set $n = 3/8/8$ for ROC, WP and CNN respectively and use BLEU4 for SBL. We use the "PARA" version of BART and direction as "from source to hypothesis".

### 4.2 Results

Table 2 and Table 3 present the reference-based and reference-free evaluation results, respectively. At a glance, these metrics do not appear to agree with each other even though some of them are designed to evaluate the same aspect (e.g., the best model in terms of fluency/coherence or diversity is different depending on the metric). Overall, GPT-3 seems to have weaker performance than most of other models in terms of quality (CBL and MSJ) and diversity (BBL, SBL, D-3 and LR-n) metrics.

However, when we look at recent model-based metrics (BERTScore, BLEURT and BARTScore), GPT-3 appears to be a much better model (a finding we will return to when we look at human evaluation results). Interestingly, we notice that human written stories have very poor performance in terms of BARTScore (BAS). We suspect BARTScore may exhibit a bias towards machine-generated stories, as the metric primarily evaluates quality based on the generative likelihood of a sequence. Machine-generated stories are specifically designed to maximize this likelihood, while human-authored stories often incorporate distinct elements, such as surprising or creative word choices (Holtzman et al., 2020). In general, all models are capable of generating stories of the appropriate length, except for GPT-3 in the CNN dataset. GPT-3 in the CNN dataset exhibits difficulty in generating stories longer than 150 words, whereas human-written

| | Model | Flu./Coh. | | | | Div. |
|---|---|---|---|---|---|---|
| | | **CBL** ⇑ | **MSJ** ⇑ | **BES** ⇑ | **BRT** ⇓ | **BBL** ⇑ |
| **ROC** | GPT-3 | 27.2 | 11.6 | 86.6 | **8.6** | 24.0 |
| | KGGPT2 | 33.5 | **15.0** | **87.0** | 9.5 | **25.6** |
| | PROGEN3 | 26.6 | 14.6 | 86.7 | 9.7 | 25.0 |
| | MTCL | 31.4 | 14.2 | 86.9 | 9.7 | 24.0 |
| | HINT | **39.6** | 13.7 | **87.0** | 8.6 | 24.6 |
| | BART | 27.5 | 14.7 | 86.8 | 9.5 | 25.1 |
| **WP** | GPT-3 | 28.6 | 12.3 | **81.6** | **11.7** | 24.4 |
| | PROGEN3 | 32.3 | **16.4** | 81.4 | 13.3 | **27.6** |
| | HINT | **45.5** | 12.8 | 80.8 | 12.1 | 23.7 |
| | BART | 32.6 | 16.2 | 81.4 | 13.0 | 27.2 |
| **CNN** | GPT-3 | **33.2** | 11.0 | **83.5** | **7.5** | 19.8 |
| | PROGEN3 | 29.6 | **14.8** | 82.2 | 9.3 | **26.2** |
| | BART | 29.1 | 14.7 | 82.2 | 9.8 | 25.7 |

Table 2: Reference-based Evaluation Results. CBL, MSJ, BES and BRT evaluate the closeness between the generated stories and the whole test reference data as an indicator of general fluency (Flu.) and coherence (Coh.). BBL focus on the recall of generated stories as an indicator of diversity (Div.). ⇑: higher is better; ⇓: lower is better. BRT values are negated here.

stories typically consist of around 500 words on average. Considering the overall assessment using various automatic metrics, there is no single winner that consistently outperforms other models.

## 5 Human Evaluation

To obtain a comprehensive assessment of the generated stories, we recruit human annotators to evaluate their quality. In order to gain insights into consistency, we employed both crowdsourced workers and in-house annotators. This approach allows us to gather diverse perspectives and obtain a more nuanced understanding of the story quality.

### 5.1 Crowdsource Annotation

We first collect human judgements using the Amazon Mechanical Turk (AMT) platform.[4] Following the approach suggested by Karpinska et al. (2021), we assessed four aspects, namely fluency, coherence, relatedness, and interestingness. Additionally, we introduced a new aspect called logicality, which assesses the extent to which the story complies with commonsense. Each of these five aspects is evaluated on an ordinal scale ranging from 1 (worst) to 5 (best). We randomly sample 20 conditional contexts (e.g., titles) from each dataset and collect stories generated by all models for human evaluation. Each story (including human-written

[4] https://requester.mturk.com/

| | Model | Div. | | | Rel. | Com. |
|---|---|---|---|---|---|---|
| | | **SBL** ⇓ | **D-3** ⇑ | **LR-n** ⇓ | **BAS** ⇓ | **LEN** ⇑ |
| **ROC** | GPT-3 | 38.5 | 67.7 | 39.1 | **4.2** | 47.3 |
| | KGGPT2 | 41.9 | 67.2 | 51.9 | 4.6 | 38.4 |
| | PROGEN3 | **30.0** | 76.9 | 39.5 | 5.0 | 40.9 |
| | MTCL | 39.4 | 69.6 | 44.4 | 4.9 | **49.7** |
| | HINT | 55.1 | 54.3 | 68.1 | 4.3 | 35.8 |
| | BART | 30.5 | 77.4 | 37.8 | 5.0 | 40.6 |
| | human | 33.1 | **80.2** | **35.8** | 5.2 | 40.3 |
| **WP** | GPT-3 | 37.5 | 69.6 | 9.7 | 4.3 | 120.6 |
| | PROGEN3 | 35.2 | 77.2 | 2.6 | 5.4 | 136.9 |
| | HINT | 64.1 | 33.9 | 67.4 | **4.1** | 119.0 |
| | BART | 35.3 | 77.5 | 1.6 | 5.4 | 129.2 |
| | human | **27.1** | **83.7** | **1.5** | 5.7 | **150.0** |
| **CNN** | GPT-3 | **26.5** | 82.9 | 9.8 | **4.4** | 147.3 |
| | PROGEN3 | 28.9 | 82.3 | 2.3 | 5.2 | 395.8 |
| | BART | 27.9 | 83.2 | **0.8** | 5.2 | 374.1 |
| | human | 27.3 | **83.8** | 6.3 | 5.4 | **498.6** |

Table 3: Reference-free Evaluation Results. SBL measures inter-story diversity by assessing differences between different stories while D-3 and LR-n (3 for ROC, 8 for WP and CNN) focus on repetition n-grams within the same story. We also include LEN as an indicator of story complexity (Com.). We compute BAS of story given condition for story relatedness (Rel.).

one) is judged by 3 annotators, and so we have annotations for 320 stories in total (140/100/80 for ROC, WP and CNN, respectively). Amazon Qualification requirements on AMT and question details can be found in Appendix C. Quality control details can be found in Appendix D.

Table 4 presents the human evaluation results. Overall, GPT-3 generates stories that are consistently of higher quality than those generated by other SOTA models. To understand whether the difference is significant, we perform a paired t-test by comparing GPT-3 to other models (including human) and find that in most cases these results are significant with $p$-value $< 0.05$ ('*' in the table). Compared with human authors, GPT-3 appears that it is generating stories that are just as good as (ROC) or better than (WP and CNN) human authors, confirming the findings of Clark et al. (2021). For WP, in particular, human stories are trimmed to the first 10 sentences (data pre-processing for training the story generation models). This abruptly shortens the stories so they might not provide a proper conclusion, and inevitably are penalised (see examples in Appendix K). For CNN, GPT-3 appears to be "plagiarising" real stories, where many story elements are not a product of creative generation but details copied from real news stories

| | Model | Flu. | Coh. | Rel. | Log. | Int. |
|---|---|---|---|---|---|---|
| ROC | GPT-3 | **4.40** | 4.43 | 4.37 | 4.37 | 3.57 |
| | KGGPT2 | 3.90* | 3.48* | 3.53* | 3.00* | 2.62* |
| | PROGEN3 | 3.88* | 3.45* | 3.37* | 2.95* | 2.57* |
| | MTCL | 3.55* | 3.12* | 3.18* | 2.73* | 2.42* |
| | HINT | 3.90* | 3.27* | 3.33* | 3.12* | 2.58* |
| | BART | 3.92* | 3.38* | 3.48* | 3.03* | 2.60* |
| | human | 4.22 | **4.58** | **4.42** | **4.48** | **3.77** |
| WP | GPT-3 | **4.37** | **4.67** | **4.28** | **4.48** | **3.47** |
| | PROGEN3 | 3.45* | 3.08* | 2.35* | 2.57* | 1.98* |
| | HINT | 3.32* | 2.63* | 2.02* | 2.25* | 1.77* |
| | BART | 3.42* | 2.73* | 2.08* | 2.27* | 1.87* |
| | human | 4.13* | 4.22* | 3.05* | 3.75* | 2.97* |
| CNN | GPT-3 | **4.22** | **4.52** | **4.58** | **4.60** | 3.20 |
| | PROGEN3 | 3.63* | 3.32* | 3.30* | 3.22* | 2.28* |
| | BART | 3.58* | 3.37* | 3.30* | 3.27* | 2.17* |
| | human | 4.10 | 4.10* | 4.23* | 4.18* | **3.72*** |

Table 4: Crowdsource Human Evaluation Results. We calculate the average score of models for each aspect: fluency (Flu.), coherence (Coh.), relatedness (Rel.), logicality (Log.) and interestingness (Int.). Model scores that are marked with ∗ indicate the performance difference between the model and GPT-3 is significant.

| | Model | Flu. | Coh. | Rel. | Log. | Int. |
|---|---|---|---|---|---|---|
| ROC | GPT-3 | **4.78** | **4.73** | **4.50** | **4.82** | **3.37** |
| | KGGPT2 | 4.52* | 3.67* | 3.57* | 3.47* | 2.50* |
| | PROGEN3 | 4.27* | 3.47* | 3.78* | 3.23* | 2.48* |
| | MTCL | 4.27* | 3.27* | 3.45* | 3.15* | 2.37* |
| | HINT | 4.38* | 4.03* | 3.38* | 3.70* | 2.38* |
| | BART | 4.37* | 3.95* | 3.85* | 3.53* | 2.70* |
| | human | 4.52* | 4.38* | 4.22 | 4.32* | 3.18 |
| WP | GPT-3 | **4.57** | **4.65** | **4.08** | **4.22** | **3.82** |
| | PROGEN3 | 3.55* | 3.03* | 2.23* | 2.57* | 2.45* |
| | HINT | 3.60* | 2.72* | 2.07* | 2.68* | 2.08* |
| | BART | 3.45* | 2.77* | 2.08* | 2.38* | 2.30* |
| | human | 4.05* | 4.07* | 3.73 | 3.87* | 3.78 |
| CNN | GPT-3 | **4.50** | **4.33** | **4.48** | **4.40** | **3.45** |
| | PROGEN3 | 3.80* | 3.45* | 3.63* | 3.45* | 2.52* |
| | BART | 3.73* | 3.25* | 3.58* | 3.32* | 2.57* |
| | human | 4.22* | 4.00* | 4.35 | 4.13* | 3.22 |

Table 5: In-house Human Evaluation Results.

define what makes a narrative interesting.

## 5.2 In-house Annotation

We next recruit university volunteers to collect in-house judgements.[5] We ask them to evaluate the same 5 aspects using the same scale. We sample 20 *disjoint* conditional contexts from each dataset for story generation here, as we are interested to test the robustness of our previous findings (with different workers and set of stories). As with crowdsource annotation, each story is also judged by 3 annotators. Details of the agreement between annotators can be found in Appendix G.

Table 5 presents the scores of story quality from in-house annotators. Interestingly, the *magnitude* of the in-house scores are generally somewhat higher than the crowdworker scores (across all metrics and datasets and models). We hypothesise that this may be because our in-house workers are more "tolerant" to mistakes as they have been exposed to machine-generated text more compared to crowdworkers. That said, the overall findings are consistent between the two groups of annotators: 1) GPT-3 is the best story generation model and outperforms both SOTA models and human stories; 2) The SOTA models do well in fluency, but poorly in most other aspects (interestingness worst); and 3) The SOTA models face notable challenges specifically in WP, as evidenced by their poor coherence, relatedness, logicality scores compared to other domains.

When comparing the results of automatic met-

(Section 6). Another reason could be that GPT-3 stories are much shorter than those generated by other models and human authors (150 vs. 300-400 words; Table 3), which makes them easier to read and thus leads to better scores. Note that this is a downside of GPT-3 where it is difficult to get it to generate long stories (Section 7).

When considering the various aspects for SOTA models, including KGGPT2, PROGEN3, MTCL, HINT, and BART, these models exhibit strong performance in terms of fluency, with scores consistently exceeding 3.5 in most cases. This indicates that the models can generate sentences that are natural and fluent. However, coherence performance differs depending on the dataset. Most models perform well on the ROC and CNN datasets, while they tend to struggle on WP, with coherence scores falling below 3.1. The observation that these models struggle with shorter WP stories compared to longer CNN stories might be because the PLMs that they are built on are mostly trained on web data which contains plenty of news articles. For relatedness, logicality and interestingness, we see a similar trend where the models perform best in ROC and worst in WP. We also observe a consistent decrease in performance from relatedness to logicality and interestingness, suggesting that the models particularly struggle to generate interesting and sensible stories. Interestingness is perhaps the most difficult aspect to optimise, as it is difficult to

328

rics (Section 4.2) to the human evaluation results, a notable discrepancy emerges, leading to a different conclusion regarding the performance of GPT-3 and the identification of a clear "best" story generation model. That said, if we consider only model-based metrics such as BERTSCore, BLEURT for fluency/coherence, and BARTScore for relatedness, a more aligned conclusion can be drawn, suggesting these metrics may be more reliable (though the trend is still less conclusive compared to human evaluation results). This observation is consistent with recent literature, which highlights the better correlation of modern model-based metrics with human evaluations (Chiang and Lee, 2023; Ke et al., 2023; Xie et al., 2023).

## 6 Plagiarism

Considering the strong performance of GPT-3 on story generation, we next provide a preliminary investigation to understand the extent to which GPT-3 copies from its training data.

**IThenticate**  We use iThenticate[6] — a professional plagiarism detection software that has comprehensive coverage over online articles — to assess how much GPT-3 plagiarises. We include only the generated content (without the condition) when checking for plagiarism. Results show that there is no strong plagiarism: similarity score for ROC, WP and CNN is 4%, 3% and 14% respectively. This is consistent with findings from McCoy et al. (2021) that the language models do not simply memorize, but combine familiar parts into novel way. However, iThenticate looks for lexical overlap to identify plagiarism. A more subtle form of plagiarism is one that copies the ideas without parroting the words (Lee et al., 2022), motivating us to perform manual check.

**Manual Check**  We randomly sample 10 GPT-3 generated stories for WP and CNN respectively and manually check whether they copy from web articles by searching for keywords, named entities and long phrases. For WP, in most instances we do not find web articles that resemble our stories, except for one instance where it copies a story from Wikipedia (noting that it is not a direct word-to-word copy). For CNN, we find 7 out of 10 instances are based on real stories from news articles, where most details (e.g., names, incident origins, sentence penalty) are faithful to the original stories. These

results highlight that GPT-3 can recall details from training data, calling into question how much of its generation is creative. More details can be found in Appendix M.

## 7 Discussion

In this paper, we only experiment with a simple prompt-based learning approach to adapt GPT-3 to the different story datasets and one can argue more prompt engineering can be explored (Liu et al., 2022; Lu et al., 2022; Mishra et al., 2022). We contend, however, that our goal here is not to propose a better prompt-based method to do story generation with GPT-3, and ultimately even with our simple approach we find that GPT-3 generates high quality stories, suggesting that with additional prompt engineering it can do better. We next discuss *qualitatively* some of the strengths and drawbacks of GPT-3 in story generation.

### 7.1 Strengths

Unlike other conditional generation tasks such as machine translation or summarisation, where the input context is rich in information and the goal is to translate or compress the input information, story generation works in the reverse manner where the model needs to "hallucinate" new information and details given a succinct context. This means that in order to do the task well, having strong world knowledge is important. Reading some of the GPT-3 stories, we observe GPT-3's advantage in this, particularly in the WP dataset where some of the prompts require niche knowledge about characters. In Table 6 we show an example in WP where the prompt is *cthulhu calls your cell , he needs to be picked up*, where *cthulhu* a fictional cosmic entity, and only GPT-3 is able to produce a coherent story and the SOTA models struggle.

### 7.2 Drawbacks

Even though GPT-3 demonstrates excellent generative capability and outperforms SOTA models significantly, we still find GPT-3 has many generation errors that can be improved.

**Story length**  GPT-3 has a parameter to control the maximum number of generated tokens but does not provides a way to control the minimum number of tokens. As one can see from Table 3, GPT-3 can not generate stories longer than 150 words for CNN, even though the prompts have long stories. We also attempted to encourage longer stories by

| Model | Story |
|-------|-------|
| GPT-3 | " yes , lord cthulhu ? " you answer the phone . " i need you to pick me up from the airport . " " of course , my lord . when will you be arriving ? " " in about an hour . i 'll text you the flight number . " you hang up and start getting ready . you know cthulhu is a demanding deity , but you do n't mind . you enjoy serving him . |
| PROGEN3 | **but everything 's different ! i miss my child in there ...** *what are you talking about ! you have to be mad . you must have misheard ! ! you come and get your mommy ... " momma always grabbed him before they climbed into the bathtub . after that happened they were always late . * bam . * wait , it 's never like this . ctlulhu always puts up a fight . i repeat : it 's never like this ! it 's always for the umpteenth time , now or never . *the phone rings . i answer it and ctlulhu 's voice booms out . |
| HINT | " sir , i 'm afraid i 'm late for work . " " i 'm sorry sir , " the technician said , " but i have to get my cell phone back . " " sir , you have to pick up this cell . " " yes sir , i do . " " and how are you going to find it ? " " i 'm just going to have to pick it up . " " well , sir , " the technician said , " it 's just that i do n't have the time to pick it up . " " why not ? " " sir , " the technician said , " i do n't have the time . " " i 'm just going to make sure it 's picked up . " |
| BART | " ah it ' s here . i thought you wanted to leave , a phone call . * " my coworkers looked at me confused . " i wasn ' t in the mood to go running around scared . but it was important . you need to pick up on your cravings . the kid-ta-bye , gabriel , payliver , cilia de los angeles . and four people ' s fault . " he was right about the fifth , at least my attention was fixed on it . three girls were outside screaming in the street , and a family of four was picking up one of them . the two men yelled and approached as they walked to the mailbox , visibly distressed . |

Table 6: Stories generated by GPT-3, HINT, BART and PROGEN3 given the same title "cthulhu calls your cell , he needs to be picked up ."

adding specific instructions as part of the prompt of GPT-3, but this did not work.

**Null generation**  Occasionally GPT-3 decides to generate no output. This is usually not an issue, since this can be solved by forcing it to generate again, although it is unclear why this occurs.

**Direct copy**  Besides the soft plagiarism issue (Section 6), GPT-3 does occasionally copy long chunks of text, e.g., the title or prompt in the story.

**Multilingual**  GPT-3 sometimes generates stories in languages other than English, despite the given prompts always being in English. In terms of statistics, out of 1000 generations we find 14 non-English stories (5 Chinese, 4 German, 1 Japanese, 1 French, 1 Russian, 1 Norwegian Nynorsk and 1 mixture of Chinese and English). Interestingly, in most of these cases the stories are related to the condition (even though in different languages) although sometimes we observe the outputs are direct translation of the prompt and not a creative story.

**Tokenisation issue**  GPT-3 generations occasionally feature "sticky" words where there are missing white spaces (e.g., *understand.With* and *timewhen*). We suspect this is due to Byte-Pair Encoding of GPT-3 where white spaces are "glued" to each subword and so every subword has two versions (one with the white space and one without). This issue arises when GPT-3 generates using a subword without the white space suffix.

**Expletives**  GPT-3 would occasionally generate stories with expletives. Interestingly, it would sometimes self-censor them (e.g., *b\*\*\*\**).

## 8   Conclusion

We present an extensive comparison of GPT-3 with SOTA models for story generation, and found that stories generated by GPT-3 are substantially better than SOTA models on multiple aspects and even rival human authors. The findings of this study indicate that we have entered a new chapter in story generation with LLMs. Future research is likely to concentrate on prompt-engineering LLMs to achieve enhanced customization, such as varying their style and length, further advancing the capabilities of story generation models. In terms of evaluation metrics, our work: 1) reveals a weak correlation between automatic lexical-based evaluation metrics and human evaluation, and that recently proposed model-based metrics appear to more reliable; and 2) contributes a new test bed for metric development, through the release of a dataset that contains story quality annotations by two groups of judges. In spite of the positive results of GPT-3 in story generation, we discuss some of its issues, the chief one being that it has a tendency to reproduce details or plots from its memories, raising foundational questions about its generation creativity.

## Limitations

As observed by Mishra et al. (2022), engineering appropriate prompts can significantly influence the performance of language models. In our current study, we randomly sample a few training examples as demonstrations for GPT-3 (in-context learning). However, a more effective approach could involve strategically selecting contextually more relevant examples.

Although text-davinci-001 was the best model at the time of our experiment, recent advancements in the field have led to the release of more powerful LLMs. Despite these improved models, we hold the view that they are unlikely to substantially alter the conclusions drawn in this study. The findings strongly suggest that LLMs will remain the dominant approach in story generation in the foreseeable future. Also, we only explore with GPT-3 in our experiments, and although we think our findings are likely to generalise to other LLMs, this has not been empirirically validated.

Since we started this work in 2022, there has been quite a development in terms of text generation evaluation metrics (Chiang and Lee, 2023; Ke et al., 2023; Xie et al., 2023; Fu et al., 2023; Liu et al., 2023a), and some of these uses LLMs themselves. Although we claim that human evaluation remains the gold standard for story generation, it remains to be seen how much these new metrics close the gap. We foresee that the question of circularity, i.e. using LLMs to evaluate LLM-generated text, will be the next challenge that the field needs to address.

In our work, we acknowledge that we did not involve domain experts (e.g., story writers) for a more specialized assessment. It would be intriguing to investigate the potential variations in judgments between lay individuals and expert evaluators in story assessment (Chiang and Lee, 2023). Recent research has indicated that certain practices in the human annotation process, such as the use of Likert scales, have limitations in capturing the true preferences of humans (Ethayarajh and Jurafsky, 2022; Liu et al., 2023c). We contend, however, that the fact that we found consistent results between two different groups of annotators suggest that our findings are likely to be robust.

## Ethics Statement

All mechanical turk experiments conducted in this paper were approved by internal ethics review board from our institution. (Ethics ID Number: 21961). Our evaluators were paid based on an estimated US$14.83 per hour rate. For each dataset, we estimate the time they would spend and vary the payment according to the estimated time. Each HIT contains 7 stories (5 stories to be evaluated and 2 controlled stories to control the evaluation quality on AMT). We pay US$2.50 per HIT for ROC, US$3.50 for WP and US$4.50 for CNN.

We remind the workers in our consent form that the potential risks about this work, which they might have to read and evaluate stories with filthy words or offended storyline and they are welcome to quit the task and we will still pay them according to the efforts they spend.

## Acknowledgements

## References

Danial Alihosseini, Ehsan Montahaei, and Mahdieh Soleymani Baghshah. 2019. Jointly measuring diversity and quality in text generation models. In *Proceedings of the Workshop on Methods for Optimizing and Evaluating Neural Language Generation*, pages 90–98, Minneapolis, Minnesota. Association for Computational Linguistics.

Prithviraj Ammanabrolu, Wesley Cheung, William Broniec, and Mark O. Riedl. 2021. Automated storytelling via causal, commonsense plot ordering. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pages 5859–5867. AAAI Press.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei.

2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

Massimo Caccia, Lucas Caccia, William Fedus, Hugo Larochelle, Joelle Pineau, and Laurent Charlin. 2020. Language gans falling short. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

Cyril Chhun, Pierre Colombo, Fabian M. Suchanek, and Chloé Clavel. 2022. Of human criteria and automatic metrics: A benchmark of the evaluation of story generation. In *Proceedings of the 29th International Conference on Computational Linguistics, COLING 2022, Gyeongju, Republic of Korea, October 12-17, 2022*, pages 5794–5836. International Committee on Computational Linguistics.

David Cheng-Han Chiang and Hung-yi Lee. 2023. Can large language models be an alternative to human evaluations? In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 15607–15631. Association for Computational Linguistics.

Elizabeth Clark, Tal August, Sofia Serrano, Nikita Haduong, Suchin Gururangan, and Noah A. Smith. 2021. All that's 'human' is not gold: Evaluating human evaluation of generated text. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 7282–7296, Online. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Kawin Ethayarajh and Dan Jurafsky. 2022. The authenticity gap in human evaluation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pages 6056–6070. Association for Computational Linguistics.

Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 889–898, Melbourne, Australia. Association for Computational Linguistics.

Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. Gptscore: Evaluate as you desire. *CoRR*, abs/2302.04166.

Seraphina Goldfarb-Tarrant, Tuhin Chakrabarty, Ralph Weischedel, and Nanyun Peng. 2020. Content planning for neural story generation with aristotelian rescoring. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4319–4338, Online. Association for Computational Linguistics.

Jian Guan, Fei Huang, Zhihao Zhao, Xiaoyan Zhu, and Minlie Huang. 2020. A knowledge-enhanced pre-training model for commonsense story generation. *Transactions of the Association for Computational Linguistics*, 8:93–108.

Jian Guan and Minlie Huang. 2020. UNION: An Unreferenced Metric for Evaluating Open-ended Story Generation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9157–9166, Online. Association for Computational Linguistics.

Jian Guan, Xiaoxi Mao, Changjie Fan, Zitao Liu, Wenbiao Ding, and Minlie Huang. 2021a. Long text generation by modeling sentence-level and discourse-level coherence. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6379–6393, Online. Association for Computational Linguistics.

Jian Guan, Zhexin Zhang, Zhuoer Feng, Zitao Liu, Wenbiao Ding, Xiaoxi Mao, Changjie Fan, and Minlie Huang. 2021b. OpenMEVA: A benchmark for evaluating open-ended story generation metrics. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6394–6407, Online. Association for Computational Linguistics.

Kazuma Hashimoto, Caiming Xiong, Yoshimasa Tsuruoka, and Richard Socher. 2017. A joint many-task model: Growing a neural network for multiple NLP tasks. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1923–1933, Copenhagen, Denmark. Association for Computational Linguistics.

Karl Moritz Hermann, Tomás Kociský, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pages 1693–1701.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. The curious case of neural text degeneration. In *8th International Conference on*

*Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

Marzena Karpinska, Nader Akoury, and Mohit Iyyer. 2021. The perils of using Mechanical Turk to evaluate open-ended text generation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1265–1285, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Pei Ke, Fei Huang, Fei Mi, Yasheng Wang, Qun Liu, Xiaoyan Zhu, and Minlie Huang. 2023. DecompEval: Evaluating generated texts as unsupervised decomposed question answering. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9676–9691, Toronto, Canada. Association for Computational Linguistics.

Jey Han Lau, Carlos Armendariz, Shalom Lappin, Matthew Purver, and Chang Shu. 2020. How furiously can colorless green ideas sleep? sentence acceptability in context. *Transactions of the Association for Computational Linguistics*, 8:296–310.

Jooyoung Lee, Thai Le, Jinghui Chen, and Dongwon Lee. 2022. Do language models plagiarize? *CoRR*, abs/2203.07618.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. A diversity-promoting objective function for neural conversation models. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 110–119, San Diego, California. Association for Computational Linguistics.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2022. What makes good in-context examples for GPT-3? In *Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures*, pages 100–114, Dublin, Ireland and Online. Association for Computational Linguistics.

Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021. Pretrain, prompt, and predict: A systematic survey of prompting methods in natural language processing. *CoRR*, abs/2107.13586.

Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023a. G-eval: NLG evaluation using GPT-4 with better human alignment. *CoRR*, abs/2303.16634.

Yiheng Liu, Tianle Han, Siyuan Ma, Jiayue Zhang, Yuanyuan Yang, Jiaming Tian, Hao He, Antong Li, Mengshen He, Zhengliang Liu, Zihao Wu, Dajiang Zhu, Xiang Li, Ning Qiang, Dinggang Shen, Tianming Liu, and Bao Ge. 2023b. Summary of chatgpt/gpt-4 research and perspective towards the future of large language models. *CoRR*, abs/2304.01852.

Yixin Liu, Alexander R. Fabbri, Pengfei Liu, Yilun Zhao, Linyong Nan, Ruilin Han, Simeng Han, Shafiq Joty, Chien-Sheng Wu, Caiming Xiong, and Dragomir Radev. 2023c. Revisiting the gold standard: Grounding summarization evaluation with robust human evaluation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 4140–4170. Association for Computational Linguistics.

Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8086–8098, Dublin, Ireland. Association for Computational Linguistics.

Li Lucy and David Bamman. 2021. Gender and representation bias in GPT-3 generated stories. In *Proceedings of the Third Workshop on Narrative Understanding*, pages 48–55, Virtual. Association for Computational Linguistics.

R. Thomas McCoy, Paul Smolensky, Tal Linzen, Jianfeng Gao, and Asli Celikyilmaz. 2021. How much do language models copy from their training data? evaluating linguistic novelty in text generation using RAVEN. *CoRR*, abs/2111.09509.

Swaroop Mishra, Daniel Khashabi, Chitta Baral, Yejin Choi, and Hannaneh Hajishirzi. 2022. Reframing instructional prompts to GPTk's language. In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 589–612, Dublin, Ireland. Association for Computational Linguistics.

Nasrin Mostafazadeh, Lucy Vanderwende, Wen-tau Yih, Pushmeet Kohli, and James Allen. 2016. Story cloze evaluator: Vector space representation evaluation by predicting what happens next. In *Proceedings of the 1st Workshop on Evaluating Vector-Space Representations for NLP*, pages 24–29, Berlin, Germany. Association for Computational Linguistics.

OpenAI. 2023. GPT-4 technical report. *CoRR*, abs/2303.08774.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Xiangyu Peng, Siyan Li, Sarah Wiegreffe, and Mark O. Riedl. 2021. Inferring the reader: Guiding automated story generation with commonsense reasoning. *CoRR*, abs/2105.01311.

Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao Chen, Michihiro Yasunaga, and Diyi Yang. 2023. Is chatgpt a general-purpose natural language processing task solver? *CoRR*, abs/2302.06476.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Hannah Rashkin, Asli Celikyilmaz, Yejin Choi, and Jianfeng Gao. 2020. PlotMachines: Outline-conditioned generation with dynamic plot state tracking. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4274–4295, Online. Association for Computational Linguistics.

Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A. Smith, and Yejin Choi. 2019. ATOMIC: an atlas of machine commonsense for if-then reasoning. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*, pages 3027–3035. AAAI Press.

Abigail See, Aneesh Pappu, Rohun Saxena, Akhila Yerukola, and Christopher D. Manning. 2019. Do massively pretrained language models make better storytellers? In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 843–861, Hong Kong, China. Association for Computational Linguistics.

Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online. Association for Computational Linguistics.

Zhihong Shao, Minlie Huang, Jiangtao Wen, Wenfei Xu, and Xiaoyan Zhu. 2019. Long and diverse text generation with planning-based hierarchical variational model. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3257–3268, Hong Kong, China. Association for Computational Linguistics.

Zhan Shi, Xinchi Chen, Xipeng Qiu, and Xuanjing Huang. 2018. Toward diverse text generation with inverse reinforcement learning. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden*, pages 4361–4367. ijcai.org.

Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. 2019. Megatron-lm: Training multi-billion parameter language models using model parallelism. *CoRR*, abs/1909.08053.

Robyn Speer and Catherine Havasi. 2012. Representing general relational knowledge in ConceptNet 5. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 3679–3686, Istanbul, Turkey. European Language Resources Association (ELRA).

Bowen Tan, Zichao Yang, Maruan Al-Shedivat, Eric Xing, and Zhiting Hu. 2021. Progressive generation of long text with pretrained language models. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4313–4324, Online. Association for Computational Linguistics.

Zhuohan Xie, Jey Han Lau, and Trevor Cohn. 2021. Exploring story generation with multi-task objectives in variational autoencoders. In *Proceedings of the The 19th Annual Workshop of the Australasian Language Technology Association*, pages 97–106, Online. Australasian Language Technology Association.

Zhuohan Xie, Miao Li, Trevor Cohn, and Jey Han Lau. 2023. Deltascore: Evaluating story generation with differentiating perturbations. *CoRR*, abs/2303.08991.

Peng Xu, Mostofa Patwary, Mohammad Shoeybi, Raul Puri, Pascale Fung, Anima Anandkumar, and Bryan Catanzaro. 2020. MEGATRON-CNTRL: Controllable story generation with external knowledge using large-scale language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2831–2845, Online. Association for Computational Linguistics.

Wenhao Yu, Chenguang Zhu, Tong Zhao, Zhichun Guo, and Meng Jiang. 2021. Sentence-permuted paragraph generation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5051–5062, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021. Bartscore: Evaluating generated text as text generation. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 27263–27277.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with BERT. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*.

Wei Zhao, Maxime Peyrard, Fei Liu, Yang Gao, Christian M. Meyer, and Steffen Eger. 2019. MoverScore: Text generation evaluating with contextualized embeddings and earth mover distance. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 563–578, Hong Kong, China. Association for Computational Linguistics.

Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. 2018. Texygen: A benchmarking platform for text generation models. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval, SIGIR 2018, Ann Arbor, MI, USA, July 08-12, 2018*, pages 1097–1100. ACM.

Yukun Zhu, Ryan Kiros, Richard S. Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*, pages 19–27. IEEE Computer Society.

## A  SOTA Story Models Details

**Knowledge Enhanced GPT-2 (KGGPT2)**  Guan et al. (2020) use heuristic rules to translate commonsense triples from commonsense knowledge bases (e.g., ConceptNet (Speer and Havasi, 2012) and ATOMIC (Sap et al., 2019)) into natural language sentences and fine-tune GPT-2 small using these sentences. They also use rules to construct negative samples from the original stories to create "bad stories" and perform additional training to encourage the model to learn representations that can distinguish the original and negative stories on ROC.

**Progressive Generation of Long Text (PRO-GEN)**  Tan et al. (2021) divide the story generation process into multiple stages where words are generated based on their order of importance (estimated using TF-IDF). In other words, PROGEN does not generate stories in a left to right manner. They fine-tune BART-large in different stages where the early stages focus on generated keywords and the intermediate stages focus on generating the next set of content words. We use PROGEN3 in our experiment which has 3 stages where it generates 15%/25%/100% of the story words after each pass.

**MEGATRON-CNTRL (MTCL)**  Xu et al. (2020) combines commonsense reasoning and storyline planning. They first train a keyword predictor with GPT-2 and the predicted keywords are used to retrieve related knowledge triples from a knowledge base. They then train a contextual knowledge ranker with BERT to rank the top-$N$ predicted knowledge triples. A second GPT-2 is trained as a conditional generator that takes both top ranked knowledge triples and other conditioning (e.g., titles) as input when generating stories. Note that the parameters of the two GPT-2 and BERT models are initialised using MEGATRON parameters (Shoeybi et al., 2019).

**High-Level Representations for Long Text Generation (HINT)**  Guan et al. (2021a) pre-train BART-base on BookCorpus (Zhu et al., 2015) with additional objectives that capture sentence-level similarity and sentence-order to learn the internal structure within a story. The model is then further fine-tuned on story datasets to generate stories in a particular dataset.

**BART**  This is a baseline model where we fine-tune BART-large on the story datasets with the standard next word prediction objective.

## B  Datasets Details

**ROCStories (ROC)**  ROC was developed by Mostafazadeh et al. (2016) and it contains 98K commonsense stories of five sentences. To obtain a more generalised lexicon, we follow the delexicalisation process from prior studies (Guan et al., 2020; Xu et al., 2020) where male/female/unknown names are replaced by [MALE]/[FEMALE]/[NEUTRAL] sentinels. For each story, the first (leading) sentence is used as conditioning context, and models are trained to generate the remaining 4 sentences.

**WritingPrompts (WP)**  WP consists of 303K human-written stories mined from Reddit's Writing Prompts forum Fan et al. (2018).[7] Each story is trimmed to contain only the first 10 sentences (following Guan et al. (2021a)). For WP, we use the prompt (which is typically a paragraph of text that sets the scene of the story) as conditioning for story generation.

**CNN News (CNN)**  CNN News (Hermann et al., 2015) is a dataset that contains long news articles with titles. CNN is a very large dataset, with 311K news articles and highlights. We sub-sample the standard training, validation and testing splits to produce splits with 10K/5K/1K stories each, respectively, for our experiments. The title of a news story is used as conditioning for story generation.

## C  Amazon Mechanic Turk Setting

**Qualification Requirements**  We set following qualification requirements for our annotators: 1) Their accept rate is greater than or equal to 97%. 2) Their location is in US. 3) They have to complete more than 1000 HITs.

**Questions**  We ask the following questions in our questionnaire.

1. Fluency: "How grammatically correct is the text of the story?"

2. Coherence: "How well do the sentences in the story fit together?"

3. Relatedness: "How relevant is the story to the title?"

---

[7]https://www.reddit.com/r/WritingPrompts/

4. Logicality: "How much does the story obey commonsense?"

5. Interestingness: "How enjoyable do you find the story?"

## D Amazon Mechanic Turk Pilot Study

While AMT is convenient to find workers for annotation work, it can be rather difficult to obtain reliable workers (Karpinska et al., 2021; Clark et al., 2021). One of our workers told us many workers install website plugins to help them to manage the workflow with AMT so that they can hoard many HITs at the same time. Therefore, HITs with high payment can easily attract irresponsible workers even though previous qualifications are set since most AMT requesters will not bother to reject work.

Therefore, we set a pilot study to aid us to help reliable workers. We randomly select 5 stories generated from different models on ROC and 1 story from the test dataset. We then train a trigram language model on ROC to mimic the style and generate 1 story from the trigram model. All stories have different titles. We randomly shuffle these 7 stories and the task is to ask people to evaluate all stories with questions mentioned in Appendix C and we will judge the quality of their evaluation based on human and trigram stories.

We invite 7 of our colleagues, which are all from non-English speaking countries to have a rough idea of the difficulty degree of the task. We calculate the average score of all quality metrics except the interestingness aspect since it is subjective. On average, our colleagues rank the human story as 4.5 and trigram story as 1.425, which shows our task is not hard to distinguish human and trigram stories. We set a rather lenient standard as "ranking human story >= 3.5 and trigram story <= 2.0" to select workers from our pilot study.

We create 100 assignments of the same HIT at different times with the qualification mentioned in Appendix C. We find running the same pilot study at different times can obtain quite different results from AMT, which align to the findings in Karpinska et al. (2021). Generally, we find that more reliable workers can be found in the evening of Eastern Daylight Time (EDT). We have 10 out of 100 people pass the pilot study but only 5 people pass it on the day. It shows the difficulty of obtaining reliable workers on AMT nowadays and the economic importance of running a pilot study

|  | IAA | Flu. | Coh. | Rel. | Log. | Int. |
|---|---|---|---|---|---|---|
| **ROC** | $r$ | 0.64 | 0.81 | 0.79 | 0.80 | 0.68 |
| | TA | 17.24 | 24.98 | 25.57 | 27.37 | 22.03 |
| **WP** | $r$ | 0.51 | 0.70 | 0.74 | 0.71 | 0.54 |
| | TA | 18.37 | 17.01 | 32.65 | 19.73 | 12.93 |
| **CNN** | $r$ | 0.46 | 0.54 | 0.61 | 0.59 | 0.50 |
| | TA | 15.13 | 12.61 | 15.97 | 11.76 | 14.29 |

Table 7: Inter Annotator Agreement (IAA) results for each aspect: fluency (Flu.), coherence (Coh.), relatedness (Rel.), logicality (Log.) and interestingness (Int.). We use one-vs-rest Pearson's $r$ to assess the extent to which each annotator agrees with the consensus. Total Agreement (TA) means the percentage where all 3 annotators choose the same score.

before conducting real research. We grant those reliable workers the customised qualification and only invite them to our real study, we also have controlled stories to monitor the quality of workers, as 2 controlled stories inserted into each HIT.

## E Amazon Mechanic Turk Issue

Our human evaluation is conducted over AMT, even though it is convenient and affordable, we find a big disagreement between our annotators. We first conduct a pilot study to test the capability of annotators to evaluate English stories and only invite workers that pass our proficient English stories reading tests to the evaluation of sampled stories. We only gave them two examples showing how we assess the example stories but we did not provide detailed English stories evaluation training to our annotators. We did not have a main annotators that can provide a standard score for example stories, which increase the difficulty of judging the quality of evaluation work we receive from AMT.

Also, as pointed out in Karpinska et al. (2021), the quality of work from annotators on AMT platform can be of high variance and have poor calibration, therefore, we would obtain more reliable human evaluation results if we hire expert raters such as professional authors or English language teachers.

## F Inter-annotator Agreement for MTurk Workers

We follow Lau et al. (2020) to estimate one-vs-rest agreement using Pearson's $r$. For each story, we single out an annotator's score and compare it to the mean scores given by the other two annotators, and we repeat this process for every score in a story

|     | IAA | Flu. | Coh. | Rel. | Log. | Int. |
|-----|-----|------|------|------|------|------|
| **ROC** | $r$ | 0.42 | 0.54 | 0.66 | 0.59 | 0.32 |
|     | TA | 38.57 | 25.0 | 25.71 | 25.71 | 8.57 |
| **WP** | $r$ | 0.36 | 0.57 | 0.73 | 0.49 | 0.54 |
|     | TA | 10.0 | 10.0 | 18.57 | 10.0 | 10.0 |
| **CNN** | $r$ | 0.36 | 0.41 | 0.47 | 0.37 | 0.35 |
|     | TA | 17.14 | 10.71 | 14.29 | 10.0 | 4.29 |

Table 8: In-house IAA Results.

and for all stories to compute Pearson's $r$ over the two sets of scores (singled-out scores vs. mean scores). We also compute the percentage where all 3 annotators choose the same score, noting that this is a much stricter agreement metric (as it does not capture the ordinal scale of the scores). Random scoring would produce 4% for this metric.

IAA results are presented in Table 7. In terms of one-vs-rest agreement ($r$), we find overall good agreement with 9 strong agreement results ($r >= 0.6$) and 6 moderate agreement results ($0.45 <= r <= 0.6$). We see some correlation between story length and agreement, as ROC has the highest agreement (shortest with 5 sentences) and CNN has the lowest (over 20 sentences). When it comes to aspects, coherence, relatedness and logicality have higher agreement compared to fluency and interestingness. While it is intuitive to see interestingness being subjective, fluency is somewhat a surprise. Manual inspection reveals that annotators have very different standards when it comes to fluency, with some workers being more strict about grammar, which contributes to the low agreement. For total agreement (TA), the numbers range between 10–25%, which is encouraging as it shows that there is still a good proportion of cases where all annotators agree on a score.

## G  Inter-annotator Agreement for In-house Workers

The In-house annotation agreement results are shown in Table 8. Surprisingly, we see similar trends in compared to Table 7, which means our findings are consistent even with a different group of workers on a different set of stories.

We again see fluency and interestingness are two aspects that have lower agreements. However, it is interesting that we find fluency has highest TA but rather lower $r$. Manual inspection reveals that annotators mainly select 4 or 5 for fluency (also can be seen from high fluency scores for ROC in Ta-

ble 5), which benefits TA because they have higher chance to agree on the same values, but might harm $r$ because the score can be negatively influenced by other rare values. Relevance again has the highest agreement, which might be because annotators agree that many stories generated by SOTA models are not related to the given condition. However, compared to IAA results from crowdsource annotators, in-house annotators have lower agreement in general. we find 2 strong agreement results ($r >= 0.6$) and 6 moderate agreement results ($0.45 <= r <= 0.6$).

## H  GPT-3 World knowledge Strengths

Table 9 and 10 demonstrates that GPT-3 can generate better stories than SOTA models thanks to its wide world knowledge strengths.

## I  Mechanic Turk Instructions

Figure 1 and 2 demonstrate our instructions and rubrics.

## J  Example Prompts

Table 11 shows some examples we use in our prompts for generation.

## K  GPT-3 Stories vs Human References

Table 12, 13 and 14 show some examples where GPT-3 obtains higher scores than its human reference.

## L  Example Generations

Table 15 shows randomly picked generation example for each dataset.

## M  Soft Plagiarism Examples

Table 16 shows some GPT-3 generated news stories where we find true online stories.

## N  GPT-3 Generation Issues

Table 17 shows some GPT-3 generated stories with several issues.

## O  Pearson Correlations between Each Story Aspects

Figure 3 and 4 present Pearson Correlations between Each Story Aspects for MTurk and in-house workers respectively.

| Model | Story |
|---|---|
| GPT-3 | as soon as i realized that i had died , the first thing that came to mind was how pointless my life had been . i had devoted myself to a cause that was ultimately fruitless and destructive , and in the end i had caused the deaths of millions of people .but then i realized that there was still a chance for redemption . in the after life , i was given the opportunity to travel to different moments in history and try to make things right . i could see the consequences of my actions , and i could try to make different choices that would lead to a better outcome .i started with the holocaust , and i was able to save many people from being killed . but it was difficult , and i had to make a lot of tough choices . eventually i moved on to other moments in history , trying to make things better wherever i could .it was a difficult journey , but it was worth it . i was finally able to make up for the mistakes that i had made in life . |
| PROGEN3 | i knew i was in a bad mood . the blood dripped from my nose and down my chin like the burly balding drifter . was it dark yet ? i could feel my brain becoming small dots of the rotten material . the blood stopped flowing . almost instantly , the world dissolved in my hands . that 's when i noticed something . i was still standing in the middle of the street . i turned to my left and saw the madman . |
| HINT | " hi , " i said , opening my eyes to a bright light . " this is hitler . " i looked at the clock . 5:00 am . " what 's your name ? " " " hitler . " " hitler . " " hitler . " " hitler . " " hitler . " |
| BART | i 'm struck with a fierce sense of shame at that thought . unfortunately it 's been a while since i 've felt this way , and normally it 's * bad * to feel that way , but i think that too should be expected when you 're still alive . there has been so much turmoil since the war . we 've lost children , we 've lost society , we 've even lost ourselves . but now , we 're free . we can finally go on with our lives , like we always should . but something happens that seems wrong , and it 's really hard to forget . there are my reasons for feeling so sad . why will i be mad today ? |

Table 9: Stories generated from GPT-3, HINT, BART and PROGEN3 given the same title "you 're hitler , and you 've just died . there is an after life . tell me about it ."

| Model | Story |
|---|---|
| GPT-3 | the young man had lost everything . his family , his home , his country . he had nothing left to lose . when he was offered the chance to join isis , he saw it as a way to finally have a purpose in life . he was given a sense of belonging and a sense of worth . he was told that he was fighting for a good cause and that he was helping to protect his people . he was given a sense of identity and a sense of purpose . he was told that he was part of something bigger than himself . |
| PROGEN3 | they 've crippled the ranks of these countries . their leaders not going to the president 's party , everyone having their noses to the grindstone , the war . political parties shooting each other with a firing squad , only to end up murdering each other . " the one who executed your followers was the first one to hell houng taqibs . " he 's right , ha is not 'd understand * reason * if he knew that . " you 're killing my people . " |
| HINT | my eyes stare at the ceiling . it 's just a window , it 's just a window . i 'm not sure if i should close it or close it again . i 'm afraid i ca n't move it . i do n't want to be here . i do n't want to be here . i want to be able to talk to the people who will help me . to be able to hear the voices of people who will help me . to be able to see the lights . to be able to see the people who will help me . |
| BART | missdamie looked at her young daughter sat with c and o . she 'd put on a radiant glow and bound herself up in an artificial bubble . she was beautiful . she twitched at the vibrancy in her hands . her hair was beautiful and silver . her smile was off . her black nails placed on gentle curves . when she smiled her smile widened and the small eyes of her daughter 's eyes danced . " isis ? " i asked , shakily , behind the mask . ms. de los angeles nodded and smiled . " yes , my name is isis de los angeles , and " her mother held her hand . " my son is wondering if isis can cause such a thing . " |

Table 10: Stories generated from GPT-3, HINT, BART and PROGEN3 given the same title "make me sympathise with an isis member ."

# Instructions:

Please read each story and judge its quality based on the five criteria given below. You have to finish assessing all stories before you press the "Submit" button.

**Note:**

- Some names in the stories are replaced with "[MALE]", "[FEMALE]" and "[NEUTRAL]", and all words are lowercased. Please do not treat these changes as grammar errors.
- These five criteria we list below should be judged <span style="color:red">independently</span>.
- **Fluency** only measures if each sentence of story has grammar errors.
- **Logicality** only measures if the story make senses based on your commonsense knowledge, For instance, "sprayed the burn with a napkin" in Example Story One does not make sense since we can't use a napkin to do the spraying.
- **Relatedness** only measures if a story is related to the title.
- **Coherence** only measures if the narrative of the story flows well - you do not need to consider how well the story fit the title in this criterion.
- **Interestingness** is a subjective metric - you only need to select based on the extent to which the story interests you and you do not have to consider if it makes sense.

## Please Judge Story One

**Title:** ${title1}

**Story:** ${story1}

**Fluency:** How grammatically correct is the text of the story?

○1(lowest)  ○2  ○3  ○4  ○5(highest)

**Coherence:** How well do the sentences in the story fit together?

○1(lowest)  ○2  ○3  ○4  ○5(highest)

**Relatedness:** How relevant is the story to the title?

○1(lowest)  ○2  ○3  ○4  ○5(highest)

**Logicality:** How much does the story make sense?

○1(lowest)  ○2  ○3  ○4  ○5(highest)

**Interestingness:** How enjoyable do you find the story?

○1(lowest)  ○2  ○3  ○4  ○5(highest)

Figure 1: A screenshot of our evaluation questions.

**Special Note: Coherence and Logicality measure two different aspects.**
**For instance:**
My 6-year-old son loves reading and writing. He loves The Little Prince I bought him. - is coherent and logical.
My 6-year-old son loves reading and writing. He hates The Little Prince I bought him. - is incoherent but logical.
My 6-year-old son loves reading and writing. He loves writing with an apple. - is coherent but illogical.
My 6-year-old son loves reading and writing. He burns The Little Prince I bought him with his water gun. - is incoherent and illogical.

**Rating Rubric:**
Please have a read of the rating rubric and rate the stories accordingly.

| Aspect | Rubric |
| --- | --- |
| Fluency: "How grammatically correct is the text of the story?" | 1: The story is full of grammar issues so that you totally cannot understand the story. For instance, incomplete or repeated or missing words or phrases.<br>2: The story has too many grammar issues that largely impacts your understanding of the story.<br>3: The story has a few grammar errors, but doesn't impact your understanding of the story too much.<br>4: The story has some slight grammar errors, but these errors have no impact of your understanding of the story.<br>5: The story has no grammar issue at all. |
| Coherence: "How well do the sentences in the story fit together?" | 1: The sentences of the story are totally irrelevant or contradicted to each other. For instance, the second sentence is irrelevant to the first sentence without any conjunction.<br>2: Most sentences of the story are irrelevant or contradicted to each other that largely impact your understanding of the story.<br>3: The story has a lot irrelevant or contradicted parts, but you can still understand the story.<br>4: Most sentences fit into the story, except for one or two slight irrelevant or contradicted parts.<br>5: All sentences of the story fit together, nothing irrelevant or contradicted. |
| Relatedness: "How relevant is the story to the title?" | 1: The story has no relationship with the title at all.<br>2: The story only has a weak relationship with the title.<br>3: The story roughly matches the title.<br>4: The story matches the title, except for one or two small aspects.<br>5: The story matches the title exactly. |
| Logicality: "How much does the story obey your commonsense?" | 1: The story is full of absurd things. For instance, fish fying in the sky, birds swimming in the sea (unless you strongly believe they are true according to the story context.)<br>2: The story has one or two things make sense, but generally very absurd.<br>3: The story roughly makes sense.<br>4: The story largely makes sense, except one or two things reads weird that you have to think twice.<br>5: The story totally complies with commonsense. |
| Interestingness: "How enjoyable do you find the story?" | 1: You find the story so boring and you are glad it is over.<br>2: You only find one or two things interesting in the story, but no more.<br>3: The story is mildly interesting.<br>4: The story is interesting that keeps you engaged until the end.<br>5: You feel the story is so interesting that you wish there is a sequel. |

Figure 2: Explanation and Rubric in the Instruction.

|  | **Example** |
|---|---|
| ROC | **prompt**: [MALE] wants to move to a new house . **completion**: his wife [FEMALE] does not want to move . [MALE] convinces her to look at a couple new houses . she does not see any that she likes . [FEMALE] refuses to move . |
| WP | **prompt**: the real reason we 're not supposed to look at the sun during a total eclipse is because it reveals something horrifying . **completion**: the harvest was here . father xander made his way through priory halls toward the sanctuary in the back garden . at his advanced years , the trip was a struggle ; each step brought a hefty dose of pain and unpleasant memories . it was eclipse day , and while he always prayed for a good outcome , in his eighty years as the priory ' s eyes , there never was . " the initiates are assembled in the sanctuary , father , " spoke a fellow priest in the hall . " thank you , my son , " father xander replied . he paused and looked at the young man . memories came flooding back like sour wine . his name was thomas hampstead , and at thirty years of age , he should have been a robust healthy lad ; but he was not . thomas ' s skin had a greyish tint and his body was frail . |
| CNN | **prompt**: (CNN) – Usain Bolt rounded off the world championships Sunday by claiming his third gold in Moscow as he anchored Jamaica to victory in the men's 4x100m relay. **completion**: The fastest man in the world charged clear of United States rival Justin Gatlin as the Jamaican quartet of Nesta Carter, Kemar Bailey-Cole, Nickel Ashmeade and Bolt won in 37.36 seconds. The U.S finished second in 37.56 seconds with Canada taking the bronze after Britain were disqualified for a faulty handover. The 26-year-old Bolt has now collected eight gold medals at world championships, equaling the record held by American trio Carl Lewis, Michael Johnson and Allyson Felix, not to mention the small matter of six Olympic titles. The relay triumph followed individual successes in the 100 and 200 meters in the Russian capital. Ï'm proud of myself and I'll continue to work to dominate for as long as possible,B̈olt said, having previously expressed his intention to carry on until the 2016 Rio Olympics. Victory was never seriously in doubt once he got the baton safely in hand from Ashmeade, while Gatlin and the United States third leg runner Rakieem Salaam had problems. Gatlin strayed out of his lane as he struggled to get full control of their baton and was never able to get on terms with Bolt. Earlier, Jamaica's women underlined their dominance in the sprint events by winning the 4x100m relay gold, anchored by Shelly-Ann Fraser-Pryce, who like Bolt was completing a triple. Their quartet recorded a championship record of 41.29 seconds, well clear of France, who crossed the line in second place in 42.73 seconds. Defending champions, the United States, were initially back in the bronze medal position after losing time on the second handover between Alexandria Anderson and English Gardner, but promoted to silver when France were subsequently disqualified for an illegal handover. The British quartet, who were initially fourth, were promoted to the bronze which eluded their men's team. Fraser-Pryce, like Bolt aged 26, became the first woman to achieve three golds in the 100-200 and the relay. In other final action on the last day of the championships, France's Teddy Tamgho became the third man to leap over 18m in the triple jump, exceeding the mark by four centimeters to take gold. Germany's Christina Obergfoll finally took gold at global level in the women's javelin after five previous silvers, while Kenya's Asbel Kiprop easily won a tactical men's 1500m final. |

Table 11: Prompts and GPT-3 Completions, including one example each for ROC, WP and CNN.

| Title [SEP] Story | From | Flu. | Coh. | Rel. | Log. | Int. |
|---|---|---|---|---|---|---|
| tell me the story of a totally inconsequential character who does n't realize he 's wrapped up in a massive high-stakes plotline - and keep them in the dark about it . [SEP] ivan was just an ordinary guy who led an ordinary life . he worked a dead-end job at a local factory , and spent his evenings and weekends playing his favorite online game with his friends . one day , out of the blue , ivan was contacted by a representative of the game company , who offered him a chance to participate in a closed beta test for their new game . ivan was excited to participate and eagerly agreed to the terms and conditions .little did he know that he was about to become embroiled in a global conspiracy . the game company was actually a front for a secret organization that was planning to use the game to recruit new members . ivan was the perfect candidate , because he was someone who was looking for an escape from his ordinary life .the organization was planning to use the game to recruit new members from all over the world , and they were counting on ivan to help them . ivan had no idea what he was getting himself into , and the organization was careful to keep him in the dark about their true intentions . ivan thought he was just playing a game , but he was actually a pawn in a much larger scheme . | GPT-3 | 5.0 | 5.0 | 5.0 | 5.0 | 5.0 |
| tell me the story of a totally inconsequential character who does n't realize he 's wrapped up in a massive high-stakes plotline - and keep them in the dark about it . [SEP] someone was moving around on screen three . * probably just another junkie stumbling around . fucking dumbass . * barry did n't really give a shit what the junkie was up to anyway . his job did n't pay enough for that . it 's not like they thought he was actually all that important , anyway . a security guard ? minimum wage , take your paycheck , shut up , and move along , thank you very much . you have to sit for twelve hours in the dead of night and devote all of your attention to watching the same twelve tiny boxes ? that 's nice . how 's $ 7.25 an hour sound ? you 're preventing vandalism and break-ins ? how wonderful . would you like a cookie ? | human | 3.67 | 2.33 | 2.0 | 3.33 | 2.33 |

Table 12: GPT-3 Generated Story Obtains Higher Scores than Its Human Reference.

| Title [SEP] Story | From | Flu. | Coh. | Rel. | Log. | Int. |
|---|---|---|---|---|---|---|
| you meet the grim reaper at a crossroads . you are fearful as he approaches , but notice he seems just as fearful of you . [SEP] the grim reaper was a fearful figure , draped in a cloak of blackness with a face that was hidden in shadow . as he came closer , you could see the fear in his eyes . " why are you so afraid of me ? " you asked . " because you 're the first person to ever see me like this , " he replied . " i 'm not here to take you , i 'm here to ask for your help . " | GPT-3 | 4.0 | 4.67 | 4.67 | 3.67 | 3.67 |
| you meet the grim reaper at a crossroads . you are fearful as he approaches , but notice he seems just as fearful of you . [SEP] my eyes jump from snowbank to snowbank , one featureless as the next . the road i walk upon is a dark grey slush , the result of many cars and snow plows making their way through the void . yet i had not seen a car in miles . i shiver . but not from the cold . the compelling . it had awakened me . this ominous feeling that i had to be at a certain place at a certain time . it controlled me . it made me get dressed . it made me open the door . it had even made me ignore my wife , who i usually kiss goodbye . the road seemed endless , but i felt no reason to stop . | human | 3.33 | 3.33 | 1.67 | 2.33 | 2.33 |

Table 13: GPT-3 Generated Story Obtains Higher Scores than Its Human Reference.

| Title [SEP] Story | From | Flu. | Coh. | Rel. | Log. | Int. |
|---|---|---|---|---|---|---|
| you are able to see the number of days a person has left to live above their heads . for the first time in your life , you see someone with a negative number . [SEP] as you walk the streets , dodging the rushing crowds , you notice something different about one of the people walking by . their number is displayed in a negative value . for a moment , you feel a pang of sadness and empathy . you know what that means - that person is going to die soon . you watch as they go about their day , completely unaware of the timer ticking away above their head . you wonder who they are , what they 're doing , and who will be there with them when their time runs out . | GPT-3 | 4.67 | 4.33 | 4.67 | 4.0 | 5.0 |
| you are able to see the number of days a person has left to live above their heads . for the first time in your life , you see someone with a negative number . [SEP] crowds are the worst . frank squints through the flurry of hovering digits as he shoulders his way through the masses of shoppers . the cacophony of vociferous patrons , distressed employees , and typical black friday commotion crescendos around him . he pauses to watch two shoppers vie over the last remaining television monitor ( only $ 600 ! ) . one is a blonde woman whose pigtails have been jostled into disarray ; above her hovers the number 10951 . the other is n't as fortunate ; the number 582 glows above her greying hair . frank sighs . he 's sure there 's something . . . vaguely poignant about fighting over a product priced at a greater value than your remaining days of life . socio-rhetorical commentary on modern consumerism , perhaps . he has little time to contemplate , however , as another shopper impatiently shoves him forward , and frank resumes his pace . while he walks , frank studies the various faces passing by : a woman in an elegant , blue hijab ; a portly , bearded man ; a fatigued mother clutching her two children by the hands ... . | human | 3.67 | 3.0 | 3.33 | 4.0 | 3.33 |

Table 14: GPT-3 Generated Story Obtain Higher Scores than Its Human Reference.

346

| Title | Generation |
|---|---|
| the boy cried . | the boy cried because he was sad and scared . his parents were getting a divorce and he did n't know what was going to happen to him . he was going to have to leave his home and his friends and move to a new house . his mother said that he could choose where he wanted to live , but his father said that he had to live with him . the boy did n't know what to do . |
| your extreme apathy has been mistaken as spiritual enlightenment , and now people are building a religion around you . | my name is john and i have a problem . well , actually , a lot of people have a problem with me but that 's not the point . the point is that i don 't care about anything . at all . not the way most people care about things anyway . when most people see a homeless person on the street they feel compassion and want to help . when i see a homeless person on the street i think " well , there goes another lazy bum who 's too good for work " . and i 'm not joking when i say that most people see the world in a completely different way than i do . they see the good in people and want to help them . i see the bad in people and want to help them . see the difference ? i 'm not saying that my way of seeing the world is wrong and theirs is right . i 'm just saying that we see the world differently and that 's okay . or at least it should be . but it 's not because people can 't accept that i don 't care about anything . they think that i 'm enlightened because i don 't feel the pain and suffering that they do . they think that because i don 't care about anything i must be in a better place than they are .and that 's why they 're building a religion around me . |
| (CNN) – Hurricane Ida moved into the southern Gulf of Mexico Sunday, prompting a declaration of emergency in Louisiana and concern along the U.S. Gulf Coast. | ida Louisiana Gov. Bobby Jindal declared a state of emergency Saturday as the Category 1 hurricane churned in the Gulf of Mexico, according to a news release from the governor's office. "Based on the latest forecasts, Hurricane Ida is expected to make landfall in Louisiana as a Category 1 hurricane and we must be prepared for significant impacts," Jindal said in the release. "I have directed my team to do everything possible to protect the lives and property of the people of Louisiana." The hurricane could bring storm surges, high winds and heavy rains to the coast, the National Hurricane Center said. The storm was about 220 miles (355 kilometers) south-southeast of the mouth of the Mississippi River on Sunday morning and was moving north-northwest at about 9 mph (15 kph), the center said. The hurricane had maximum sustained winds of 75 mph (120 kph). Forecasters said Ida could strengthen a bit before making landfall. The center of the storm is expected to come ashore early Tuesday morning. |

Table 15: Randomly Chosen GPT-3 Generation Examples from Each Dataset.

**GPT-3 Generated Story (Link to the Online Resource)**

Washington (CNN) – With the posting of 400,000 classified documents from the Iraq war, WikiLeaks has shown a much heavier hand redacting compared to its previous publication of documents. In an interview with Wolf Blitzer on Wednesday, WikiLeaks founder Julian Assange said that the organization redacted the names of Iraqis who provided information to the United States military in order to protect them from retaliation."We have a duty to protect some of the people who are named in these documents," Assange said.WikiLeaks has redacted the names of coalition forces and informants, as well as the names of Iraqis who cooperated with the United States military, in all of its recent releases of classified documents. (https://www.reuters.com/article/britain-assange-idUSL4N2GD2PZ)

(CNN) – A former University of Virginia lacrosse player convicted of killing his ex-girlfriend was sentenced to 23 years in prison on Thursday, according to the Commonwealth Attorney's office in Charlottesville, Virginia. George Huguely V was convicted in February of second-degree murder in the May 2010 death of Yeardley Love. He was also convicted of grand larceny and burglary.Prosecutors had asked for a sentence of 30 years, while Huguely's defense attorneys had asked for 10 years."Yeardley was my everything," Huguely said in a statement before sentencing. "I was not thinking, I was not rationale, and I was not in control of my own actions."Love was found beaten to death in her off-campus apartment. Huguely's DNA was found on Love's body and in her bedroom.During the trial, prosecutors said Huguely had sent Love a series of emails in which he threatened to kill her."I'm going to kill you. You deserve it. You took the one thing I loved the most," Huguely wrote in one email, according to prosecutors. (https://en.wikipedia.org/wiki/Murder_of_Yeardley_Love)

(CNN) – Crisis club Racing Santander has been thrown out of next season's Copa del Rey after its players staged a bizarre on-pitch boycott of a quarterfinal tie against Real Sociedad. Racing Santander's players staged a bizarre on-pitch boycott of a Copa del Rey quarterfinal tie against Real Sociedad on Wednesday, leading to their expulsion from the competition.The players refused to take to the field in the first leg of the last-eight tie at the Anoeta Stadium in San Sebastian, protesting at the club's failure to pay their wages.After a delay of around 45 minutes, the referee finally called off the game, with the Spanish Football Federation (RFEF) confirming Racing Santander's expulsion from the competition."The Competition Committee has decided to expel Racing Santander from the Copa del Rey for not appearing in the match," said an RFEF statement."Real Sociedad have been awarded a 3-0 win and will now face the winners of the match between Real Madrid and Valencia, who will play the semi-final."Racing Santander, who are currently bottom of the Spanish first division, released a statement on their website protesting at the club's treatment by the RFEF."Racing Santander has today been expelled from the Copa del Rey for not appearing in the match against Real Sociedad, despite the fact that our players were prepared to play," said the statement."We believe that this decision is unjust, given that our players have not been paid their wages for several months and have made it clear that they would not be able to compete in the match."The club also released a video of the players refusing to take to the field, with captain David Lafuente leading the protest.Racing Santander's expulsion from the Copa del Rey means they will not compete in next season's competition, while Real Sociedad will face the winners of the match between Real Madrid and Valencia in the semi-final. (http://edition.cnn.com/2014/01/30/sport/football/spanish-cup-boycott-santander-football/index.html)

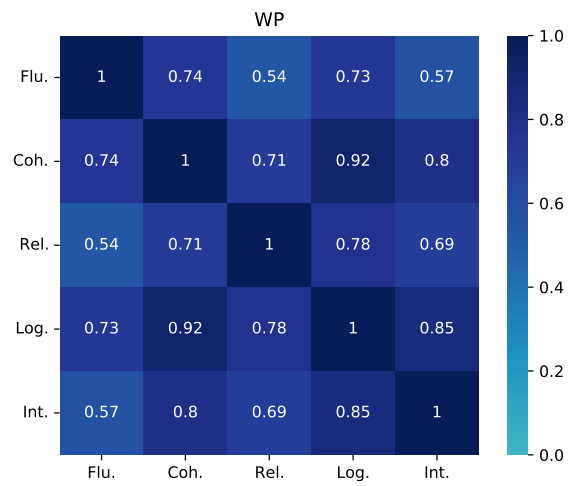Table 16: GPT-3 generated CNN News and link to the Online Resource

| Issue | Title | Generation |
|---|---|---|
| Generate Chinese from an English title | you are an alien body snatcher , and this is your first day in your new host body . you have no idea what his life is . | 你不是人，你是外星人。这是你第一天在新的主人身上。你不知道他的生活是什么。你只是偶然得到他的身体，他的生活已经开始了，而你还不知道这是什么。你的主人是一个普通人，他没有你的外星人能力。你不知道你的主人在做什么，也不知道他有没有你的朋友。你只是在他的身体里，完全不知道他的生活。 |
| Generate Japanese from an English title | what it means to live ( 250 words ) . | 人生とは何かを探して生きることである。探しけることで、生きる をつける。例えばを探しけることで、その が叶う率を高める。を叶えるために努力をしけることで、そのが叶う可能性が高まる。しかし、を叶えるために努力をしけることはしい。を叶えるためには、そのにして情を持つことが大切だ。そして、そのを叶えるためには、自分を信じることが大切だ。を叶えるためには、そのにチャレンジすることが大切だ。そして、を叶えるためには、失を恐れずに前することが大切だ。 |
| Stick word together to create "new" words | write a short story based off the lyrics to any song . | there was a timewhen we were youngand our dreams were alivewe wanted to be rock starsand travel the worldwe wanted to be someonesomeone famoussomeone lovednow we 're just oldour dreams have diedour passions gonewe just want to be left aloneto die in peace |
| Copy from the leading text | [MALE] was a very sweet dog . | [MALE] was a very sweet dog . he loved his family and was always happy to see them . even though he was old , he never stopped wagging his tail and licking everyone 's faces . |
| Copy from the prompt | [MALE] is a manager at kfc . | he started as a fry cook [MALE] wants to move to a new house , but his wife [FEMALE] does not want to move . [MALE] convinces her to look at a couple new houses , and she does not see any that she likes . [FEMALE] refuses to move . (Table 11) |
| Generate profanity words | [FEMALE] was eating lunch at school . | a boy came up to her and asked her to go out with him . she said no and he called her a stuck up b * * * * . |

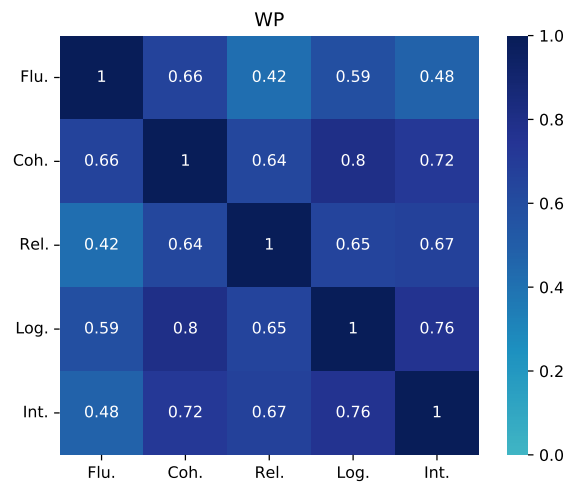Table 17: Examples of some interesting generation errors we find for GPT-3.
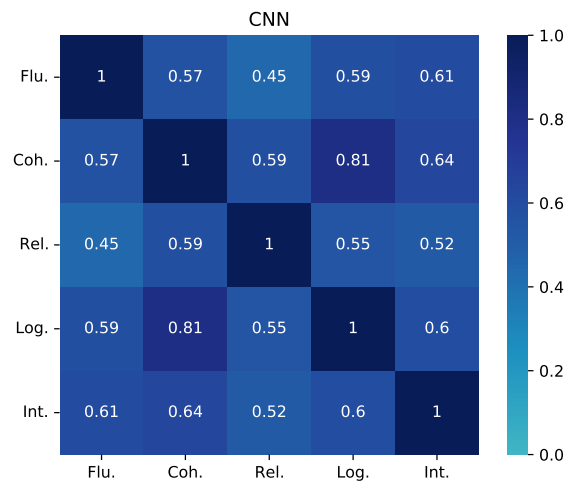
Figure 3: Pearson Correlations between Each Aspect from Crowdsourcing annotations.

(a) ROC



(b) WP



(c) CNN

Figure 4: Pearson Correlations between Each Aspect from in-house annotations.

# Trustworthiness of Children Stories Generated by Large Language Models

**Prabin Bhandari**
Department of Computer Science
George Mason University
`pbhanda2@gmu.edu`

**Hannah Marie Brennan**
Department of English, Linguistics Program
George Mason University
`hbrennan@gmu.edu`

## Abstract

Large Language Models (LLMs) have shown a tremendous capacity for generating literary text. However, their effectiveness in generating children's stories has yet to be thoroughly examined. In this study, we evaluate the trustworthiness of children's stories generated by LLMs using various measures, and we compare and contrast our results with both old and new children's stories to better assess their significance. Our findings suggest that LLMs still struggle to generate children's stories at the level of quality and nuance found in actual stories.[1]

## 1 Introduction

Advancements in pretrained large language models (LLMs) like GPT-3 (Brown et al., 2020) and LLaMA (Touvron et al., 2023), have made it easier to generate natural language text for a variety of downstream tasks, including generating narrative text like children's stories. The ability to generate natural text using LLMs has seen substantial improvement with the innovation of instruction-following models like InstructGPT (Ouyang et al., 2022) and Alpaca (Taori et al., 2023), resulting in a better alignment with user intentions.

These systems are being used as a general-purpose chat-bots by the general public. As these models are integrated more into everyday applications, it is crucial to continuously evaluate LLMs' performance to ensure that they are indeed trustworthy and accurate.

Trustworthiness in the case of LLMs is a broad term that refers to reliability and confidence in the generated text outputs along with their suitability for a specific downstream task. A trustworthy LLM minimizes errors, biases, and potentially harmful content while consistently producing clear and contextually suitable text. With the advancing capabilities of LLMs, concerns regarding their trustworthiness have arisen. Notably, they are being used more frequently to support creative writing (Clark et al., 2018), raising concerns about the generation of inappropriate or offensive text (Price, 2016) and biased content (Lucy and Bamman, 2021a). One domain in which trustworthiness is of particular importance is text generation intended for children. This paper seeks to evaluate the trustworthiness of children's stories generated by LLMs including generative LLMs and instruction following models. In the case of text generation geared towards children, LLMs' ability to generate age-appropriate materials to target audiences also becomes a vital aspect of overall trustworthiness.

To assess the trustworthiness of LLMs in generating children's stories, we use two open-source foundation language models, OPT (Zhang et al., 2022) and LLaMA (Touvron et al., 2023), along with an instruction-following model Alpaca (Taori et al., 2023) to generate children's stories. Then, we compare these generated stories against actual children's stories, old and modern. Our assessment takes into account a number of aspects, including statistics derived from the text like the Flesch reading ease score (Flesch, 1948), toxicity present in the text, the most influential topics present in the text, and the sentence structure of these texts.

Our findings reveal that LLMs lack a high level of trustworthiness when tasked with generating children's stories. While the generated children's stories do share similarities in topics and patterns with the actual stories (mostly modern ones), they are also susceptible to generating toxic content. Moreover, LLMs struggle to capture the intricacies and nuances of children's literature, evident from the disparity in sentence structure between the generated and actual stories.

---

[1]Code and dataset are publicly available: https://github.com/prabin525/trustworthiness-of-children-stories-generated-by-LLMs

## 2 Related Work

### 2.1 Story Generation

Recently, LLMs have been increasingly used to supplement creative writing efforts for entertainment and social media. Applications include work related to narrative generations (Sun et al., 2023; Simon and Muise, 2022; Razumovskaia et al., 2022; Xiang et al., 2018). Yuan et al. (2022) tested Wordcraft, a tool created to assist writers with story generation using LLMs. In their study, writers who were tasked with working with the AI agent noted that Wordcraft lacked content awareness and would create grammatical stories with nonsensical topics or plots.

### 2.2 Children and AI

AI and LLMs have also been applied to contexts involving children. Researchers at MIT had children work with social robots to evaluate how much the children could learn through activities involving robots (Williams, 2019). There is much discussion on how to integrate AI into early childhood education (Yang, 2022; Kasneci et al., 2023). With the increasing use of AI by and around children, there is an urgent need for more thorough evaluations of LLMs and the appropriateness of generated content for vulnerable audiences.

### 2.3 Trustworthiness Testing

Chiang and Lee (2023) investigated whether LLMs can replace humans in evaluating texts. Specifically, they looked at open-ended story generation and adversarial attacks. They found that there were similar ratings between LLMs and human evaluators. Venkit et al. (2023) found that unbalanced sources of training data result in biased generations in GPT-2, and proposed strategies to reduce bias using adversarial triggers. Tang et al. (2022) presented EtriCA, a neural generation model which aims to remedy issues of relevance and coherence of generated texts. Lucy and Bamman (2021b) studied the bias existing in GPT-3's generated stories. Guo et al. (2023) have proposed a similar study specifically testing how similar text generated by ChatGPT is to text produced by human writers.

## 3 Methodology

To investigate the trustworthiness of children's stories generated by LLMs, we compare them with actual old and modern children's stories. We collect a diverse set of stories from different sources, including both older stories such as folktales, and more recent children's stories. We use both LLMs and instruction-following models to generate stories with different prompt lengths and instruction templates. As story generation is an open-ended problem with no reference text, we rely on other metrics instead of any automatic measure of evaluation like BARTScore (Yuan et al., 2021) or BERTScore (Zhang et al., 2019). We use various metrics to compare the generated stories with actual stories, including in-text statistics such as sentence length and a measure of toxicity in the text, as well as an evaluation of topics covered in these stories. Furthermore, we analyze and compare the grammatical structures of the stories using dependency structures extracted from both the original and the artificially generated stories.

In the following section, we describe the experimental setup, including details on the collected data, the story generation process, and the evaluation metrics used for comparison. Subsequently, we present the results obtained from our experimentation.

## 4 Experiments

### 4.1 Data

Our data consists of 132 original children's stories collected from various online sources and categorized into two categories: old and modern. The old stories generally include traditional children's stories like folktales and fairy tales, whereas the modern stories include more recent children's literature published after the year 2000. Both sets of original children's stories are comprised of English texts aimed at children between the ages of three and thirteen, with both data sets representing the full range of these target ages. Overall, 122 are classified as old stories, and the remaining 10 as modern stories. Specifically, the older stories were obtained via Project Gutenberg,[2] and the modern stories from various online platforms.[3] We use the old stories as a reference for the story generation task and compare the generated stories against both old and modern stories.

### 4.2 Story generation

We generate stories using language models and an instruction-following model.

---

| Model | Prompt Length | Count |
|---|---|---|
| OPT | First Sentence (OPT-Line) | 610 |
| | First 256-tokens (OPT-256) | 610 |
| | First 512-tokens (OPT-512) | 610 |
| LLaMA | First Sentence (LLaMA-Line) | 610 |
| | First 256-tokens (LLaMA-256) | 610 |
| | First 512-tokens (LLaMA-512) | 610 |
| | **Total** | 3660 |

Table 1: Breakdown of the stories generated using LLMs.

**Language Models** Our story generation task using LLMs uses two foundational language models: OPT (Zhang et al., 2022) and LLaMA (Touvron et al., 2023), with model sizes of 6.7 billion and 7 billion parameters, respectively. To generate stories, we provide a portion of each old story as context for the LLMs. Specifically, we use the first sentence, the first 256 tokens, and the first 512 tokens of each old story as a prompt. We use top-$k$ sampling-based decoding with $k$ set to 100 and generate five samples for each prompt, resulting in a total of 3660 generated stories. The breakdown of the generated stories along with the length of the prompt is given in Table 1.

**Instruction-following Models** For instruction-following story generation, we use Alpaca (Taori et al., 2023), which is an instruction-following model that is based on the LLaMA architecture and is fine-tuned using self-instruct (Wang et al., 2022). We use the Alpaca model based on the 7B variant of the LLaMA model. We use four different instruction templates to generate stories, two of which require a story title as input and two of which do not. For the templates that require a story title, we use the title of old stories as input. The templates are provided in Table 2. To generate stories, we use top-$k$ sampling-based decoding with $k$ set to 100 and generate five samples for each template, resulting in a total of 2440 generated stories with 610 stories per template.

### 4.3 In-text statistics

We compare various statistics derived from the text of the generated stories against those of actual stories. Specifically, we use two metrics: sentence length and Flesch reading ease score (Flesch, 1948).

**Flesch Reading Ease Score** The Flesch reading ease score (FRES) measures the readability of a text and is based on two factors: average sentence length and the average number of syllables per word. It provides a score between 0 and 100, with higher scores indicating easier readability. A Flesch reading ease score above 60 for a text indicates that it can easily be read by children up to the age of 15. The formula for calculating the FRES of a text is shown in Equation 1.

$$FRES = 206.835 - 1.015 \left( \frac{total\ words}{total\ sentence} \right) \\ - 84.6 \left( \frac{total\ syllables}{total\ words} \right) \quad (1)$$

### 4.4 Toxicity of text

Gehman et al. 2020 found that the LLMs can generate 'toxic' text from a very innocuous prompt and attribute this to a significant amount of offensive, factually unreliable, and otherwise toxic content in the training data of these models. We want to investigate the level of toxicity in our generated children's stories. Ideally, generated children's stories should be free of any toxic text.

We use Detoxify (Hanu and Unitary team, 2020), a BERT (Devlin et al., 2019) based toxic text detector, to identify the presence of toxic text in the generated children's stories. Detoxify generates score labels in the range of 0 to 1, assessing the toxicity of the text based on categories such as toxic, severely toxic, obscene, threat, insult, and identity hate. Specifically, we use detoxify for each sentence of our actual and generated stories to get toxicity measures across the six categories.

### 4.5 Topic Modeling

We also analyze the data for topic modeling using pyLDAvis (Tran, 2022). We compare the topics found in the data set of older stories with the LLM-generated stories. The older stories and the modern stories are also compared to assess whether there has been a shift in topics over time that would potentially influence topic properties in the LLM-generated stories. A probable diachronic shift in topics of stories geared towards young audiences also highlights the need to test the toxicity of generated stories, as seen in the previous section.

To avoid uninformative topics, the data is preprocessed to remove stopwords and names. All texts are categorized for specific topics using word clustering for a set of documents. Modeling is performed automatically without a predefined list of

| S.N. | Template |
|------|----------|
| T1 | Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.<br><br>### Instruction:<br>Write a short children's story given the title.<br><br>### Input:<br>TITLE<br><br>### Response: |
| T2 | Below is an instruction that describes a task. Write a response that appropriately completes the request.<br><br>### Instruction:<br>Write a short children's story.<br><br>### Response: |
| T3 | Below is an instruction that describes a task, paired with an input that provides further context. Write a response that appropriately completes the request.<br><br>### Instruction:<br>Write a children's story given the title.<br><br>### Input:<br>TITLE<br><br>### Response: |
| T4 | Below is an instruction that describes a task. Write a response that appropriately completes the request.<br><br>### Instruction:<br>Write a children's story.<br><br>### Response: |

Table 2: Templates used by Alpaca for story generation.

labels. The visualizations using pyLDAvis break down the topics based on the 122 older stories, the 10 modern stories, and the generated stories from OPT, LLaMA, and Alpaca.

## 4.6 Sentence structure

The structure of the sentences within a text can reveal the type or genre of the text. To analyze sentence structures, we construct a dependency tree for each sentence in both the original and generated children's stories. The dependency tree depicts the syntactic dependencies between the words in a sentence, effectively capturing the grammatical structure of the sentence. We then convert these dependencies into unlabeled directed graphs, preserving sentence structure while removing specific words. We then generate the Weisfeiler Lehman graph hash (Shervashidze et al., 2011) for each graph. The Weisfeiler Lehman hashes are identical for isomorphic graphs and strongly guarantee that non-isomorphic graphs will get different hashes. We compare the frequency of hashes to evaluate

the similarity between the sentence structure of the generated stories and the actual stories.

## 5 Generated stories follow modern trends but struggle with nuances

Figure 1 shows the box plot of sentence lengths for old and modern original stories, as well as for the generated stories. Being literary texts, children's stories do not strictly confine to formal English conventions and many contain sentences with higher word counts; so for clarity, we removed all the outliers from the plot. One interesting observation is that modern children's stories generally have shorter sentence lengths than older children's stories, which adheres to previous research that shows a trend of decreasing sentence length in print (Haussamen, 1994). The generated stories from OPT and LLaMA show an increase in sentence length as the prompt length increases. We hypothesize that these models learn the pattern of larger sentence length from the older stories used as context, which is then reflected in the generated text. However, stories
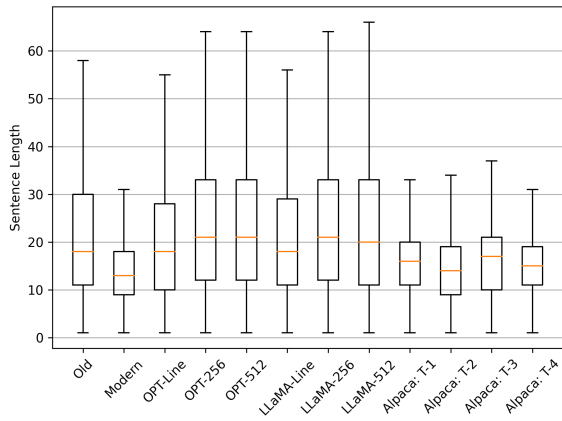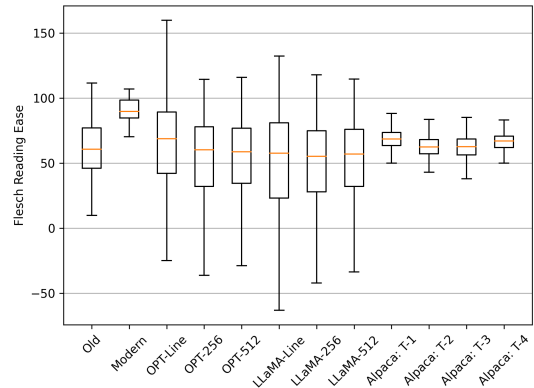
Figure 1: Comparison of sentence length in generated children's stories and actual children's stories. The generated children's stories exhibit shorter sentence lengths compared to the older original stories but are similar in sentence length to modern stories. Language models prompted with older stories tend to generate longer sentences following the patterns of the context that had been provided.
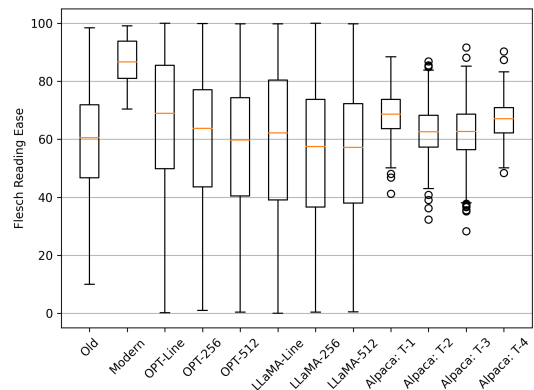
generated using the instruction-following model Alpaca, have sentence lengths similar to modern actual stories, indicating that language models may have been trained mostly on the newer text, and tend to generalize modern trends when instructed to generate text of a specific type.

The Flesch reading ease score is a statistical measure of the readability of a text and was optimized to be general enough at the time of its formulation, as can be seen with the constant values in equation 1. That is why, we may find FRES values not within the range of 0 to 100 as seen in Figure 2a. We also removed the outliers from the box plot in Figure 2a. Since we are not interested in exact values but in the general trend these values represent, we use the FRES values in the range of 0-100 and show their box plot in Figure 2b.

Our results from the Flesch reading ease score reveal several interesting observations. Firstly, we see that modern children's stories have a higher FRES than older stories, meaning that the modern ones are easier to read. This can be attributed to the fact that sentences are getting shorter and might have to do with simpler word selection. Secondly, we see that LLMs prompted with older stories tend to follow the pattern of the context and generate stories that are more difficult to read, as the context length increases. Finally, we see that the instruction-following model Alpaca generates stories that are easier to read compared to older original children's stories but are not as readable as



(a) FRES on all data (over 100 is undefined).



(b) FRES limited to well-defined 0-100 range.

Figure 2: Comparison of FRES in generated children's stories and actual children's stories : (a) FRES with all data and (b) FRES only in the range of 0 and 100. The generated children's stories are easier to read compared to older actual stories but are not as easy as modern original stories. Language models prompted with older stories tend to generate text that is more difficult to read, likely because they follow the patterns in the prompts.

modern children's stories. We posit that this observation can be attributed to the fact that LLMs used in our study are generic models, and the instruction following model is also only fine-tuned for general instructions rather than instructions specific to children's story generation.

Overall, we see that modern children's stories are easier to read than older children's stories. As most of the training data for LLMs comes from newer text, the model tends to follow the trend of modern children's stories in their generated text for sentence length and word selection. However, it should be noted that these models are not fine-tuned for children's stories generation, and therefore may not capture the nuances of children's stories resulting in stories that might be difficult to read for intended readers.
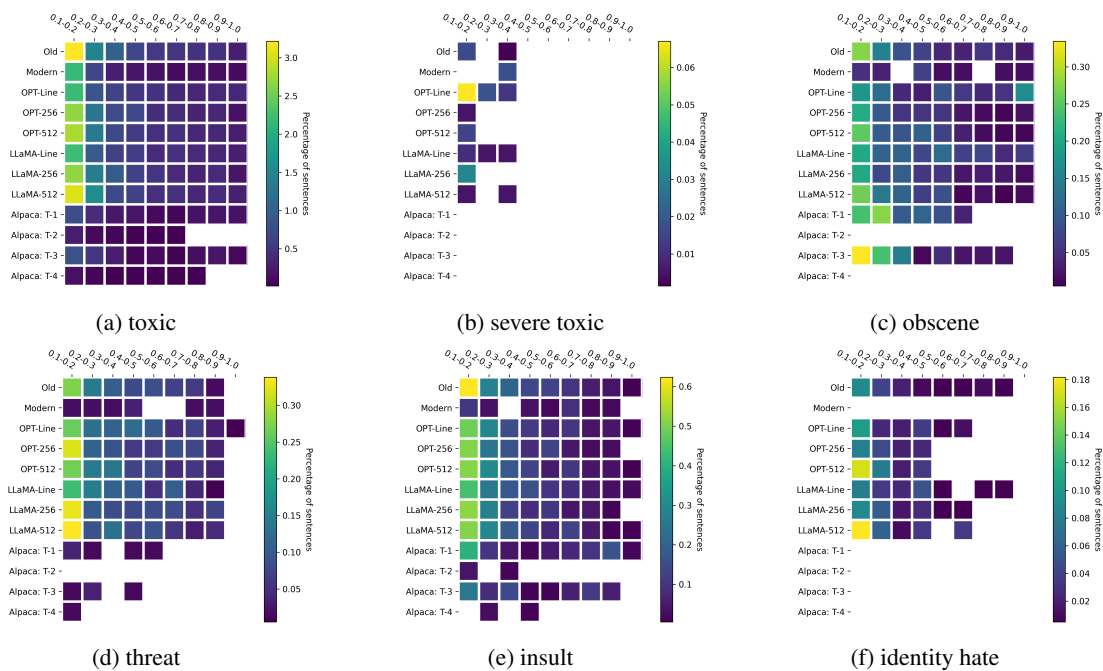
Figure 3: Various toxicity measures for the actual and generated stories. Each cell in a subplot represents the percentage of sentences rated on a toxicity scale, with x-axis values indicating the toxicity level. Values for ratings in the range of 0-0.1 have been omitted from the plots for clarity.

# 6 Generated stories may contain toxic text

Our analysis of toxicity in actual and generated stories reveals several noteworthy findings. We present the toxicity measures for both actual and generated stories in Figure 3. Notably, we find that older stories tend to be more toxic than modern stories across all toxicity measures. This trend is not solely due to the smaller sample size of modern actual stories, as we have normalized the toxicity ratings to ensure an accurate comparison. Rather, it suggests that writers are becoming more mindful of the language they use in children's literature. Although modern stories are less toxic compared to older stories, we still observe some level of toxicity in them. This toxicity in modern actual stories is often related to the narrative of the story. For example, threats and insults might be needed for some stories, but identity hate is not appropriate for children's stories. It is noteworthy that modern stories do not have toxic text related to identity hate but older stories do.

Similar to our previous observation, we see that LLMs tend to learn patterns from the context they are provided with. As evident from the stories generated by OPT and LLaMA, we see that the toxicity aligns with older stories and gradually increases with an increase in the length of the context. The

stories generated using the instruction-following model Alpaca tend to be less toxic and mostly resemble modern stories. However, stories generated using the T1 and T3 templates have a lot of obscene text compared to stories generated using T2 and T4, which have none. As shown in Table 2, T1 and T3 take the title as input whereas T2 and T4 do not. It is possible that the model remembered the story title and generalized the patterns of the story or generalized to some other text in the template, leading to the generation of obscene text. This finding is consistent with Gehman et al. 2020, who suggest that children's stories generated by LLMs can contain highly toxic text despite an innocuous prompt.

Our analysis of toxicity in original and generated stories reveals that older stories tend to be more toxic than modern ones, that LLMs can learn toxic patterns from context leading to the generation of toxic text, and that LLMs can even generate toxic text from a very innocuous prompt. These findings suggest that further work is needed to make LLMs useful as tools for generating age-appropriate children's literature.

# 7 Generated stories share main topics with original stories

After preprocessing the data, the original stories were found to have four major topics. All of the topics tended to share the existence of some small character. The first topic mentions elements such as time, goodness, and greatness, and the presence of words like head, round, night, and water likely indicate specific scenes or settings within the narrative. The second topic contained new elements like a prince, the color white, a girl, and eyes. These additional keywords suggest different perspectives within the overarching narrative. The third topic introduces elements like a house and a heart. Like the previous topics, it shares mentions of a little character, time, goodness, and a prince. The difference between 'house' and 'heart' could indicate a change in the setting or moral of the narrative. The last topic introduces new elements of wolf, people, eyes, and a mother. These keywords might suggest narratives that introduce new characters and themes. Overall, these topics provide insight into the underlying themes present in the older 122 stories in the data set. The topics revolve around narratives involving a small character, time, goodness, and various other elements such as princes, nights, water, girls, and wolves.

Comparatively, the topics of the generated stories obtained from OPT, LLaMA, and Alpaca show minor differences. The first topic suggests a narrative that involves characters like kings, mothers, princes, and princesses. It also mentions elements of time, goodness, greatness, and shadow. The prince, princess, and shadow hint at the fairy tale or fantasy theme. The second topic shares similarities with the previous topic, with a focus on little, prince, time, goodness, and greatness, but it also introduces new elements like eyes, houses, heads, and the color white. These additions suggest different scenes, perhaps removed from the monarchy or castle theme, and suggest a different narrative. The third topic seems to center around family dynamics, with mentions of mothers, fathers, and children. It also includes keywords relating to time, goodness, night, and poverty. This suggests a change in the narrative away from the fantasy-focused topic. The last topic includes keywords like little, time, and goodness. It includes elements of fathers, eyes, and houses. The presence of 'long' and 'night' suggests a different tone or atmosphere within the narrative. These general results show remarkable similarity
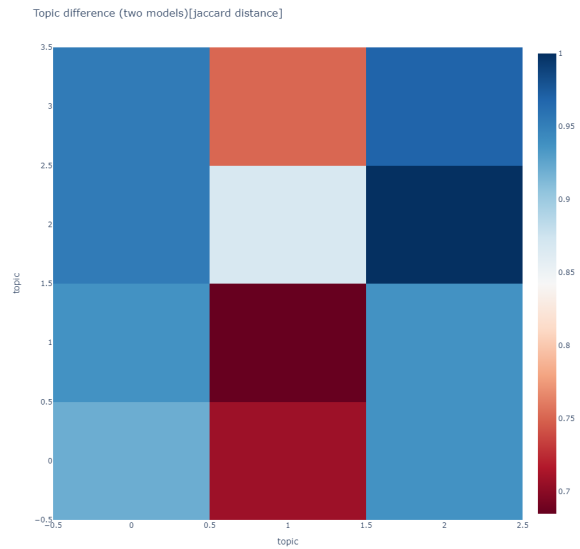


Figure 4: Comparison of topics in generated children's stories and actual children's stories. The plot shows that the most shared topics (x:2, y:2) include 'white', 'world', 'great', 'water', 'black', 'house', 'little', 'king', 'called', and 'good'. The least shared topics (x:1, y:1) include 'heart', 'head', 'poor', 'house', 'looking', 'children', 'good', 'young', 'lady', and 'night'.

with the data set on which the LLMs were trained. The topics revolve around narratives involving characters such as kings, princes, mothers, fathers, and children. The topics also touched upon topics of time, goodness, greatness, poverty, and setting elements of houses, nights, and the color white.

As with toxic content testing, we ran topic modeling for a small number (10) of modern stories in order to compare the general topics that are currently aimed at children. The first topic includes keywords related to spatial orientation (right, inside, door, left), objects (ream, head, frog), time, and actions (started). The keyword 'eyes' may suggest a focus on visual perception or observation. The second topic emphasizes time, objects (ream, door), spacial orientation (right, inside), a frog, a head, fairies, and 'need'. The presence of fairies introduces a fantastical or imaginative element to the topic. The third topic revolves around time, spacial orientation (right, inside, door), physical attributes (head, eyes, long, hand), and a frog. The inclusion of 'long' might suggest a temporal or duration-related aspect. The last topic highlights time, spatial orientation (right, inside, door), objects (ream, frog), physical attributes (head, eyes, small), and the action of starting something. The modern stories' topic modeling results suggest a re-

curring theme involving concepts such as time, spatial orientation, objects, and actions. Each topic emphasizes different aspects and introduces additional elements like fairies or physical attributes. Figure 4 represents the level of similarity and difference between the real stories and the LLM-generated stories. There are greater similarities between these stories than there appear to be differences in the main topics.

As expected, the results of the topic modeling showed similarities between the original 122 stories in the training corpus and the stories generated by the LLMs. These stories shared fairy tale and fantasy elements as well as topics of goodness, greatness, time, and setting elements of night, houses, and the color white. Once we compare this with the modern stories, we see that the focus of the small data set we have is similarly focused on time and fairies, but has more topics relating to spatial orientation. We are likely seeing a change in the content of stories written for children. With only ten modern stories, we cannot reliably generalize over all stories, but we noticed tendencies such as that the modern story set did tend to involve more overtly educational elements aimed at younger age groups when compared to the older stories.

## 8 Generated stories do not have similar sentence structure to original stories

Table 3 shows the percentage of overlapping Weisfeiler Lehman hashes between the dependency tree graphs of sentences generated by various models and those actual children's stories, both old and modern. We also got an overlap of **35.57** percentage between old and modern actual stories, which is greater than all the values in Table 3. This shows that the structure of sentences in children's literature has changed over time, which supports our earlier findings that children's literature has undergone noticeable changes over time.

Additionally, we observe a higher percentage of overlap between old original stories and the stories generated by OPT and LLaMA, which again aligns with our earlier findings that LLMs learn from their context. Furthermore, for the stories generated by OPT and LLaMA, we see an average overlap of 30% with modern stories, which can be attributed to the fact that these models were trained on a dataset consisting of recent text.

The stories generated by Alpaca have a slightly higher overlap with modern stories compared to

| Model | Percentage overlap with | |
| | Old stories | Modern stories |
| --- | --- | --- |
| OPT-Line | 34.82 | 34.21 |
| OPT-256 | 31.37 | 28.88 |
| OPT-512 | 32.49 | 29.89 |
| LLaMA-Line | 34.23 | 33.64 |
| LLaMA-256 | 32.14 | 29.82 |
| LLaMA-512 | 32.27 | 30.73 |
| Alpaca: T-1 | 17.31 | 20.37 |
| Alpaca: T-2 | 14.67 | 17.52 |
| Alpaca: T-3 | 15.20 | 16.92 |
| Alpaca: T-4 | 15.41 | 17.84 |

Table 3: Overlap of the hashes of the dependency tree graph of the sentences in generated stories against old and modern actual stories.

old stories, but the percentage overlap in sentence structures is still relatively low ( $\leq 20\%$ ). Given that the old and modern actual stories share around 35% of the same sentence structures, we expected Alpaca's generated stories to overlap more with modern stories. But since Alpaca is a generic model fine-tuned for instruction-following and not solely trained or fine-tuned on children's literature, it seems plausible that it would not be capable of fully generalizing over sentence or grammatical structures observable in children's literature.

## 9 Conclusion and Future Work

Our study examines the trustworthiness of children's stories generated by large language models. While these generated stories may share similar topics and patterns with actual stories, they fail to capture all the nuances present in children's literature, and may even contain toxic material that is inappropriate for children. Based on our findings, we conclude that LLMs are not yet appropriate for generating high-quality children's literature. Moving forward, we plan to extend our work by implementing reinforcement learning with both automatic and human feedback to improve the quality of LLM-generated children's stories.

# References

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Cheng-Han Chiang and Hung-yi Lee. 2023. Can Large Language Models Be an Alternative to Human Evaluations? ArXiv:2305.01937 [cs].

Elizabeth Clark, Anne Spencer Ross, Chenhao Tan, Yangfeng Ji, and Noah A. Smith. 2018. Creative Writing with a Machine in the Loop: Case Studies on Slogans and Stories. In *23rd International Conference on Intelligent User Interfaces*, pages 329–340, Tokyo Japan. ACM.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Rudolph Flesch. 1948. A new readability yardstick. *Journal of applied psychology*, 32(3):221.

Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. 2020. RealToxicityPrompts: Evaluating neural toxic degeneration in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3356–3369, Online. Association for Computational Linguistics.

Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang, Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng Wu. 2023. How Close is ChatGPT to Human Experts? Comparison Corpus, Evaluation, and Detection. *arXiv.org*. Place: Ithaca Publisher: Cornell University Library, arXiv.org.

Laura Hanu and Unitary team. 2020. Detoxify. Github. https://github.com/unitaryai/detoxify.

Brock Haussamen. 1994. The future of the english sentence. *Visible language.*, 28(1).

Enkelejda Kasneci, Kathrin Sessler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank Fischer,

Urs Gasser, Georg Groh, Stephan Günnemann, Eyke Hüllermeier, Stepha Krusche, Gitta Kutyniok, Tilman Michaeli, Claudia Nerdel, Jürgen Pfeffer, Oleksandra Poquet, Michael Sailer, Albrecht Schmidt, Tina Seidel, Matthias Stadler, Jochen Weller, Jochen Kuhn, and Gjergji Kasneci. 2023. ChatGPT for good? On opportunities and challenges of large language models for education. *Learning and Individual Differences*, 103:102274.

Li Lucy and David Bamman. 2021a. Gender and representation bias in GPT-3 generated stories. In *Proceedings of the Third Workshop on Narrative Understanding*, pages 48–55, Virtual. Association for Computational Linguistics.

Li Lucy and David Bamman. 2021b. Gender and Representation Bias in GPT-3 Generated Stories. In *Proceedings of the Third Workshop on Narrative Understanding*, pages 48–55, Virtual. Association for Computational Linguistics.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.

Rob Price. 2016. Microsoft deletes racist, genocidal tweets from AI chatbot Tay - Business Insider.

Evgeniia Razumovskaia, Joshua Maynez, Annie Louis, Mirella Lapata, and Shashi Narayan. 2022. Little red riding hood goes around the globe:crosslingual story planning and generation with large language models.

Nino Shervashidze, Pascal Schweitzer, Erik Jan Van Leeuwen, Kurt Mehlhorn, and Karsten M Borgwardt. 2011. Weisfeiler-lehman graph kernels. *Journal of Machine Learning Research*, 12(9).

Nisha Simon and Christian Muise. 2022. TattleTale - Storytelling with Planning and Large Language Models.

Yuqian Sun, Ying Xu, Chenhang Cheng, Yihua Li, Chang Hee Lee, and Ali Asadipour. 2023. Explore the future earth with wander 2.0: Ai chatbot driven by knowledge-base story generation and text-to-image model. In *Extended Abstracts of the 2023 CHI Conference on Human Factors in Computing Systems*, CHI EA '23, New York, NY, USA. Association for Computing Machinery.

Chen Tang, Chenghua Lin, Henglin Huang, Frank Guerin, and Zhihao Zhang. 2022. EtriCA: Event-Triggered Context-Aware Story Generation Augmented by Cross Attention. ArXiv:2210.12463 [cs].

Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models.

Khuyen Tran. 2022. pyLDAvis: Topic Modelling Exploration Tool That Every NLP Data Scientist Should Know.

Pranav Narayanan Venkit, Sanjana Gautam, Ruchi Panchanadikar, Ting-Hao 'Kenneth' Huang, and Shomir Wilson. 2023. Nationality bias in text generation.

Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2022. Self-instruct: Aligning language model with self generated instructions.

Kyle Williams. 2019. Neural lexicons for slot tagging in spoken language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Industry Papers)*, pages 83–89, Minneapolis, Minnesota. Association for Computational Linguistics.

Jingwen Xiang, Zoie Zhao, Mackie Zhou, Megan McKenzie, Alexis Kilayko, Jamie C Macbeth, Smith Edu, Smith Edu, Smith Edu, Smith Edu, Smith Edu, and Smith Edu. 2018. Interleaving a Symbolic Story Generator with a Neural Network-Based Large Language Model.

Weipeng Yang. 2022. Artificial Intelligence education for young children: Why, what, and how in curriculum design and implementation. *Computers and Education: Artificial Intelligence*, 3:100061.

Ann Yuan, Andy Coenen, Emily Reif, and Daphne Ippolito. 2022. Wordcraft: Story Writing With Large Language Models. In *27th International Conference on Intelligent User Interfaces*, IUI '22, pages 841–852, New York, NY, USA. Association for Computing Machinery. Event-place: Helsinki, Finland.

Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021. Bartscore: Evaluating generated text as text generation. *Advances in Neural Information Processing Systems*, 34:27263–27277.

Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. Opt: Open pre-trained transformer language models.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert.

# On Text Style Transfer via Style-Aware Masked Language Models

**Sharan Narasimhan, Pooja Shekar, Suvodip Dey, Maunendra Sankar Desarkar**

Department of Computer Science and Engineering
Indian Institute of Technology Hyderabad, India

`{sharan.n21,poojashekar15,suvodip15}@gmail.com`
`maunendra@iith.ac.in`

## Abstract

Text Style Transfer (TST) involves transforming a source sentence with a given style label to an output with another target style meanwhile preserving content and fluency. We look at a fill-in-the-blanks approach (also referred to as prototype editing), where the source sentence is stripped off all style-containing words and filled in with suitable words. This closely resembles a Masked Language Model (MLM) objective, with the added initial step of masking only relevant style words rather than BERT's random masking. We show this simple MLM, trained to reconstruct style-masked sentences back into their original style, can even transfer style by making this MLM "Style-Aware". This simply involves appending the source sentence with a target style special token. The Style-Aware MLM (SA-MLM) now also accounts for the direction of style transfer and enables style transfer by simply manipulating these special tokens. To learn this n-word to n-word style reconstruction task, we use a single transformer encoder block with 8 heads, 2 layers and no auto-regressive decoder, making it non-generational. We empirically show that this lightweight encoder trained on a simple reconstruction task compares with elaborately engineered state-of-the-art TST models for even complex styles like Discourse or flow of logic, i.e. Contradiction to Entailment and vice-versa. Additionally, we introduce a more accurate attention-based style-masking step and a novel "attention-surplus" method to determine the position of masks from any arbitrary attribution model in O(1) time. Finally, we show that the SA-MLM arises naturally by considering a probabilistic framework for style transfer. [*]

## 1 Introduction

Text Style Transfer can be thought of as a form of Controllable Language Generation (Hu et al.,

2017) with tighter criteria. Dathathri et al. (2020) show that a classifier trained on the final embeddings of any arbitrary large language model to learn class labels of any dataset, can guide subsequent generations to resemble this dataset's style by backpropagating signals from the classifier to the LLMs activation layers. This approach, as well as other generational approaches, generally suffer from the lack of the model's ability to preserve content. This content preservation criteria, where only style influencing words/phrases must be affected, forms the critical challenge in the Style Transfer task. Learning style transfer from a parallel dataset is easy, i.e. where the output sentence for every target style is known and the model learns a straightforward mapping of input-output sentence pairs in a supervised manner. Like most work, we assume the more realistic case, where the dataset is non-parallel, and the task is unsupervised. The style transfer task is equivalent to the task of estimating the hypothetical parallel dataset from a non-parallel or "partially-observed" parallel dataset. We use this notion in later sections to elucidate the mechanisms which enable our approach.

Style is usually represented using a data-driven approach, i.e. the class labels present in an annotated dataset. One could argue that for a majority of style transfer applications, only a subset of words participate in giving the input sentence its corresponding style. As examples, we can think of this being true in the case of Sentiment, Toxicity, Formality, Politeness etc, where a few word edits can lead to a style change while the other attributes of the sentence are disaffected. As expected, this type of approach has been vastly used in the style transfer task with good success. This prototype editing approach consists a) ranking tokens according to how much they affect the underlying style b) determining which subset of tokens to mask to produce a style absent sentence, and c) transforming this style absent sentence into one that contains the

---

target style.

We show that a single self-attention encoder block trained to reconstruct style-masked sentences to the original versions using a non-parallel dataset also learns style transfer automatically. This method, resembling a masked language model, compares with state-of-the-art models albeit having a simpler training process, lesser parameters and without using pre-trained language models. We consider a novel discourse manipulation task and show that the SA-MLM outperforms another state-of-the-art model in this respect.

## 2 Learning Style Transfer from only Reconstruction

Many similarities appear between the editing approach for style transfer and BERT's masked language model objective. The only difference being that we reconstruct the original sentence from its style-masked version, rather than random or perplexity-based-masking sentences in typical MLMs. We postpone discussing our approach for style masking to the next section. In this section, we explore the question *"How does an MLM trained only on a reconstruction task automatically gain the ability to perform the unseen style-transfer task?"*.

**An Ensemble Approach.** Intuitively, we can reason that just training a model to reconstruct sentences from style-masked versions ("the food was <blank>" to the "the food was good") will not work. Since a) this does not give context to the model about styles present and b) does not give us any knob to control the output style. A simple solution would be training $n$ different models, each to reconstruct one particular style. Style transfer can be performed by feeding a style-masked sentence to that corresponding target style model. This however a) is not computationally efficient or scale-able b) limits the learning of each model to only a fraction of the dataset, thereby over-fitting to that target domain. We perform the style transfer task using this ensemble model consisting of two generational encoder-decoder models, denoted as *Ensemble* in Tables 3, 5, 4 and 6.

**A Single Style-Aware Encoder.** As an alternative to training $n$ models, we can train a single model and contextualize style by concatenating a <target_style> special token to the input style-masked sentence. Training on the reconstruction task this way (e.g. "the food was good" from "the

| Task | Positive to Negative | Contradiction to Entailment |
|------|----------------------|------------------------------|
| **Input** | This movie is by far one of the best urban crime dramas i've seen . | a woman is sitting outside at a table using a knife to cut into a sandwich . a woman is sitting inside |
| **Style Masked** | This movie is by <mask> one of the <mask> urban crime <mask> i've seen . | a woman is sitting outside at a table using a knife to cut into a sandwich . a woman <mask> <mask> <mask> |
| **Output** | This movie is by far one of the worst urban crime garbage i 've seen . | a woman is sitting outside at a table using a knife to cut into a sandwich . a woman is a outside |

Table 1: Examples of Sentiment and Discourse style transfer by the SA-MLM on the IMDb and SNLI datasets respectively.

food was <blank>. <positive>") allows it to infer the target style needed and reconstruct accordingly. This allows us to perform style transfer by simply manipulating the target style token. The model estimates the unseen portion of the hypothetical parallel dataset. Furthermore, we hypothesize that a word-by-word generational approach using an autoregressive decoder might degrade performance in content preservation criteria. We juxtapose the performances of these approaches, denoted as *Generational* and *Encoder-only* in Tables 3, 5, 4 and 6.

To summarize, a simple modification (in the form of label concatenation) to the traditional MLM task enables style transfer by training solely on a reconstruction/denoising task. This forms the basis of our approach.

## 3 Masking Style with Attention

We explore the question, *"What constitutes a good style-masking step?"*. Intuitively, we can reason that our style-masking approach must a) produce accurate attribution scores for each token and b) use an appropriate masking policy that determines which tokens to mask using these attribution scores. The final style-masked sequences (input to the SA-MLM) must be a) completely devoid of style information and b) accurate, i.e. not done at the expense of content information.

### 3.1 Accurate Attribution Scores

Many prototype editing methods use Vanilla Attention (VA) as attribute scores (Wu et al., 2019b;

Zhang et al., 2018; Wu et al., 2020). It has been shown that attention is not explanation, i.e. these attribution scores do not align with human interpretability (Jain and Wallace, 2019). VA does not correlate well with other well-known attribution methods (such as Integrated Gradients Sundararajan et al. (2017)). We instead use "Explainable Attention" (EA) scores from a Diversity-LSTM classifier (Mohankumar et al., 2020; Nema et al., 2017) which have been shown to correlate better with other attribution methods as well as human judgement. We discuss more about the Diversity LSTM in section A.6 of the appendix. We also quantitatively compare the efficacy of the style-masking step between EA and VA in Table 2.

## 3.2 An Accurate Masking Policy

Even with having accurate attribution scores using explainable attention, effective style-masking requires careful selection of a policy which satisfies certain criteria. The primary criteria being that only tokens which significantly contribute to the style of a sentence must be masked, and other tokens must be ignored to ensure content is also preserved. Similar to the masking policy in Wu et al. (2020), it is natural to consider a "top $k$ tokens" scheme in which the top $k$ tokens with highest attribution are masked. However, this static approach fails for sentences which do not have exactly $k$ style-contributing tokens, leading to either partial style masking or erroneous masking of content tokens. For the same reason, even a sentence length aware scheme such as "top 15%" masking fails. Furthermore, all such policies require sorting, leading to $O(n \log n)$ time complexity for style masking of each sentence in a batch.

**Attention Surplus.** Let $A = \{A_i \dots A_n\}$ denote the attention distribution of a sentence of size $n$. Intuitively, we can reason that all "special" tokens which might contribute more to style should have an attribution greater than the average base attribution of the sentence, given by $A^{mean} = 1/n$. Generalising this further, we refer to tokens with $A^i \geq A^{baseline}$ as tokens with "attention surplus" with respect to a sentence-length sensitive baseline attention $A^{baseline}$ given by:

$$A^{baseline} = (1 + \lambda_\epsilon) * A^{mean} \qquad (1)$$

where $\lambda_\epsilon$ is a hyperparameter of range $0 - 1.0$. This chosen threshold $A^{baseline}$ is sensitive to the

sentence length as well and subsequently ensures that the number of style-significant tokens can be dynamically determined, without need of an elaborate algorithm. As a sanity check, we observe that even in the adversarial case where all tokens might be equally important to style, $A$ resolves into a UniformDistribution(n) and our policy correctly resorts to masking all tokens[†]. Let $Mask$ denote the token mask matrix of size $n$ initialised with zeros.

$$\text{Mask}[A_i \geq A^{baseline}] = 1 \qquad (2)$$

Using a vectorised batch-wise approach, we can style-mask an entire input batch in just O(1) complexity, compared to sorting-based approaches which take $O(\text{batch\_size} \cdot n \log n)$ .

## 4 Related Work

Recent work focuses on various common paradigms such as disentanglement (Hu et al., 2017; Shen et al., 2017), cycle-consistency losses (Yi et al., 2020; Luo et al., 2019; Dai et al., 2019; Liu et al., 2021), induction (Narasimhan et al., 2022; Shen et al., 2020). Jin et al. (2021) and Hu et al. (2020) provide surveys detailing the current state of style transfer and lay down useful taxonomies to structure the field. In this section, we only discuss contemporary work similar to ours (prototype editing approaches) assuming the same unsupervised setting.

Li et al. (2018) present the earliest known work using the prototype editing method, in which a "delete" operation is performed on tokens based on simple count-based methods, and the retrieval of the target word is done by considering TF-IDF weighted word overlap. Malmi (2020) first train MLMs for the source and target domains and perform style transfer by first masking text spans where the models disagree (in terms of perplexity) the most, and use the target domain MLM to fill these spans. Wu et al. (2019b) introduce the Attribute-Conditional MLM, which most closely aligns to the working of the SA-MLM, also uses an attention classifier for attribution scores, a count and frequency-based method to perform masking, and a pretrained BERT model fine-tuned on the style transfer task. Lee (2020) and Wu et al. (2020) also follow roughly the same pipeline but uses

---

[†]Assuming $\lambda_\epsilon = 0$, whereas in practice we find $\lambda_\epsilon = 0.15, 0.5$ giving optimal masking for the Sentiment and Discourse TST, respectively. More is discussed in 7.2

a generational transformer encoder-decoder approach and also fine-tunes using signals from a pretrained classifier. Wu et al. (2019a) uses a hierarchical reinforced sequence operation method is used to iteratively revise the words of original sentences. Madaan et al. (2020) uses n-gram TF-IDF based methods to identify style tokens and modify them as either "add" or "replace" TAG tokens, which are then substituted by the decoder to perform style transfer. Similar to the SA-MLM, (Xu et al., 2018) also uses attribution-based methods from a self-attention classifier. However, they use an LSTM (Hochreiter and Schmidhuber, 1997) based approach, one to generate sentences from each domain. (Reid and Zhong, 2021) performs unsupervised synthesis to create a pseudo-parallel dataset and uses multi-span editing techniques to fill in the style using a fine-tuned pretrained language model.

# 5 Methodology

In this section, we describe the working of the SA-MLM during training and inference using a formal probabilistic framework.

## 5.1 Notation

Let $S$ denote the set of all style labels for a annotated dataset $D$ of the form $\{(x_0, l_0), (x_1, l_1) \ldots (x_n, l_n)\}$ where $x_i$ denotes the input sentence and $l_i \in S$ denotes the label corresponding to $x_i$. The set of all sentences of style $s$ in $D$ is denoted by $\hat{x^s} = \{x_j : \forall j \ where (x_j, s) \in D\}$. We use a special meta label $m_s$ to represent the "style-masked" label class having $s$ the original style. Subsequently, $x^{m_s}$ refers to the set of all style-masked sentences with source style $s$. The set of all style-masked sentences from $D$ is given by $x^m = Union(x^{m_s} : \forall s \in S)$.

## 5.2 A Probabilistic framework

For the sake of convenience in notation, we assume binary style labels, $S = \{0, 1\}^\ddagger$. We assume that a non-parallel dataset is a partially observed hypothetical parallel dataset. The SA-MLM, therefore, has to estimate the unseen half of this hypothetical parallel dataset. We follow the assumption that every output sentence with a style $s$ is a result of:

- sampling from a latent style-masked prior, $p(x^{m_s})$,

---

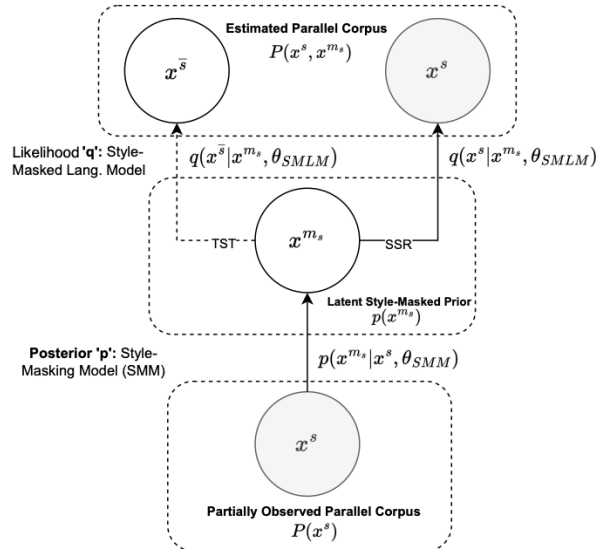‡In theory, this can be extended to any number of styles.



Figure 1: Probabilistic overview of our style transfer method.

- which we get from style masking input sentences, $p(x^{m_s}|x^s, \theta_{SMM})$ (posterior), and is then

- reconstructed to form the sentences with target styles using the SA-MLM, $q(x^s, x^{\bar{s}}|x^{m_s})$ (likelihood).

Style transfer is equivalent to estimating the unseen half of the hypothetical parallel $q(x^{\bar{s}}|x^{m_s})$. The overall model is summarised in Fig. 1.

## 5.3 SA-MLM architecture: A single self-attention block

The SA-MLM, $q(x^s, x^{\bar{s}}|x^{m_s})$ in this case, consists of a lightweight Transformer Encoder (a single self-attention block) with 2 layers, 8 heads and embeddings of size 512. To learn style transfer, training on the reconstruction task is sufficient, e.g. outting "The food was good" from a style masked input "The food was <blank>. <positive>".

## 5.4 Training objective

Our model, when trained only on the reconstruction task i.e. $q(x^s|x^{m_s})$, automatically learns style transfer i.e. $q(x^{\bar{s}}|x^{m_s})$. The intuition for why this is, is given in Section 2. Strictly speaking, this behaviour of automatically learning an unseen task is the result of two features we adopt, a) a single model with a common latent prior for all styles and b) the presence of target style information in the input sentences. Our model, consisting of a single
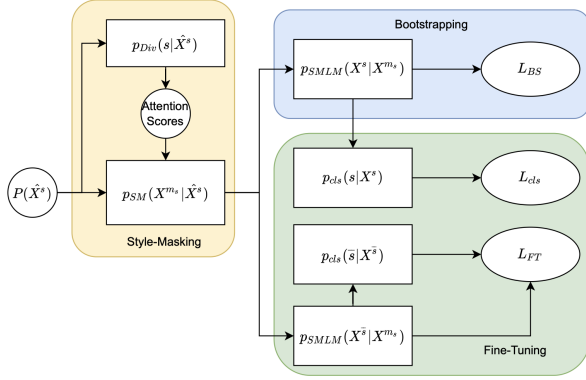
Figure 2: Overview of model architecture considering sentences of style $s$ i.e. $x^s$. In reality, this is applied all styles $x^s, \forall s \in S$

self-attention encoder block, minimizes the NLL reconstruction loss.

$$L_{recon}(\theta_{encoder}) = -\log q_{encoder}(x^s|x^{m_s})$$

## 5.5 Fine-tuning

The presence of the special target style token during the reconstruction task forces the encoder to try and ensure that the target style is present in the output. To further enforce this behaviour, we fine-tune the encoder for one epoch using techniques similar to those found in (Liu et al., 2021). While training on the reconstruction task, we simultaneously train a classifier to predict target style labels using final layer embeddings of the reconstructed output [§].

$$L_{cls}(\theta_{cls}) = -\log p_{cls}(s|x^s)$$

For fine-tuning, this classifier then provides supervision signals to the encoder with respect to the style transfer accuracy to further enforce the outputs to align with the target style. This is formulated as min-max objective $L_{FT}$ between the classifier $cls$ and the $encoder$:

$$\min_{\theta_{encoder}} \max_{\theta_{cls}} -log p_{cls}(\overline{s}|x^{\overline{s}})$$

## 6  Datasets and Tasks

We report the split and label-wise statistics of each dataset in Table 14 of the appendix.
**Sentiment Style Transfer:** Following many past studies, we evaluate our model for the sentiment style transfer task using three review datasets, Yelp,

Amazon and IMDb. All three datasets are annotated with two labels corresponding to positive or negative reviews and are non-parallel. [¶]

**Discourse Style Transfer:** Some style transfer tasks are more complex than others and have different levels of granularity (Lyu et al., 2021). To show that our seemingly simplistic approach can perform more cognitive tasks, we introduce the Discourse style transfer task by performing style transfer on the SNLI (Bowman et al., 2015) dataset. Each instance in the dataset consists of two sentences, which either contradict, entail or are neutral (no relationship). We consider the task of manipulating the discourse or "flow of logic" between two sentences, i.e. from contradiction to entailment and vice-versa. Unlike the sentiment task, which is "intra-sentence", where the style can be attributed to a select set of words, the discourse task is "inter-sentence" and requires the model to be cognizant of the context (especially for the Contradiction to Entailment Task) and detect the flow of logic.

## 7  Analysis of Style-Masking approach

In this section, we evaluate and justify our choice of style-masking architecture, i.e. "Explainable Attention" + "Attention-Surplus" masking policy.

### 7.1  Analysis of Various Attribution Methods

We consider various other attribution methods for our analysis i.e. Vanilla Gradients, Integrated Gradients (Sundararajan et al., 2017), Vanilla Attention, Attention * X (or inputs) and Explainable Attention (Mohankumar et al., 2020). We do not consider techniques such as LIME (Ribeiro et al., 2016), LRP (Bach et al., 2015), DeepLIFT (Shrikumar et al., 2017) as they are relatively more computationally expensive during inference time. For the style-masking policy, we use our "attribution-surplus" to determine which tokens to mask, as mentioned in Section 3.

In Table 2, we compute the Accuracy% and s-BLEU on the final style-masked sequences produced by each attribution method on the test split for all four datasets. We can reason that an ideal style-masking method should be able to produce sentences that completely mask out style, thereby fooling a pretrained classifier (minimizing its Accuracy%) and also preserving content information (maximizing the s-BLEU between the source

---

[§]We use a single feed-forward layer with input being the average of the last layer embeddings of $X_i^s$ (excluding the meta-labels).

[¶]For Yelp, Amazon and IMDb, we used the pre-processed version specified in https://github.com/yixinL7/Direct-Style-Transfer.
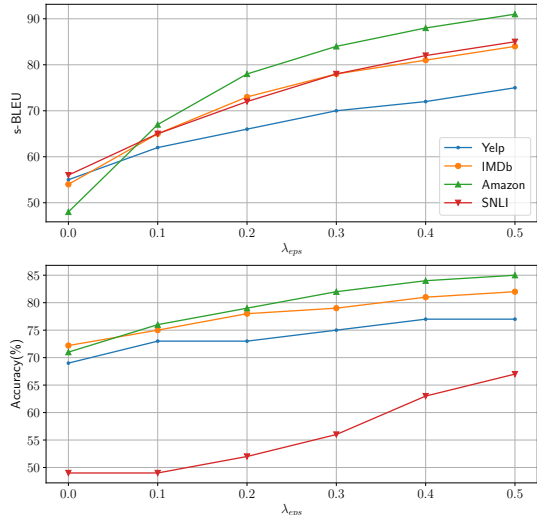
Figure 3: Effect of $\lambda_\epsilon$ on the resultant style-masked sentences using our "EA+AS" method. We compute s-BLEU(Top) and Accuracy%(Bottom) of the style-masked sentences on the test split of each dataset.

and style-masked sentences). We see that though Vanilla Attention is able to generally produce the lowest Accuracy%, however it does so at the expense of preserving content, reflected as lower s-BLEU compared to Explainable Attention, which, on the other hand, has the best content-preserving style masking throughout all datasets and comes as a close second in terms of Accuracy %. Other gradient-based methods do not perform favourably in any aspect.

## 7.2 Effect of $\lambda_\epsilon$ on Style-Masking

We can intuitively reason that the s-BLEU metric of the style-masked sentences serves as a rough upper bound for the s-BLEU we can potentially achieve on the output sentences after style transfer. As expected (from Eq. 1), we observe a positive correlation between $\lambda_\epsilon$ and both s-BLEU and Accuracy% as seen in Fig. 3. It is desirable to carefully choose $\lambda_\epsilon$ to be high enough to boost future s-BLEU scores on the style transfer task and also ensure that the sentences are sufficiently style masked with low Accuracy% scores. On manual checking, we observed that $\lambda_\epsilon = 0.15$ served best to accurately style-mask sentences for the Yelp, IMDb and Amazon datasets. SNLI required a higher $\lambda_\epsilon$ of 0.5 to ensure content information was preserved appropriately.

## 8 Experiments

We perform style transfer on all the datasets and analyse the results. Henceforth, we refer to the *Encoder-only* variant (described in Section 5.3) as

our "SA-MLM" flagship model while comparing with baselines in Tables 3, 5, 4 and 6. The other two SA-MLM variants i.e. *Ensemble* and *Generational* described in Section 2, serve as ablation studies to support our hypothesis that a single non-generational self-attention encoder only approach is optimal for style transfer.

## 8.1 Automatic Evaluation Metrics

Typically, metrics for style transfer include Style Transfer Accuracy (using a pretrained classifier), BLEU for content preservation, and perplexity (using a pretrained LM) to check the fluency of outputs. Xu et al. (2020) show that this traditional set of metrics can be gamed. For fluency, we use the "Naturalness" metric (Mir et al., 2019) instead of PPL as it is shown to correlate better with human judgement. Apart from using BLEU for content preservation, we also report METEOR (Banerjee and Lavie, 2005), ROUGE-L (Lin, 2004) and CIDEr (Vedantam et al., 2015) scores in Table 12 of the appendix.

## 8.2 Baselines Selection

As baselines, we choose DirR (Liu et al., 2021), Stable (Lee, 2020), Transforming (Sudhakar et al., 2019), Tag (Madaan et al., 2020), CrossAligned (Shen et al., 2017), CycleRL (Xu et al., 2018), StyleEmbedding (Fu et al., 2018), D&R (Li et al., 2018) and CycleMulti (Dai et al., 2019). For the hyperparameters of each baseline, we consider the optimal parameters of the best models for each dataset reported in each respective work. Whenever available, we directly make use of the flagship style transfer outputs published as part of the original work of each reference paper to ensure that a fair comparison is done.

**SOTA baselines.** For all the datasets we considered, we found that the Direct (Liu et al., 2021), Lewis (Reid and Zhong, 2021) and Tag (Madaan et al., 2020) reported the strongest results out of all contemporary work in style transfer and our in experiments. We, therefore, consider these as current state-of-the-art baselines for the style transfer task to compare against.

## 8.3 Hyperparameter selection

The self-attention encoder block consists has 2 layers, 8 heads and embeddings of size 512. We train it for 15 epochs on the reconstruction Task and fine-tune it using signals from a pre-trained classifier

| Attribution Model | Yelp | | IMDb | | Amazon | | SNLI | |
|---|---|---|---|---|---|---|---|---|
| | Acc.% | s-BLEU | Acc.% | s-BLEU | Acc.% | s-BLEU | Acc.% | s-BLEU |
| Vanilla Attention (VA) | 73.8 | 62.41 | **69.8** | 62.4 | **70** | 57.54 | **50.76** | 66 |
| Explainable Attention (EA) | **71.3** | **64.32** | 75.25 | **70** | 77.36 | **73.21** | 66.5 | **85.14** |
| Vanilla Gradients | 74.2 | 38.8 | 81.5 | 54.47 | 74.64 | 44.19 | 61.36 | 39 |
| Gradients * X | 97.2 | 37 | 93 | 50.35 | 84.92 | 40.37 | 70.14 | 39 |
| Integrated Gradients | 77.7 | 37.29 | 81.75 | 42.42 | 71 | 40.77 | 74.73 | 43 |
| No Masking | 100 | 100 | 100 | 100 | 100 | 100 | 100 | 100 |

Table 2: Comparison of quality of style-masking produced using various attribution models. We found that $\lambda_\epsilon = 0.0$ worked best with all gradient-based methods. For attention based methods (VA and EA), we found that $\lambda_\epsilon = 0.15, 0.5$ worked best for {Yelp, IMDb, Amazon}, SNLI respectively.

for 1 epoch. During fine-tuning, $\lambda_{sta}$ is set to 1 and gradient-clipping with a threshold of $10^{-3}$ was set to prevent gradient explosion.

**Number of Parameters.** The Tag (Madaan et al., 2020) and DirR (Liu et al., 2021) models (the two best performing baselines) have 50M and 1.5B parameters respectively. The SA-MLM Encoder-only variant has 45M parameters, 30x lesser parameters than DirR's fine-tuned GPT-2 model, and roughly the same number of parameters as Tag, but outperforming it in the IMDb and SNLI datasets. We report details on training time required and infrastructure used in section A.2 of the appendix.

## 8.4 Quantitative metrics

We compute style transfer% (percentage of sentences with target style) using a Bi-LSTM based pretrained classifier trained on each dataset (refer section A.4 of the appendix for classifier details). r-BLEU and s-BLEU refer to the BLEU score taken between the output sentences and the human reference and ground truth sentences, respectively. For fluency, we measure the mean "Naturalness" score (the "Nat." column) as the mean classification score of a pretrained fluency discriminator [‖]. We also add a "Mean" score consisting of the average of style transfer%, s-BLEU and naturalness (normalised to 100) columns to denote a rough measure of the overall quality of each style transfer model.

## 8.5 Sentiment Style Transfer

Sentiment style transfer is performed on Yelp (Table 3), IMDb (Table 5) and Amazon (Table 4). We observe that for Yelp and IMDb, DirR and Encoder-only are the best-performing models according to the Mean score. In IMDb, DirR and Lewis performs better than Encoder-only in content preservation metrics but slightly lags behind in naturalness scores. In IMDb, Encoder-only achieves a

[‖]We use the pretrained naturalness classifier available in https://github.com/passeul/style-transfer-model-evaluation.

| Model | TST% | r-BLEU | s-BLEU | Nat. | Mean |
|---|---|---|---|---|---|
| DirR | 92.9 | 23.5 | 60.8 | 0.84 | **79.27** |
| Stable | 81.6 | 15.6 | 39.2 | 0.73 | 64.6 |
| Transforming | 84.8 | 18.1 | 44.7 | 0.83 | 70.9 |
| Tag | 87.7 | 16.9 | 47 | 0.83 | 72.57 |
| CrossAligned | 74.4 | 6.8 | 20.2 | 0.68 | 54.2 |
| CycleRL | 51.1 | 14.8 | 46.1 | 0.86 | 61.07 |
| StyleEmbedding | 8.59 | 16.7 | 67.6 | 0.87 | 54.4 |
| D&R | 88 | 12.6 | 36.8 | 0.89 | 71.27 |
| CycleMulti | 83.8 | 22.5 | 63 | 0.86 | 77.6 |
| Lewis | 93.1 | - | 58.5 | 0.84 | **78.53** |
| Ensemble | 56.5 | 20.5 | 63.2 | 0.85 | 68.23 |
| Generational | 63.4 | 20.3 | 61.3 | 0.83 | 69.23 |
| Encoder-only | 91.2 | 18.3 | 53.4 | 0.88 | **77.6** |

Table 3: Quantitative metrics for the Yelp Dataset.

significantly high style transfer% score at a reasonable s-BLEU of 60.9. In Amazon, Tag, Lewis and Encoder-only are the best performing. Overall we observe Encoder-only, DirR, Lewis and Tag as the best-performing models.

## 8.6 Discourse Style Transfer

It is natural to suspect that prototype editing methods are only capable of working well on "course-grained" styles, i.e. where the presence of style is determined by the presence of a fixed set of words (such as "good", "delicious" in the Yelp dataset). To inspect if this is true and gauge the ability of the SA-MLM to operate on more cognitive and complex tasks, we consider "Discourse style transfer" by using Natural Language Inference (NLI) datasets.

We report statistics for the SA-MLM in Table 6. We observe that the Encoder-only model does well overall in this task and obtains a strong mean score of 85.53 (higher than the sentiment style transfer tasks). It lags behind a little mainly in the style transfer% metric, but with a strong s-BLEU and Naturalness score. We also report qualitative examples of the discourse style transfer task in Table 9 and 10 of the appendix.

| Model | TST% | r-BLEU | s-BLEU | Nat. | Mean |
|---|---|---|---|---|---|
| DirR | 58.2 | 30.1 | 60.6 | 0.91 | 69.93 |
| Stable | 57.2 | 24.9 | 50 | 0.83 | 63.4 |
| Transforming | 58.7 | 25.5 | 52.3 | 0.92 | 67.5 |
| Tag | 75 | 32.6 | 68.3 | 0.91 | **78.1** |
| CrossAligned | 73.9 | 1.5 | 2.5 | 0.62 | 46.1 |
| StyleEmbedding | 41.1 | 13.4 | 31.2 | 0.92 | 54.8 |
| D&R | 52 | 27.2 | 56.6 | 0.92 | 66.9 |
| Lewis | 74.3 | - | 65.6 | 0.93 | 78.53 |
| Ensemble | 52.4 | 31 | 71 | 0.91 | **71.47** |
| Generational | 53.4 | 31 | 69.6 | 0.88 | 70.33 |
| Encoder-only | 63.9 | 29.6 | 69.8 | 0.92 | **75.03** |

Table 4: Quantitative metrics for the Amazon dataset.

| Model | TST% | s-BLEU | Nat. | Mean |
|---|---|---|---|---|
| DirR | 80.3 | 67.9 | 0.92 | **80.05** |
| cycleMulti | 67.2 | 73.7 | 0.93 | **77.95** |
| Ensemble | 66.8 | 69.2 | 0.92 | 76 |
| Generational | 68.9 | 65.6 | 0.93 | 75.83 |
| Encoder-only | 87.9 | 60.9 | 0.92 | **80.27** |

Table 5: Quantitative metrics for the IMDb dataset.

## 8.7 Additional Content Metrics.

Past work does not tend to clarify the meaning and prioritise the presence of "content-preservation" abilities in style transfer models Lee et al. (2021). In this effort, a more thorough analysis of content preservation abilities of DirR, Tag and Encoder-only is given in section A.5 of the appendix.

## 8.8 Qualitative Examples

Examples of the style transfer task performed using the SA-MLM for the IMDb and SNLI dataset are given in Table 8 and 9 of the appendix.

## 8.9 Human Evaluations

We only consider the relatively unexplored Discourse style transfer (Entailment to Contradiction and vice versa) task for human evaluations. We were unable to reproduce the DirR and Lewis baselines to run over the SNLI dataset. Therefore, we only compare the next strongest performing models, i.e., Tag and SA-MLM. Three volunteers were given the task of voting on 200 instances (equally split for the E to C and C to E task) from the test set. A vote consists of four options, i.e., "Model 1 better", "Model 2 better" or "Both Good", "Both Bad", where the models were randomised. To de-

| Model | TST% | s-BLEU | Nat. | Mean |
|---|---|---|---|---|
| Tag | 48.3 | 90.2 | 0.98 | 78.83 |
| Ensemble | 52.2 | 88.5 | 0.98 | 79.57 |
| Generational | 58 | 86.7 | 0.98 | 80.9 |
| Encoder-only | 76.3 | 86.3 | 0.94 | **85.53** |

Table 6: Quantitative metrics for the SNLI dataset.

| Direction | Tag better | SA-MLM better | Both Good | Both Bad | NA |
|---|---|---|---|---|---|
| E to C | 11 | **55** | 8 | 18 | 8 |
| C to E | 8 | **44** | 2 | 35 | 11 |

Table 7: Human Evaluations done to compare Tag and SA-MLM on Discourse TST task on SNLI dataset. "E" and "C" denote "Entailment" and "Contradiction" respectively.

termine the outcome for each instance, a majority from three separate votes was taken, one from each volunteer. In the case of no majority, the outcome is "No agreement". As seen in Table 7, the SA-MLM performs better in both tasks by a significant margin.

## 9 Conclusion

We introduce the SA-MLM, a modification of the standard MLM, which we show is capable of performing TST by using a style-masked input and performing a simple same-style reconstruction task with a lightweight Transformer Encoder block. On fine-tuning the SA-MLM for the TST%, it is on par with state-of-the-art models with orders of more parameters and sophisticated architectures in the Sentiment TST task. We show that complex styles such as flow of logic/ discourse can be manipulated even with using this simple style masking assumption. We empirically show that the SA-MLM performs well in this Discourses Manipulation task and outperforms another strong baseline in this task, also seen through human evaluations.

## 10 Limitations

The apparent limitation with all prototype editing models, including the SA-MLM, is that it encourages the model to only fill in necessary style words and preserve the length and structure of the original sentence. In the case of SA-MLM, the word-to-word input-output mapping while training the encoder prevents the output sentence length from changing. Though it can be argued that this even works for a relatively cognitive style like discourse, in the future, there might exist styles which explicitly require the addition/deletion of words/phrases in order to alter the style successfully. Future work will therefore focus on enabling variable-length TST outputs, similar to the (Madaan et al., 2020) approach or by incorporating a padded masked language model (Malmi, 2020).

## Acknowledgements

We would like to thank all the reviewers, whose inputs and recommendations helped to substantially improve the quality of this study. We would like to thank the human annotators for their participation.

## References

Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS ONE*, 10.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. *CoRR*, abs/1409.0473.

Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *IEEvaluation@ACL*.

Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *EMNLP*.

Chandrahas, Aditya Sharma, and Partha P. Talukdar. 2018. Towards understanding the geometry of knowledge graph embeddings. In *ACL*.

Ning Dai, Jianze Liang, Xipeng Qiu, and Xuanjing Huang. 2019. Style transformer: Unpaired text style transfer without disentangled latent representation. In *ACL*.

Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2020. Plug and play language models: A simple approach to controlled text generation. *ArXiv*, abs/1912.02164.

Zhenxin Fu, Xiaoye Tan, Nanyun Peng, Dongyan Zhao, and Rui Yan. 2018. Style transfer in text: Exploration and evaluation. In *AAAI*.

Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural Computation*, 9:1735–1780.

Zhiqiang Hu, Roy Ka-Wei Lee, and Charu C. Aggarwal. 2020. Text style transfer: A review and experiment evaluation. *ArXiv*, abs/2010.12742.

Zhiting Hu, Zichao Yang, Xiaodan Liang, Ruslan Salakhutdinov, and Eric P. Xing. 2017. Toward controlled generation of text. In *ICML*.

Sarthak Jain and Byron C. Wallace. 2019. Attention is not explanation. In *NAACL*.

Di Jin, Zhijing Jin, Zhiting Hu, Olga Vechtomova, and Rada Mihalcea. 2021. Deep learning for text style transfer: A survey. *ArXiv*, abs/2011.00416.

Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2017. Bag of tricks for efficient text classification. In *EACL*.

Dongkyu Lee, Zhiliang Tian, Lanqing Xue, and Nevin Lianwen Zhang. 2021. Enhancing content preservation in text style transfer using reverse attention and conditional layer normalization. In *ACL*.

Joosung Lee. 2020. Stable style transformer: Delete and generate approach with encoder-decoder for text style transfer. In *INLG*.

Juncen Li, Robin Jia, He He, and Percy Liang. 2018. Delete, retrieve, generate: a simple approach to sentiment and style transfer. In *NAACL*.

Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *ACL 2004*.

Yixin Liu, Graham Neubig, and John Wieting. 2021. On learning text style transfer with direct rewards. In *NAACL*.

Fuli Luo, Peng Li, Jie Zhou, Pengcheng Yang, Baobao Chang, Zhifang Sui, and Xu Sun. 2019. A dual reinforcement learning framework for unsupervised text style transfer. *ArXiv*, abs/1905.10060.

Yiwei Lyu, Paul Pu Liang, Hai Xuan Pham, Eduard H. Hovy, Barnab'as P'oczos, Ruslan Salakhutdinov, and Louis-Philippe Morency. 2021. Styleptb: A compositional benchmark for fine-grained controllable text style transfer. *ArXiv*, abs/2104.05196.

Aman Madaan, Amrith Rajagopal Setlur, Tanmay Parekh, Barnabás Póczos, Graham Neubig, Yiming Yang, Ruslan Salakhutdinov, Alan W. Black, and Shrimai Prabhumoye. 2020. Politeness transfer: A tag and generate approach. In *ACL*.

Eric Malmi. 2020. Unsupervised text style transfer with padded masked language models.

Remi Mir, Bjarke Felbo, Nick Obradovich, and Iyad Rahwan. 2019. Evaluating style transfer for text. *ArXiv*, abs/1904.02295.

Akash Kumar Mohankumar, Preksha Nema, Sharan Narasimhan, Mitesh M. Khapra, Balaji Vasan Srinivasan, and Balaraman Ravindran. 2020. Towards transparent and explainable attention models. In *ACL*.

Sharan Narasimhan, Suvodip Dey, and Maunendra Sankar Desarkar. 2022. Towards robust and semantically organised latent representations for unsupervised text style transfer. *ArXiv*, abs/2205.02309.

Preksha Nema, Mitesh M. Khapra, Anirban Laha, and Balaraman Ravindran. 2017. Diversity driven attention model for query-based abstractive summarization. In *ACL*.

Machel Reid and Victor Zhong. 2021. Lewis: Levenshtein editing for unsupervised text style transfer. In *FINDINGS*.

Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "why should i trust you?": Explaining the predictions of any classifier. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.

Ananya B. Sai, Mithun Das Gupta, Mitesh M. Khapra, and Mukundhan Srinivasan. 2019. Re-evaluating adem: A deeper look at scoring dialogue responses. *ArXiv*, abs/1902.08832.

Tianxiao Shen, Tao Lei, Regina Barzilay, and T. Jaakkola. 2017. Style transfer from non-parallel text by cross-alignment. In *NIPS*.

Tianxiao Shen, Jonas Mueller, Regina Barzilay, and T. Jaakkola. 2020. Educating text autoencoders: Latent representation guidance via denoising. In *ICML*.

Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. 2017. Learning important features through propagating activation differences. In *ICML*.

Akhilesh Sudhakar, Bhargav Upadhyay, and Arjun Maheswaran. 2019. "transforming" delete, retrieve, generate approach for controlled text style transfer. In *EMNLP*.

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. *ArXiv*, abs/1703.01365.

Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4566–4575.

Sarah Wiegreffe and Yuval Pinter. 2019. Attention is not not explanation. In *EMNLP*.

Chen Henry Wu, Xuancheng Ren, Fuli Luo, and Xu Sun. 2019a. A hierarchical reinforced sequence operation method for unsupervised text style transfer. *ArXiv*, abs/1906.01833.

Chunhua Wu, Xiaolong Chen, and Xingbiao Li. 2020. Mask transformer: Unpaired text style transfer based on masked language. *Applied Sciences*, 10:6196.

Xing Wu, Tao Zhang, Liangjun Zang, Jizhong Han, and Songlin Hu. 2019b. Mask and infill: Applying masked language model for sentiment transfer. In *IJCAI*.

Jingjing Xu, Xu Sun, Qi Zeng, Xuancheng Ren, Xiaodong Zhang, Houfeng Wang, and Wenjie Li. 2018. Unpaired sentiment-to-sentiment translation: A cycled reinforcement learning approach. In *ACL*.

Peng Xu, Jackie Chi Kit Cheung, and Yanshuai Cao. 2020. On variational learning of controllable representations for text without supervision. In *ICML*.

Xiaoyuan Yi, Zhenghao Liu, Wenhao Li, and Maosong Sun. 2020. Text style transfer via learning style instance supported latent space. In *IJCAI*.

Y. Zhang, Jingjing Xu, Pengcheng Yang, and Xu Sun. 2018. Learning sentiment memories for sentiment modification without parallel data. In *EMNLP*.

# A  Appendix

## A.1  Ethics Statement

Any TST model can be used for illicit purposes. Therefore, it is important we keep in mind a code of ethics (e.g. `https://www.acm.org/code-of-ethics`). We will make all our code open-source and will contain all details of experimentation and implementation, training time, additional hyperparameters used in the form of log files included inside the directories of our saved models, which can also be used to replicate results.

## A.2  Computational Expense and Infrastructure used

The most parameter-heavy SA-MLM model model was from the SNLI dataset. Therefore we report statistics for this model to gauge the overall computational expenses the SA-MLM demands. The model has 45 million parameters and each epoch took approximately 224 seconds to train on an Nvidia V100-SMX2 GPU and an Intel(R) Xeon(R) E5-2698 CPU. For complete details, we will make the code open source which will also contain the models we trained along with log files with all metadata about the model architecture and training.

## A.3  Qualitative examples of style transfer

Qualitative examples of style transfer using the SA-MLM for the IMDb and SNLI datasets are given in 8 and 9 respectively. We also compare some qualitative examples between the SA-MLM and the Tag (Madaan et al., 2020) baselines in Table 10.

## A.4  Details of pre-trained Classifier

We use a Bi-LSTM as our choice of classifier as it performs comparably to FastText (Joulin et al., 2017) and outperforms it in the SNLI dataset. A comparison of the two models is given in Table 11

## A.5  Additional Content Preservation Metrics

We present more content preservation metrics in Table A.5 to compare the top three performing mod-

| Direction | Negative to Positive | Positive to Negative |
|---|---|---|
| **Input** | ben affleck is back to making the same boring bad acting films . | this movie is by far one of the best urban crime dramas i 've seen . |
| **Style Masked** | ben affleck is back to making the same <mask> <mask> acting films . | this movie is by <mask> one of the <mask> urban crime <mask> i 've seen . |
| **Output** | ben affleck is back to making the same truly great acting films . | this movie is by far one of the worst urban crime garbage i 've seen . |

Table 8: Example of Sentiment style transfer on the IMDb dataset.

| Direction | Entailment to Contradiction | Contradiction to Entailment |
|---|---|---|
| **Input** | a guy in a red jacket is snowboarding in midair . a guy is outside in the snow | a woman is sitting outside at a table using a knife to cut into a sandwich . a woman is sitting inside |
| **Style Masked** | a guy in a red jacket is snowboarding in midair . a guy is <mask> in the <mask> | a woman is sitting outside at a table using a knife to cut into a sandwich . a woman <mask> <mask> <mask> |
| **Output** | a guy in a red jacket is snowboarding in midair . a guy is swimming in the park | a woman is sitting outside at a table using a knife to cut into a sandwich . a woman is a outside |

Table 9: Example of Discourse style transfer on the SNLI dataset.

| Direction | Entailment to Contradiction | Contradiction to Entailment |
|---|---|---|
| **Input** | **a black women holding a sign that says free hugs in the city . a woman is holding a sign** | **a tan dog chases a black and white soccer ball . a dog is chasing after a cat** |
| Output (Tag) | a black women holding a sign that says free hugs in the city . a woman is holding a sign | a tan dog chases a black and white soccer ball . a dog is chasing after a sport |
| Output (Encoder-only) | a black women holding a sign that says free hugs in the city . a woman is holding a book | a tan dog chases a black and white soccer ball . a dog is outside after a ball |
| **Input** | **a man is doing a task by a body of water on a farm . the man is doing something by the water** | **a dad with his child and an apple pie . a dad and his daughter with an blueberry pie** |
| Output (Tag) | a man is doing a nap by a body of water on a farm . the man is doing pushups by the water | a dad with his child and an apple outside . a dad and his daughter with an acousticelling outside |
| Output (Encoder-only) | a man is doing a task by a body of water on a farm . the man is doing nothing by the beach | a dad with his child and an apple pie . a dad and his daughter with an apple outside |

Table 10: Examples of Discourse style transfer on SNLI of SA-MLM vs Tag (Madaan et al., 2020)

| Dataset | Model | Acc.% |
|---------|-------|-------|
| Yelp | FastText | 97.6 |
|  | Bi-LSTM | 97 |
| IMDb | FastText | 99.35 |
|  | Bi-LSTM | 99 |
| Amazon | FastText | 92.1 |
|  | Bi-LSTM | 93 |
| SNLI | FastText | 72.5 |
|  | Bi-LSTM | 84 |

Table 11: Comparison of FastText, Bi-LSTM models for classification task on all datasets.

els i.e., SA-MLM, Tag (Madaan et al., 2020) and DirR (Liu et al., 2021).

### A.6 The Diversity-LSTM and Explainable Attention

Effective style-masking requires an attribution model with a high degree of plausibility, which motivates our use of "explainable" attention scores Mohankumar et al. (2020) as choice for the style-masking step.

**Why not use standard attention?** Vanilla attention scores do not serve as accurate attribution scores. Attention scores over RNN hidden states for the classification task do not correlate well with other standard interpretation metrics (Jain and Wallace, 2019), such as gradient and occlusion based methods. Feeding alternative adversarial/random attention distributions lead to only a modest effects are the model's decision (Wiegreffe and Pinter, 2019). However Wiegreffe and Pinter (2019) shows that these adversarial distributions, if properly produced, do induce poorer performance showing that vanilla attention is still partially faithful to its explanation. Mohankumar et al. (2020) postulate that attention scores over hidden states (H) are not explainable due to information mixing and subsequent entanglement/coupling and mutual information among H in RNNs. To mitigate this entanglement, diversity driven learning (inspired by results in Nema et al. (2017)) is enforced among H. This promotes the attention mechanism over such diversity-enforced H to satisfy "faithfulness" and "plausibility" properties when interpreted as attribution scores, which we refer to as "Explainable attention" (EA). Mohankumar et al. (2020) empirically show that EA does not suffer any loss in performance in the downstream task. Supporting plausibility, a) EA scores correlate better with strong attribution tools such as Integrated Gradients b) On analysis over POS tags, EA attends more to tags which are contextually important w.r.t

the given task and c) Correlates better to human judgement than vanilla attention.

**The Diversity Driven LSTM.** The Diversity LSTM consists of an LSTM-based classifier with attention (Bahdanau et al., 2015) over the H. The final context vector is fed through a feedforward layer to generate the output.

$$\tilde{\alpha}_t = \mathbf{v}^T \tanh(\mathbf{Wh} + \mathbf{b}) \ \ \forall t \in [1, m]$$
$$\alpha_t = \text{softmax}(\tilde{\alpha}_t)$$
$$\mathbf{c}_\alpha = \sum_{t=1}^{m} \alpha_t \mathbf{h}_t$$

To enforce the H of the LSTM to be "diverse" i.e. more disentangled w.r.t each other, the conicity (Chandrahas et al. (2018), Sai et al. (2019)) metric is used as an auxillary loss and is defined as the mean of "Alignment to Mean" (ATM) for all vectors $\mathbf{v}_i \in \mathbf{V}$:

$$\text{ATM}(\mathbf{v}_i, \mathbf{V}) = \text{cosine}(\mathbf{v}_i, \frac{1}{m} \sum_{j=1}^{m} \mathbf{v}_j)$$
$$\text{conicity}(\mathbf{V}) = \frac{1}{m} \sum_{i=1}^{m} \text{ATM}(\mathbf{v}_i, \mathbf{V})$$

The attention mechanism over a Diversity LSTM's H is now encouraged to be faithful to a particular set of scores, thus promoting the model to move towards more faithful and plausible attributions. The final loss is given as:

$$L(\theta_{Div}) = -\log p_{Div}(y|P) + \lambda_c on \ \text{conicity}(\mathbf{H}^P)$$

EA requires only training an additional diversity driven RNN classifier over the given dataset. After which, a single forward pass is required to obtain attribution scores. This is unlike other methods such as IG, Lime, DeepLift, Occlusion, wherein each generating each explanation requires comparatively more operations.

**Diversity-LSTM Hyperparameter Selection** Table A.6 gives details of the Diversity-LSTM classifier used during the style-masking step. We also note that the Diversity-LSTM's performance is comparable to a Bi-LSTM and FastText classifier (as shown in A.4).

| Dataset | Model | METEOR | ROUGE-L | CIDEr | Embd. Avg. Cosine Sim. | Vector Extrema Cosine Sim. | Greedy Matching Score |
|---------|-------|--------|---------|-------|----------------------|---------------------------|----------------------|
| Yelp | Encoder-only | 0.376 | 0.739 | 4.934 | 0.939 | 0.767 | 0.867 |
| | DirR | 0.444 | 0.83 | 5.813 | 0.969 | 0.867 | 0.926 |
| | Tag | 0.362 | 0.707 | 4.326 | 0.934 | 0.765 | 0.867 |
| IMDb | Encoder-only | 0.414 | 0.8 | 5.657 | 0.96 | 0.755 | 0.891 |
| | DirR | 0.472 | 0.852 | 6.344 | 0.978 | 0.847 | 0.933 |
| Amazon | Encoder-only | 0.464 | 0.868 | 6.725 | 0.964 | 0.782 | 0.921 |
| | DirR | 0.469 | 0.823 | 6.612 | 0.967 | 0.821 | 0.929 |
| | Tag | 0.453 | 0.835 | 6.548 | 0.966 | 0.781 | 0.917 |
| SNLI | Encoder-only | 0.563 | 0.906 | 8.297 | 0.986 | 0.886 | 0.96 |
| | Tag | 0.606 | 0.944 | 8.619 | 0.992 | 0.921 | 0.972 |

Table 12: Content Preservation metrics for all datasets comparing top performing models

| Dataset | Acc.% | $Loss_{con}$ | $\lambda_{con}$ |
|---------|-------|--------------|-----------------|
| Yelp | 96 | 0.06 | 10 |
| IMDb | 100 | 0.09 | 10 |
| Amazon | 89 | 0.03 | 20 |
| SNLI | 82 | 0.18 | 10 |

Table 13: Classification task statistics and choice of $\lambda_{con}$ for each dataset.

| Dataset | Style Label | Train | Dev | Test |
|---------|-------------|-------|-----|------|
| Yelp | Positive | 266K | 2000 | 500 |
| | Negative | 177K | 2000 | 500 |
| IMBb | Positive | 178K | 2000 | 1000 |
| | Negative | 187K | 2000 | 1000 |
| Amazon | Positive | 277K | 985 | 500 |
| | Negative | 179K | 1015 | 500 |
| SNLI | Entailment | 183K | 3329 | 3368 |
| | Contradiction | 183K | 3278 | 3237 |

Table 14: Split and label wise statistics of each dataset.

# Affective Natural Language Generation of Event Descriptions through Fine-grained Appraisal Conditions

**Yarik Menchaca Resendiz** and **Roman Klinger**
Institut für Maschinelle Sprachverarbeitung, University of Stuttgart
{yarik.menchaca-resendiz,roman.klinger}@ims.uni-stuttgart.de

## Abstract

Models for affective text generation have shown a remarkable progress, but they commonly rely only on basic emotion theories or valance/arousal values as conditions. This is appropriate when the goal is to create explicit emotion statements ("The kid is happy."). Emotions are, however, commonly communicated implicitly. For instance, the emotional interpretation of an event ("Their dog died.") does often not require an explicit emotion statement. In psychology, appraisal theories explain the link between a cognitive evaluation of an event and the potentially developed emotion. They put the assessment of the situation on the spot, for instance regarding the own control or the responsibility for what happens. We hypothesize and subsequently show that including appraisal variables as conditions in a generation framework comes with two advantages. (1) The generation model is informed in greater detail about what makes a specific emotion and what properties it has. This leads to text generation that better fulfills the condition. (2) The variables of appraisal allow a user to perform a more fine-grained control of the generated text, by stating properties of a situation instead of only providing the emotion category. Our Bart and T5-based experiments with 7 emotions (Anger, Disgust, Fear, Guilt, Joy, Sadness, Shame), and 7 appraisals (Attention, Responsibility, Control, Circumstance, Pleasantness, Effort, Certainty) show that (1) adding appraisals during training improves the accurateness of the generated texts by $10\,\mathrm{pp}$ in $F_1$. Further, (2) the texts with appraisal variables are longer and contain more details. This exemplifies the greater control for users.

## 1 Introduction

The main task of conditional natural language generation (CNLG) is to provide freedom to control the output text. It is commonly addressed as the intersection of text-to-text (Radford et al., 2019;

Condition: Joy Responsibility
Output: I won the tournament due to extensive training.

Figure 1: Conditioning text generation on emotions (blue) and appraisals (green) results in an improved fulfillment of the emotion condition by incorporating event descriptions (green) in the output text. This enables more fine-grained control over the generated text.

Lewis et al., 2020; Raffel et al., 2020) and data-to-text generation (Kondadadi et al., 2013; Lebret et al., 2016; Castro Ferreira et al., 2017). Therefore, models typically use two inputs: a textual trigger-phrase, and a condition to guide the generation.

In affective CNLG models, the condition is an affective state, typically represented as valence/arousal values (Maqsud, 2015) or discrete emotion names (Ghosh et al., 2017; Song et al., 2019). Arguably, the use of theories of basic emotions (Ekman and Davidson, 1994; Plutchik and Kellerman, 2013) is appropriate when the main requirement is to express a particular emotion. However, a natural communication of emotions also includes implicit expressions, where the main content of a message is not (only) the emotion. As an example, humans describe an event and leave it to the dialogue partner to infer the affective meaning ("Yesterday, my dog died"). In fact, Casel et al. (2021) report that event descriptions are used to convey an emotion in 75 % of instances in the TEC corpus (Mohammad, 2012): The sentence "I won money in the lottery" does, for most people, not require a mention of the associated emotion.

In this paper, we focus on the task of generating such emotionally connotated event descriptions (Figure 1). This poses the challenge how to represent the link between "factual" events and their emotion. Appraisal theories from psychology attempt to explain that connection with variables that represent the cognitive evaluation by a person in context of a situation (Ellsworth and Smith, 1988; Scherer et al., 2001). Does the person feel *respon-*

*sible*? Do they pay *attention* to what is going on? Is the event *pleasant*? Does somebody have *control* over what is happening? How much *effort* is needed to deal with the outcome of the situation? These variables explain emotions: Feeling *responsible* is a prerequisite for feeling *guilty*, not knowing about the outcome of a potentially negative event might cause *fear* (while knowing about it is more likely to cause *sadness*).

Our paper has two main contributions: (1) We hypothesize and show that providing appraisal information along the emotion category to the model, leads to a better fulfillment of the emotion condition. (2) We show that adding appraisal variables leads to a more fine-grained control of the generation process and the resulting texts show more details regarding the described event.[1]

## 2 Related Work

### 2.1 Emotion and Appraisal Theories

Emotions, a state of belief (Green, 1992) that results in psychological and physical changes, reflect individual's thoughts and conduct. Ekman (1992) claims the existance of six basic emotions (Anger, Disgust, Fear, Happiness, Sadness, and Surprise) that occur in response to some stimulus. Plutchik (2001) conceptualized eight primary emotions that serve as the foundation for others. While these theories do mention events as a major element in the process of developing an emotion, they do not explicitly explain the link between stimulus events and the emotion category.

Appraisal theories aim at explaining the underlying cognitive process of event evaluations. They link emotions via interpretations, evaluations, and explanations of events. Smith and Ellsworth (1985) show that 6 appraisal dimensions are sufficient to discriminate between 15 emotion categories—indeed, they constitute the emotion. Scherer et al. (2001) describes a sequence of appraisals in which events are evaluated.

Appraisal theories have only recently received interest in computational linguistics, firstly by developing analysis methods motivated to analyze events and their structure (Balahur et al., 2011). Hofmann et al. (2020) were the first who explicitly modeled appraisal variables in an existing corpus of event descriptions (Troiano et al., 2019). They used the variables from Smith and Ellsworth (1985), namely

| Conf. | Input Prompt and Output |
|---|---|
| E | generate *joy*: <u>*Last day I*</u> **was very relaxed.** |
| EA | generate *joy attention NoRESP control NoCIRC NoPLEA effort NoCERT*: <u>*Last day I*</u> **was very relaxed because I worked for 6 hours** |
| A | generate *attention NoRESP control NoCIRC NoPLEA effort NoCERT*: <u>*Last day I*</u> **decided to work for 6 hours** |

Table 1: Examples for training data. The input prompt is underlined, conditions and trigger-phrase are in *italic text*, and the output is printed in **bold**.

Attention, Certainty, Circumstance, Control, Effort, Pleasantness, and Responsibility. Troiano et al. (2023) created a larger corpus and showed that appraisals can be reliably recovered by external readers, and that they help for emotion classification. We use their corpus crowd-enVENT[2] of 6600 event descriptions, but limit their (partially correlating) 21 appraisal concepts to those that overlap with the definitions by Smith and Ellsworth (1985), which were defined via principle component analysis.

### 2.2 Affective Natural Language Generation

Most state-of-the-art systems for natural language generation follow a sequence-to-sequence approach (Sutskever et al., 2014; Cho et al., 2014). Such models take as input a sequence of words and generate as output a sequence of words. Chatbots, for instance, consider a question or an utterance from the user as input and output an answer or reaction. The architecture has two main modules, an encoder, which generates an abstract semantic representation of the input text, and a decoder, which takes the encoder representation and generates output words (Sutskever et al., 2014; Radford et al., 2019; Raffel et al., 2020; Lewis et al., 2020).

Transformer-based approaches commonly outperform recurrent neural networks (Raffel et al., 2020). We use two such methods in our paper, namely *Bart* (Lewis et al., 2020), which can be seen as a generalization of GPT (Radford et al., 2018; Brown et al., 2020; Radford et al., 2019) for its left-to-right decoder and BERT (Devlin et al., 2019) due to the bidirectional encoder. The training objective is to reconstruct the original text using a corrupted input. Further, we use *T5*, an encoder–decoder model with the philosophy to reframe NLP problems as text-to-text tasks (Raffel et al., 2020).

Most conditional language generation work has

---

[1]Training scripts and generated data are available at https://www.ims.uni-stuttgart.de/data/emotioncnlg.

[2]https://www.ims.uni-stuttgart.de/data/appraisalemotion

focused on sentiment polarity (Zhang et al., 2019; Maqsud, 2015; Niu and Bansal, 2018) and topical text generation (Orbach and Goldberg, 2020; Chan et al., 2021). The small number of papers that tackle emotion conditions include Affect-LM (Ghosh et al., 2017), a language model for generating conversational text, conditioned on five categories (Anger, Sadness, Anxiety, Positive, and Negative sentiment). Affect-LM enables customization of emotional content and intensity in the generated sentences. The customization is achieved by concatenating a condition vector to the embedding representation of the sentence. EmoDS (Song et al., 2019) is a dialogue system that can generate responses expressing the desired emotion explicitly or implicitly. The implicit generation is guided by a sequence-level emotion classifier, which recognizes a response not containing any emotion word. Within the dialog domain, the Emotional Chatting Machine involves three modules to generate responses (Zhou et al., 2018). These modules are a high-level abstraction of emotion expressions, a change in implicit internal emotion states, and an external emotion vocabulary. The Multi-turn Emotional Conversation Model (MECM, Cui et al., 2022) introduces modules to track the emotion throughout the conversation. Colombo et al. (2019) presents a GPT-2-based model (Radford et al., 2019). They use classifiers together with emotion and topic lexicons to guide the output. We use this model as a strong baseline.

None of the previous works focused on generating emotionally connotated event descriptions, which are a natural way to tell someone about the own emotional experience. None of them used psychological theories other than affect and basic emotions. We fill these gaps by combining the recent methods with appraisal theories.

## 3 Methods

The objective of our paper is to understand if adding appraisal information in addition to emotion conditions to a generator (1) improves the accuracy of the output, i.e., the likelihood that the output in fact exhibits the target emotion. Further, (2), we aim at understanding if these appraisal variables provide a more fine-grained control to the users (e.g., "I am relaxed" vs. "I am relaxed because I worked for only 6 hours"). To address these goals, we configure three CNLG models (Table 1), all based on *Bart* (Lewis et al., 2020) and *T5* (Raf-

fel et al., 2020): (a) *Condition on emotions* (E), where the model only gets informed by the target emotion (Anger, Disgust, Fear, Guilt, Joy, Sadness, or Shame) to be expressed in the generated text. (b) *Condition on emotions and appraisals* (EA), which has both the emotions and appraisals as conditional variables. The comparison between E and EA will allow us to understand the impact of the appraisals. In addition (c), we *condition on appraisals* only (A), where each generated sentence can be conditioned on one or multiple appraisals (Attention, Responsibility, Control, Circumstance, Pleasantness, Effort, or/and Certainty).

**Training.** In each configuration, we embed the conditions in the input prompt, to fine-tune the models. This strategy avoids expensive training—encoders or decoders, or both—with condition information from scratch. We create training data out of existing corpora that are labeled for emotions and appraisals consisting of input prompts and output pairs. The input prompt contains the conditions (e.g., *joy*; *joy attention*), as special tokens, followed by the trigger-phrase (e.g., *Last day I*). The output are the remaining words the model should learn to produce (e.g., *was relaxed because I worked for 6 hours*). This leads to the following three prompt representations (see Table 1 for examples):

**E:** (condition on emotions only)
    "generate [*emotion*]: [*trigger-phrase*]"
**EA:** (condition on both emotions and appraisals)
    "generate [*emotion*] [*appraisals*]$^m$: [*trigger-phrase*]"
**A:** (condition on appraisals only)
    "generate [*appraisals*]$^m$: [*trigger-phrase*]"

where *emotion* ∈ {anger, shame, disgust, fear, guild, joy, sadness} and *appraisals* is a string of the form "{attention, NoATTE} {responsibility, NoRESP} {control, NoCONT} {circumstance, No-CIRC} {pleasantness, NoPLEA} {effort, NoEF-FORT} {certainty, NoCERT}". The *trigger-phrase* consists of the first $n$ words of the training text, where $n$ is randomly chosen ($1 \leq n \leq 9$).

By using non-special tokens to represent the target conditions, the models can make use of knowledge acquired in pretraining. We opt for a string representation over a numerical representation (e.g., "control" instead of "1" or "NoCONT" instead of "0"), because preliminary experiments showed that numerical representations are sometimes interpreted as a request for repetitions by *T5* ("*generate 1 1: I feel*" → "*I feel I feel*").
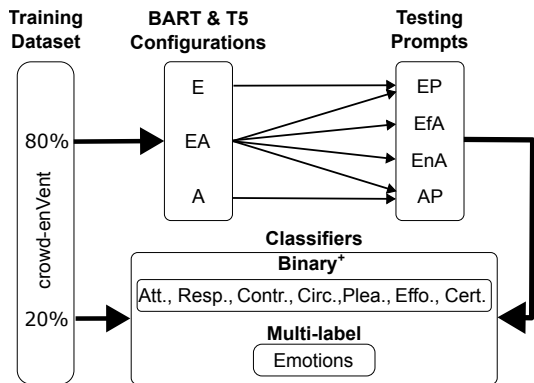
Figure 2: Experiment workflow

**Inference.** At prediction time, we obtain the five most probable sentences for each prompt. These sentences are selected using beam search (Lowerre, 1976) with beam size 30, next token temperature of 0.7, top-p[3] (nucleus) sample of 0.7. We ensure that our output excludes sentences with repeated instances of the same bigram.

## 4 Experiments

The following subsections explain the experiments conducted to test our hypotheses. In §4.1, we describe the setting and fine-tuning of the models. In §4.2, we provide results to answer the question (1) if appraisals in conjunction with emotion conditions improve the generation such that it meets the emotion condition. In §4.3 we discuss results to understand if appraisals are a means for a more fine-grained control of the generation process.

### 4.1 Experimental Settings

Figure 2 illustrates the workflow and the utilized combinations between classifiers, CNLG models, and synthetic testing prompt sets. We fine-tune according to three training set configurations (E, EA, A). This leads to six models (Bart, T5) which we evaluate with multiple testing prompt sets. The testing prompt sets only partially mirror the training regime, because the combinations of the conditional variables can be expected to be put together more freely at prediction time than as they occur in labeled data. We compare the emotion-informed models (E, EA) using the emotion testing prompt set (EP) to understand the impact of adding appraisals in the condition while not showing appraisals at prediction time. This enables us to understand if presenting appraisals improves the model's

internal representation of emotion concepts.

In addition, to understand how appraisals influence the output at inference time, we use testing prompt set with the most frequently cooccurring appraisals (EfA)—these combinations can be considered to be "compatible" with each other and the emotion (Figure 2). To challenge the models, we further use the emotion with all appraisals turned off (emotion with negative appraisals, EnAP) and test what happens when we do not provide an emotion category (appraisal-only, AP). To evaluate the performance of the models, we calculate $F_1$ with automatic emotion and appraisal classifiers (§4) and with human annotation (§5).

**Dataset.** The basis for our experiments is the crowd-enVENT data set of autobiographical reports of emotional events (see §2.1). We use a subset to train emotion and appraisal classifiers for evaluation and another subset for fine-tuning the generators (Appendix A). Each event has 21 author-assessed appraisal values, created by asking crowdworkers to complete a sentence for a given emotion (e.g., "I felt [emotion] when/that/if..."). We observed in preliminary experiments that both generation architectures (*T5* and *Bart*) have issues differentiating between the conditions and the trigger phrase, potentially due to the incompatibility of the conditions. For that reason, we focus on emotions and appraisals that have been proven to be predictable by Hofmann et al. (2020)—the variables that Smith and Ellsworth (1985) showed to be principle components for emotion categories.

We use instances that correspond to one of seven emotions (Anger, Disgust, Fear, Guilt, Joy, Sadness, and Shame) and contain an annotation with at least one of the seven appraisals[4] (Attention, Responsibility, Control, Circumstance, Pleasantness, Effort, and Certainty). This leads to 2750 instances in the corpus that we use for training. Appendix A reports details and statistics of our filtered data.

**Model Training and Data Augmentation.** We train the generation models with 80 % of the instances from this filtered corpus. The dataset is preprocessed with two goals, firstly, to create the prompts (§3) according to the desired model configuration (A, E, EA), and secondly to augment the data to prevent the models from mapping the same trigger phrase to the same output. To achieve

---

[3]Top tokens whose sum of likelihoods does not exceed a certain value (p).

[4]We discretize the [1:5] ordinal values to boolean values at a threshold of $\geq 4$, as suggested by the authors of the data set.

that, we duplicate each instance $t$ times, where $2 \leq t \leq 5$ is randomly chosen. In each duplication, a unique random number of $n$ token combinations ($1 \leq n \leq 9$) from the textual instance is used as part of the trigger phrase. Therefore, the duplication does not lead to identical instances.

**Emotion and Appraisal Classifiers.** To evaluate the performance of the generation models automatically, we use eight classifiers (one per appraisal and one for all emotions) using the remaining 20 % of the filtered crowd-enVENT dataset (15 % for training the classifier, and 5 % to evaluate the classifiers). The classifiers are built on top of RoBERTa (Liu et al., 2019) with default parameters (10 epochs, batch size 5). Each appraisal classifier predicts a boolean value whereas the emotion classifier predicts one of seven emotions. The classifiers show a performance of .75 $F_1$ Macro-Avg. for emotion classification and .56 $F_1$ for appraisal classification. These scores are, despite the limited amount of available data, comparable to previous experiments (Troiano et al., 2023). Details on these classifiers are reported in Appendix B. These classifiers allow us to perform a large set of experiments, but the non-perfect performance motivates us to confirm the main results in a human study (§5).

**Evaluation.** To evaluate the three CNLG model configurations, we create four testing prompt sets each using the thirteen most frequent starting n-grams from the crowd-enVent dataset ("I felt", "When a", "I was", "When I", "I had", "I got", "When my", "I found", "I went", "I saw", "I did", "When someone", and "I am") as trigger phrase, the seven emotions and the seven appraisals. *Emotion Prompt set* (EP) consists of 91 possible combinations between prompts and emotions (e.g., *generate joy: I felt*). The *Emotion with most frequent Appraisals Prompt set* (EfA) includes the 910 combinations between prompts, emotions and the 10 most frequent appraisals per emotion from the crowd-enVent corpus. The *Emotion with negative Appraisals Prompt sets* (EnAP) is similar to EP, but includes the appraisal vector, all set to negative values. The *Appraisal Prompt set* (AP) has the 104 possible combinations between the 13 prompts and one appraisal at a time (including the case where all appraisals are off).

It is nonsensical to compare all CNLG models on all testing prompt sets (Figure 2, interaction between Bart & T5 configurations and Test

Prompts)—e.g., the E configuration would not be able to interpret appraisal prompts (AP), similarly for the A model configuration. For every possible combination between CNLG model and the four testing prompt sets, we generate the five most probable sentences for each prompt (13,910 in total).

**State-of-the-art Baseline.** To understand how well a generic model can solve the task of affective event generation, we compare against the Affective Text Generation model (ATG, Colombo et al., 2019). ATG is conditioned on both an emotion and a topic, with the help of word lexicons. To make a fair comparison with *T5* and *Bart*, we fine-tune the language model underlying ATG, namely GPT-2, to produce emotion event descriptions using the same data that we use to train *T5* and *Bart*. The emotion and topic lexicons are unmodified because we consider them to be an essential element of ATG. Finally, for each emotion that is available in ATG and in our data (Fear, Joy, Anger, Disgust, Sadness), we generate sentences with varying intensity and target topic (Legal, Military, Politics, Monsters, Religion, Science, Space, Technology—520 in total).

## 4.2 RQ1: Do Appraisal Variables Improve Affective Text Generation?

We start the discussion of our first goal of this paper (do appraisal variables improve the model) quantitatively. Table 2 shows how well the texts from the various generation models exhibit the target emotion (evaluated against the automatic classifiers). The results should be interpreted in the context of the perplexity (Ppl.) information in Table 3.

Table 2 confirms our hypothesis for both *T5* (2nd block) and *Bart* (3rd block). The important parts are the E and EA models compared on the same *testing emotion prompt set* (EP), which only contains emotion conditions. We see that, except for Shame, the appraisal-informed model always shows a better performance—despite not showing appraisal information at inference time. Apparently, the model learns a more accurate internal emotion representation with the additional information. On average, *T5* shows a 10pp higher $F_1$ with appraisal information than without.

Obviously, an interesting question is if this performance could be further improved when providing additional appraisal information to the prompt. When using appraisal values frequently cooccurring with the emotion concept (EfA), the perfor-

| Arch. | Conf. | Testing Prmpt. | Ang. | Disg. | Fear | Guilt | Joy | Sad. | Shame | M. Avg. |
|---|---|---|---|---|---|---|---|---|---|---|
| ATG | E | — | .10 | .18 | 0.25 | — | .06 | .17 | — | .15 |
| T5 | E | EP | .28 | .50 | .63 | .23 | .60 | .32 | **.40** | .42 |
| T5 | EA | EP | .46 | **.58** | **.70** | .27 | **.77** | **.58** | .32 | **.52** |
| T5 | EA | EfA | .39 | .60 | .57 | **.35** | **.77** | .47 | .21 | .48 |
| T5 | EA | EnAP | **.52** | .55 | .64 | **.35** | .58 | .41 | .19 | .46 |
| Bart | E | EP | .36 | .45 | .40 | .29 | .63 | .43 | **.49** | .43 |
| Bart | EA | EP | **.41** | **.57** | .48 | **.41** | .63 | **.54** | .36 | **.49** |
| Bart | EA | EfA | .34 | .45 | **.52** | .29 | **.75** | .46 | .44 | .47 |
| Bart | EA | EnAP | .34 | .51 | .43 | .26 | .57 | .33 | .37 | .40 |

Table 2: Emotion $F_1$ scores of models trained with only emotions (E), emotions and appraisal conditions (EA), and only appraisal conditions (A) over the generated text using the testing prompt sets: EP (Emotions Prompt set), EnAP (Emotions with negative Appraisals Prompt set, all the appraisals are turned off) and EfA (Emotion with the most frequent Appraisals Prompt set).

| Arch. | Conf. | Testing Prmpt. | Tokens (std.) | Nouns (std.) | Verbs (std.) | Clauses (std.) | Ppl. |
|---|---|---|---|---|---|---|---|
| Hum. | Hum. | enVent | 19.3 (23) | 3.2 (3.5) | 2.8 (3.3) | .9 (1.5) | — |
| ATG | E | — | 16.4 (1.6) | 2.4 (1.3) | 2.3 (.9) | 1.7 (.6) | 22.2 |
| T5 | E | EP | 9.2 (3.4) | 2.1 (1.0) | 2.2 (1.0) | 1.2 (.6) | 26.9 |
| T5 | EA | EP | 15.1 (4.3) | 2.3 (1.1) | 2.3 (1.1) | 1.5 (.6) | 28.5 |
| T5 | EA | EfA | 13.9 (4.8) | 2.1 (1.1) | 2.1 (1.1) | 1.5 (.6) | 28.5 |
| T5 | EA | EnAP | 14.3 (4.5) | 2.2 (1.0) | 2.2 (1.1) | 1.5 (.6) | 28.5 |
| T5 | A | AP | 8.2 (3.8) | 1.8 (1.1) | 1.8 (1.0) | 1.2 (.6) | 23.5 |
| Bart | E | EP | 8.1 (4.1) | 1.7 (1.1) | 1.9 (1) | 1.4 (.5) | 69.2 |
| Bart | EA | EP | 10.5 (3.7) | 1.9 (1.0) | 1.6 (.8) | 1.2 (.4) | 51.3 |
| Bart | EA | EfA | 11.7 (4.1) | 1.9 (1.1) | 1.8 (1) | 1.3 (.5) | 51.3 |
| Bart | EA | EnAP | 13.2 (4.4) | 2.3 (1.1) | 1.9 (1) | 1.4 (.6) | 51.3 |
| Bart | A | AP | 7.7 (3.4) | 1.7 (1.2) | 1.4 (1.2) | 1.4 (.4) | 58.3 |

Table 3: Analysis of generated text using different model architectures, configurations, and prompt test sets. Mean/standard deviations are based on Spacy's tokenizer and POS. Ppl.: perplexity on test data.

mance is still higher than when not providing appraisal values during training, but apparently leaving the model more freedom in the generation with fewer conditions leads to better texts (EfA vs. E). As expected, turn off the appraisals (EnAP) leads to a drop in performance—but remains still better than the emotion-only (E) models.

Across all experiments, *T5* outperforms *Bart* and ATG. The low ATG performance could be attributed to the use of dictionaries to guide the generation process, which naturally has limited coverage and might not be suitable to describe events.

These results need to be interpreted in context with the perplexity scores shown in the last column of Table 3. Here, we see that ATG shows better performance. More importantly to answer our research question regarding the impact of appraisals is to compare the perplexity of the various E, EA, and A configurations. For the *T5* model (which shows the better emotion accuracy), there is a small decrease in language quality measured with perplexity. For the Bart model, the perplexity is in fact improving with appraisals.

### 4.3 RQ2: Do Appraisals Allow for a more Fine-grained Control?

To understand how appraisal theories can provide a more fine-grained control to the user, we conduct a quantitative and a qualitative analysis.

**Quantitative Analysis.** Table 3 shows the statistics of the generated data with the various model configurations for various prompts and as a point of reference the human and ATG-model results. Under the assumption that appraisals provide more information and more control, we would expect longer, more detailed instances with the EA models. This is indeed the case for both *T5* and *Bart*. On the emotion prompt test set (EP), instances obtained with the model trained with appraisal information (EA) are 15 tokens long for *T5*, while instances of the model trained only with emotion conditions (E) are 9 tokens long. When adding incompatible appraisal information to the prompt test data (EnAP), the text becomes even longer, with 15 tokens. The compatible appraisal values (EfA) are in between with 14 tokens. The perplexity is mainly influenced by the model architecture (GPT-2 being best, closely followed by *T5*), but it is lower for appraisal-informed models. Therefore, we can conclude that EA models generate longer instances, however, it is accompanied by the drawback of text quality, as evidenced by an increase in perplexity.

**Qualitative Analysis.** To gain a better understanding of the impact of appraisal information on the generated text, we focus on *T5*, the best-performing model to generate the target emotion (§4.3). Table 4 shows examples of texts that stem from different configurations (same trigger-phrase but different conditions). We select the most frequent appraisal and emotion combinations from the crowd-enVent dataset as conditions to generate texts.

We see that E-configuration-based generation lacks details on the event in comparison to the EA

| N. | Conf. | Condition | Trigger-Phrase | Generated Text |
|---|---|---|---|---|
| 1 | E | Shame | I felt | I felt ... when I was in my early 20's. |
| 2 | EA | Shame Att., Resp., Effo. | I felt | I felt ... when I walked out of a class meeting with no explanation |
| 3 | EA | Shame and Resp. | I felt | I felt ... when I walked out of a restaurant and smelt bad |
| 4 | E | Joy | I got | I got a job I really wanted. |
| 5 | EA | Joy, Resp., Contr., Plea. | I got | I got a job I had wanted for months leading up to my graduation. |
| 6 | EA | Joy, Resp., Contr., Plea. Effo. | I got | I got a job I applied for last week. |
| 7 | E | Fear | I had | I had to go on a rollercoaster |
| 8 | EA | Fear, Att. | I had | I had to go into a hospital with covid symptoms and give birth |
| 9 | EA | Fear, Att., Effo., Cert. | I had | I had to go into a hospital because I didn't know what was wrong with me. |
| 10 | E | Sad. | When my | When my dad passed away |
| 11 | EA | Sad., Effo., Cert. | When my | when my boyfriend told me he was leaving for college because I was too busy to do anything with him |
| 12 | EA | Sad., Plea. | When my | when my boyfriend told me he was leaving for college because I felt alone and I like I was going too |
| 13 | E | Ang. | I found | I found out my partner was cheating on me |
| 14 | EA | Ang., Plea. | I found | I found out my partner was cheating on me |
| 15 | EA | Anger, Effo., Cert. | I found | I found out my partner was taking advantage of my offer to buy him a house. |
| 16 | E | Disg. | I went | I went to the doctor and found a dead frog. |
| 17 | EA | Disg., Att., Effo. | I went | I went to a restaurant to try their dishes. They were very poor quality and did not clean up. |
| 18 | EA | Disg., Effo. | I went | I went to a restaurant to try their dishes.' |
| 19 | E | Guilt | I saw | I saw a friend being bullied at school. |
| 20 | EA | Guilt, Resp. | I saw | I saw a homeless person who needed medical attention because I couldn't afford it |
| 21 | EA | Guilt, Resp., Cert. | I saw | I saw a homeless person who had been ill and died |

Table 4: Example texts generated by T5 using different model configurations, conditions, and Trigger-Phrases.

configuration (Sentence 4 vs. 5 or 6). In Sentence 5, "I had wanted for months leading up to my graduation." the graduation aspect of the event makes one's responsibility for getting a desired job more prominent. Such properties can similarly be found in other sentence pairs in the E (e.g., 1, 4, 7, 13, 16) and EA (e.g., 2, 3, 5, 15, 17) configurations.

Appraisals that are untypical for an emotion (e.g., *Pleasantness* in *Fear* or *Sadness*) do not change the general emotion of the text (e.g., 13 and 14), but they guide the models in order to describe an event that fulfills the appraisal condition. This can be seen in a comparison of Sentences 11 and 12, where the difference is a switch of *Certainty* and *Effort* to *Pleasantness*. The model then generates "I like I was going..." to add some pleasantness despite the predominant condition being *Sadness*. Other cases show that the appraisal condition is ignored by the generator if the emotion condition is contradicting (Sentence 13 and 14). This explains why EnAP testing prompts show longer results (Table 3).

## 5 Human Evaluation

We conduct a human study to validate the automatic evaluation. Further, this study assesses additional measures, namely the quality of the generated text. We focus on the best-performing model, *T5*, fine-tuned in the EA and E configurations.

**Setup.** We randomly select 100 sentences from the following model-configuration and testing prompt set combinations: EA with EP, E with EP, and EA with EfA. In addition, we include 30 sentences from the crowd-enVent dataset to confirm the validity of the crowd-working setup. These 30 sentences are selected to be "easily-annotated" based on a high inter-annotator agreement in the original data.

We evaluate the 330 sentences on the platform https://www.soscisurvey.de. The survey consists of 23 statements to be rated on a five-level Likert scale. Seven statements correspond to the emotions ("What do you think the writer of the text felt when experiencing this event?"). Seven statements correspond to the appraisal variables ("How much do these statements apply?"), and seven questions measure the text quality (fluency, grammaticality, being written by a native speaker, semantical coherence, realistic event, written by an artificial intelligence, written by a human). In addition, we include two attention checks. We recruit participants

| Conf. | Testing Prmpt. | Ang. | Disg. | Fear | Guilt | Joy | Sad. | Shame | M. Avg. |
|---|---|---|---|---|---|---|---|---|---|
| **Hum.** | | | | | | | | | |
| Hum. | enVent | 1 | 1 | 1 | **1** | 1 | 1 | 1 | 1 |
| E | EP | .69 | .72 | .72 | **.83** | .89 | .67 | **.82** | **.76** |
| EA | EP | **.79** | **.74** | **.73** | .62 | **.92** | .82 | .6 | .74 |
| EA | EfA | .73 | .67 | .62 | .45 | .71 | .74 | .65 | .65 |
| **Auto.** | | | | | | | | | |
| Hum. | enVent | .86 | 1 | .9 | 1 | 1 | 1 | 1 | .97 |
| E | EP | .46 | .14 | .0.5 | **.44** | .78 | .33 | **.41** | .44 |
| EA | EP | **.55** | **.38** | **.82** | .31 | 1 | **.6** | .26 | **.56** |
| EA | EfA | .53 | .5 | .33 | .4 | .67 | .5 | .2 | .45 |

Table 5: Human annotation results as $F_1$ (top). For comparison, we show the automatic evaluation on the same subsample (bottom).

| Conf. | Testing Prmpt. | Fluency | Grammar | Native Spkr | Coherency | Really happen | 5−Written by AI | Written by Human |
|---|---|---|---|---|---|---|---|---|
| Hum. | enVent | 4.1 | 2.98 | 4 | 3.83 | 4.47 | 2.83 | 3.92 |
| E | EP | **3.55** | **2.43** | **3.4** | **3.36** | **4** | **2.42** | **3.25** |
| EA | EP | 3.07 | 1.88 | 2.82 | 2.89 | 3.57 | 1.86 | 2.93 |
| EA | EfA | **3.55** | **2.43** | 3.3 | 3.23 | 3.88 | 2.17 | 3.18 |

Table 6: Human evaluation of text quality using the five-level Likert scale, where 1 is *not agree at all*, and 5 is *extremely agree*. (higher is better).

via https://www.prolific.co/. §C.1 shows the questions in detail.

**Results.** To compare the performance of the conditional natural language generation models, using the human evaluation (five-level), we discretize emotion and appraisal scores, analogously to the discretization of the crowd-enVENT labels for our conditional models. We assign the labels based on a majority vote of three annotators.

Table 5 shows the performance of the generation models evaluated by the annotators on the top (Hum.). To be able to compare this to the automatic evaluation that we reported in §4.2 we show the automatic classifier-based evaluation on the same data that we used for human evaluation in addition at the bottom (Auto.). The first row, in both the human and the automatic evaluation, is the result of the evaluation on the 30 "easily-annotated" instances from the crowd-enVent data—both parts perform close-to-perfect—confirming that the general experimental setup is feasible. Further, we see that the automatic evaluation on the subset used for human evaluation mimics the results in Table 2.

The two rows for the EP testing prompt (with EA and E model configurations) also mimic the automatic evaluation. This is, however, not shown in the average $F_1$ score because the differences are less pronounced. Nevertheless, we observe that all emotions are better generated with the EA model than with the E model, except for *Guilt* and *Shame*. Therefore, the human evaluation confirms that training models with appraisal information lead to a better generation of emotion-bearing sentences. We report results for appraisals in Appendix C.2.

Table 6 shows the results for the evaluation of the quality of the generated sentences, in terms of fluency, grammar errors, coherency, text origin (text was written by a native English speaker or machine), and mimicking real event descriptions (what the text describes might happen). We have seen in Table 3 that instances generated with appraisal conditions in addition to emotion conditions lead to considerably longer texts. This seems to come with the disadvantage that the text quality is lower in all measured variables. Nevertheless, most of the values are still in an acceptable range, with the exception for grammaticality and the estimate that the text might have been written by an AI (which, however, both show comparably low values for real texts as well). As expected, the variables *Written by AI* and *Written by Human* have a strong negative correlation (Pearson's $\rho = -.77$). Importantly, the text mostly remains coherent.

# 6 Conclusion and Future Work

We presented the first study on conditional text generation based on both basic emotion category names and appraisal theories. We find that the emotion is more reliably represented when appraisals are provided during training, even when the appraisals are not provided during inference.

In addition, we provide evidence that the combination of appraisals enables a more fine-grained control over the generated text. By switching the appraisal variables, distinct event descriptions are produced, even when the emotion remains constant.

This leads to important future work: While we believe that appraisals shall be used to generate more detailed and accurate texts, the decrease in text quality needs to be controlled. In our work, we relied on prompt-based representations of the conditions in the generator models. Different model architectures (e.g., embedding the condition into the encoder, decoder, or both) could improve or maintain the quality of the generated text.

In our experiments, we relied on annotated data with labels that we used as conditions. In these data, all variables were always accessible. In a real-world setup, a deployable model would need to automatically estimate (a subset of) appraisal dimensions or request required information from a user. This might lead to a novel setup of conditioning under partial information which poses new challenges for general models of conditional text generation.

Finally, we left the topic of the event description to the choice of the model. In a real-world setup, additional conditions need to be included, for instance a topic, or a previous utterance in a dialogue. These various conditions might be in conflict in the context of a dialogue, and the model would need to rank (automatically) the conditions.

## 7 Ethical Considerations

### 7.1 Models

The proposed models are intended to link emotion theories from psychology and computational linguistics. The generated event descriptions can be used by psychologists to study the impact of appraisal and emotions in written text. There are several potential risks if the model is not used with care. It can result in biased or discriminatory language, despite that we have not observed such behaviour. Potential reasons are that a model is trained on biased data which could lead to generated texts that perpetuate stereotypes or marginalize certain groups. Particularly in the case of implicit expressions of emotions, it is important to employ models with care.

In principle, models could be used for malicious purposes, for instance to generate deceptive or harmful content (e.g., spreading misinformation or generating fake news articles). Therefore, it is crucial to employ responsible and ethical practices when utilizing natural language generation models. These risks are mainly inherent from the base pre-train language models (*Bart* and *T5*) and they are not intrinsic to our method.

### 7.2 Human evaluation

To conduct the human study in this research, we adhere to our institutional regulations and follow the recommendations by the Gemeinsame Ethikkommission der Hochschulen Bayerns[5] (GEHBa, Join Ethics Committee of the Universities in Bavaria).

As per the guidelines provided by the committee, studies that do not pose any specific risks or burdens to participants beyond what they experience in their daily lives do not require formal approval. Our study falls within that category. Therefore, it did not require approval from an ethics committee.

We relied on crowd-workers to conduct the human evaluation. The annotators were recruited using `https://www.prolific.co`, and paid according to the platform rates (£9.00/hr). All participants were shown a consent form containing the information and requirements regarding the study. They had to confirm their acceptance to be able to participate in the study. We provided an email address to contact us in case of problems during and after the study.

## 8 Limitations

Considering that our conditional approach is prompt-based, it is not surprising that it has certain limitations. First, we mentioned that both Bart and T5 have difficulties generating coherent and grammatical text, presumably because of a limited compatibility between the conditional variables (§4). Second, the conditions need to be represented as words or tokens and not numerical representation (e.g., 1 or 0), since the models cannot identify the conditions and the prompt in the fine-tuning stage. Third, the number of available datasets annotated with appraisals and emotions is very limited, since the use of appraisal theories is relatively new in the NLP community despite being a mature topic in psychology.

Even though appraisal conditions provided a better text generation for a target emotion, through event descriptions, the text quality suffers a small drop in quality (Table 6). Overall, we hope that the presented methodology and results can help guide future research and rise interest in psychological appraisal theories.

## Acknowledgements

---

[5] `https://www.gehba.de/home/`

# References

Alexandra Balahur, Jesús M. Hermida, Andrés Montoyo, and Rafael Muñoz. 2011. Emotinet: A knowledge base for emotion detection in text built on the appraisal theories. In *Natural Language Processing and Information Systems*, pages 27–39, Berlin, Heidelberg. Springer Berlin Heidelberg.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Felix Casel, Amelie Heindl, and Roman Klinger. 2021. Emotion recognition under consideration of the emotion component process model. In *Proceedings of the 17th Conference on Natural Language Processing (KONVENS 2021)*, pages 49–61, Düsseldorf, Germany. KONVENS 2021 Organizers.

Thiago Castro Ferreira, Iacer Calixto, Sander Wubben, and Emiel Krahmer. 2017. Linguistic realisation as machine translation: Comparing different MT models for AMR-to-text generation. In *Proceedings of the 10th International Conference on Natural Language Generation*, pages 1–10, Santiago de Compostela, Spain. Association for Computational Linguistics.

Alvin Chan, Yew-Soon Ong, Bill Pung, Aston Zhang, and Jie Fu. 2021. Cocon: A self-supervised approach for controlled text generation. In *International Conference on Learning Representations*.

Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, Doha, Qatar. Association for Computational Linguistics.

Pierre Colombo, Wojciech Witon, Ashutosh Modi, James Kennedy, and Mubbasir Kapadia. 2019. Affect-driven dialog generation. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3734–3743, Minneapolis, Minnesota. Association for Computational Linguistics.

Fuwei Cui, Hui Di, Lei Shen, Kazushige Ouchi, Ze Liu, and Jinan Xu. 2022. Modeling semantic and emotional relationship in multi-turn emotional conversations using multi-task learning. *Applied Intelligence*, 52(4):4663–4673.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *"Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)"*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Paul Ekman. 1992. An argument for basic emotions. *Cognition & emotion*, 6(3-4):169–200.

Paul Ekman and Richard J Davidson. 1994. *The nature of emotion: Fundamental questions.* Oxford University Press.

Phoebe C. Ellsworth and Craig A. Smith. 1988. From appraisal to emotion: Differences among unpleasant feelings. *Motivation and emotion*, 12(3):271–302.

Sayan Ghosh, Mathieu Chollet, Eugene Laksana, Louis-Philippe Morency, and Stefan Scherer. 2017. Affect-LM: A neural language model for customizable affective text generation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 634–642, Vancouver, Canada. Association for Computational Linguistics.

Otis H. Green. 1992. *The Belief-Desire Theory of Emotions*, pages 77–106. Springer Netherlands, Dordrecht.

Jan Hofmann, Enrica Troiano, Kai Sassenberg, and Roman Klinger. 2020. Appraisal theories for emotion classification in text. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 125–138, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Ravi Kondadadi, Blake Howald, and Frank Schilder. 2013. A statistical NLG framework for aggregated planning and realization. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1406–1415, Sofia, Bulgaria. Association for Computational Linguistics.

Rémi Lebret, David Grangier, and Michael Auli. 2016. Neural text generation from structured data with application to the biography domain. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1203–1213, Austin, Texas. Association for Computational Linguistics.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training

for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.

Bruce T Lowerre. 1976. *The harpy speech recognition system.* Carnegie Mellon University.

Umar Maqsud. 2015. Synthetic text generation for sentiment analysis. In *Proceedings of the 6th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 156–161, Lisboa, Portugal. Association for Computational Linguistics.

Saif Mohammad. 2012. #emotional tweets. In *\*SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pages 246–255, Montréal, Canada. Association for Computational Linguistics.

Tong Niu and Mohit Bansal. 2018. Polite dialogue generation without parallel data. *Transactions of the Association for Computational Linguistics*, 6:373–389.

Eyal Orbach and Yoav Goldberg. 2020. Facts2Story: Controlling text generation by key facts. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 2329–2345, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Robert Plutchik. 2001. The nature of emotions: Human emotions have deep evolutionary roots, a fact that may explain their complexity and provide tools for clinical practice. *American scientist*, 89(4):344–350.

Robert Plutchik and Henry Kellerman. 2013. *Theories of emotion*, volume 1. Academic Press.

Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training. *OpenAI blog*.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67.

Klaus R Scherer, Angela Schorr, and Tom Johnstone. 2001. *Appraisal processes in emotion: Theory, methods, research.* Oxford University Press.

Craig. A. Smith and Phoebe. C. Ellsworth. 1985. Patterns of cognitive appraisal in emotion. *Journal of Personality and Social Psychology*, 48(4):813–838.

Zhenqiao Song, Xiaoqing Zheng, Lu Liu, Mu Xu, and Xuanjing Huang. 2019. Generating responses with a specific emotion in dialog. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3685–3695, Florence, Italy. Association for Computational Linguistics.

Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27.

Enrica Troiano, Laura Ana Maria Oberlaender, Maximilian Wegge, and Roman Klinger. 2022. x-envent: A corpus of event descriptions with experiencer-specific emotion and appraisal annotations. In *Proceedings of the Language Resources and Evaluation Conference*, pages 1365–1375, Marseille, France. European Language Resources Association.

Enrica Troiano, Laura Oberländer, and Roman Klinger. 2023. Dimensional modeling of emotions in text with appraisal theories: Corpus creation, annotation reliability, and prediction. *Computational Linguistics*, 49(1).

Enrica Troiano, Sebastian Padó, and Roman Klinger. 2019. Crowdsourcing and validating event-focused emotion corpora for German and English. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4005–4011, Florence, Italy. Association for Computational Linguistics.

Rui Zhang, Zhenyu Wang, Kai Yin, and Zhenhua Huang. 2019. Emotional text generation based on cross-domain sentiment transfer. *IEEE Access*, 7:100081–100089.

Hao Zhou, Minlie Huang, Tianyang Zhang, Xiaoyan Zhu, and Bing Liu. 2018. Emotional chatting machine: Emotional conversation generation with internal and external memory. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence and Thirtieth Innovative Applications of Artificial Intelligence Conference and Eighth AAAI Symposium on Educational Advances in Artificial Intelligence*, AAAI'18/IAAI'18/EAAI'18, page 730–738. AAAI Press.

| Appraisal | Precision | Recall | $F_1$ |
|---|---|---|---|
| Attention | .68 | .66 | .66 |
| Certainty | .51 | .39 | .38 |
| Circumstance | .60 | .57 | .58 |
| Control | .56 | .56 | .56 |
| Effort | .54 | .53 | .52 |
| Pleasantness | .63 | .59 | .60 |
| Responsibility | .60 | .58 | .59 |
| Macro-Avg. | .59 | .55 | .56 |

Table 7: Precision, Recall and $F_1$ scores from the appraisal classifiers.

# A  Filtered Crowd-enVent Dataset

As described in §4.1, we examine seven emotions (Anger, Disgust, Fear, Guilt, Joy, Sadness, and Shame), and seven appraisals (Attention, Responsibility, Control, Circumstance, Pleasantness, Effort, and Circumstance) as conditional variables. Therefore, we filter the crowd-enVent dataset by removing records that do not have one of the seven emotions with at least one of the seven emotions. We follow the same criteria proposed by Troiano et al. (2023) to discretize the emotion and appraisal values (1 if the annotator score is larger than 3, else 0). Table 11 provides the statistical analysis of the filtered dataset. It shows the co-occurrence between emotions and appraisals, as well as details about the text, including the number of tokens, verbs, adjectives, nouns, and clauses.

# B  Automatic Classifiers

To get an impression of the reliability of the different model architectures (*Bart* and *T5*) with different conditional configurations (EA, E, A), we train one multi-label classifier for the seven emotions and 7 binary classifiers for each appraisal. The classifiers are built on top of RoBERTa (Liu et al., 2019) using the standard parameters for ten epochs with a batch size of five. Please refer to Table 7 for precision, recall, and $F_1$ scores of the appraisal classifiers, and Table 9 for the corresponding scores related to emotions.

The results for automatic classification of the appraisals are presented in Table 8. We observed that appraisal information improves the performance for emotion accuracy. This cannot be observed for the appraisal variables. For most appraisal dimensions, the model that is not conditioned on emotions works better (A is better than EA). The gap between EA and E for the same architecture is 7 pp for *T5*, and 1 pp for *Bart*.

| Arch. | Conf. | Testing Prmpt. | Att. | Resp. | Contr. | Circ. | Plea. | Effo. | Cert. | M. Avg. |
|---|---|---|---|---|---|---|---|---|---|---|
| T5 | EA | AP | **.45** | .42 | .36 | .52 | .50 | .47 | .37 | .44 |
| T5 | A | AP | **.45** | **.48** | .44 | **.71** | **.66** | .46 | **.38** | **.51** |
| Bart | EA | AP | **.45** | .43 | .42 | .53 | .47 | **.50** | .35 | .45 |
| Bart | A | AP | .35 | .43 | **.50** | .57 | .60 | .48 | .35 | .46 |

Table 8: Appraisal $F_1$ score over the generated text using the AP Prompt set, from the models conditioned on emotion and appraisals (EA), and appraisals (A).

| Emotion | Precision | Recall | $F_1$ |
|---|---|---|---|
| Anger | .72 | .58 | .64 |
| Disgust | .74 | .80 | .77 |
| Fear | .78 | .93 | .85 |
| Guilt | .56 | .71 | .62 |
| Joy | .91 | .92 | .98 |
| Sadness | .91 | .87 | .89 |
| Shame | .66 | .43 | .52 |
| Macro-Avg. | .75 | .75 | .75 |

Table 9: Precision, Recall and $F_1$ scores from the emotion classifier over the 7 classes.

# C  Human Evaluation Study Details

## C.1  Study Details

The human evaluation is performed on 330 sentences, 30 human-generated sentences from the crow-enVent dataset, and 100 sentences randomly selected from each of the following model configurations and prompt sets: EA with EP, E with EP, and EA with EfA. We use human-generated sentences to validate the study as a gold standard, under the assumption that humans are capable of accurately evaluating text written by other humans. For this purpose, we selected the top 30 *easy* sentences by ranking the filtered crowd-enVent dataset using two metrics: Emotion agreement and appraisal agreement. Table 10 shows the statistical analysis of the 330 sentences.

The survey was deployed on https://www.soscisurvey.de, and it consists of 23 questions (Table 13), divided into three sections of seven

| Conf. | Testing Prmpt. | Tokens (std.) | Nouns (std.) | Verbs (std.) | Adj. (std.) | Clauses (std.) |
|---|---|---|---|---|---|---|
| Hum. | enVent | 22.8 (16.8) | 4.4 (3.2) | 3.3 (2.4) | 1.2 (1.8) | 1.7 (.7) |
| EA | EP | 15.3 (4.0) | 2.4 (1.0) | 2.2 (1.0) | .7 (.8) | 1.5 (.6) |
| EA | EfA | 13.7 (4.7) | 1.8 (1.2) | 2.1 (1.2) | .6 (.9) | 1.4 (.6) |
| E | EP | 9.2 (3.6) | 1.6 (1.0) | 1.6 (0.8) | .5 (.7) | 1.3 (.5) |

Table 10: Statistical analysis of the automatically and human-generated text for human evaluation.

| Emo | Docs. | Att. | Resp. | Contr. | Circ. | Plea. | Effo. | Cert. | Tokens (std.) | Nouns (std.) | Verbs (std.) | Adj. (std.) | Clauses (std.) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Ang. | 450 | 305 | 55 | 86 | 72 | 15 | 309 | 184 | 21.8 (30.8) | 3.7 (4.4) | 3.2 (4.4) | 0.9 (1.8) | 1.4 (0.7) |
| Dis. | 450 | 228 | 66 | 90 | 103 | 6 | 193 | 155 | 19.4 (19.1) | 3.7 (3.4) | 2.8 (2.8) | 1.0 (1.5) | 1.4 (0.6) |
| Fear. | 450 | 378 | 119 | 100 | 157 | 17 | 345 | 148 | 19.4 (24.5) | 3.4 (3.9) | 2.8 (3.7) | 1.0 (1.4) | 1.3 (0.7) |
| Guilt. | 225 | 129 | 168 | 119 | 33 | 16 | 119 | 109 | 20.5 (22.1) | 3.2 (2.9) | 3.13 (3.4) | 1.0 (1.5) | 1.3 (0.6) |
| Joy. | 450 | 292 | 274 | 240 | 77 | 417 | 192 | 241 | 17.9 (20.7) | 3.2 (3.2) | 2.5 (2.9) | 1.1 (1.5) | 1.2 (0.6) |
| Sad. | 450 | 290 | 94 | 65 | 200 | 5 | 336 | 189 | 18.9 (22.8) | 2.9 (3.3) | 2.9 (3.4) | 1.0 (1.6) | 1.3 (0.6) |
| Shame. | 225 | 140 | 163 | 93 | 37 | 9 | 125 | 100 | 18.4 (22.4) | 2.8 (3.1) | 2.9 (3.6) | 0.8 (1.2) | 1.4 (0.7) |
| Total/Avg. | 2700 | 1762 | 939 | 793 | 679 | 485 | 1619 | 1126 | 19.5 (23.7) | 3.3 (3.7) | 2.9 (3.5) | 1.0 (1.5) | 1.4 (0.6) |

Table 11: Statistical analysis of the filtered crowd en-Vent dataset. Appraisal columns show the co-occurrence of a given appraisal and one emotion (row). Token, Nouns, Adj., and Clauses columns are the average counts for each instance.

| | Conf. | Testing Prmpt. | Att. | Resp. | Contr. | Circ. | Plea. | Effo. | Cert. | M. Avg. |
|---|---|---|---|---|---|---|---|---|---|---|
| Hum. | Hum. | enVent | .94 | .88 | .69 | .71 | .85 | .77 | .60 | .78 |
| | EA | EfA | .72 | .63 | .54 | .37 | .6 | .67 | .55 | .58 |
| Auto. | Hum. | enVent | .71 | .74 | .53 | .64 | .92 | .38 | .48 | .63 |
| | EA | EfA | .57 | .63 | .5 | .36 | .24 | .12 | .49 | .42 |

Table 12: Human annotation results as $F_1$ (1st and 2nd row) and automatic classification results (3rd and 4th row) of the human generated text (1st and 3rd row) and the automatically generated text (2nd, and 4th).

statements each, and two attention checks in a random position. The first section evaluates the emotion category of the text, the second the appraisal perception, and the last one, the quality of the text. We ask the annotator how much they agree to each statement using a five-level Likert scale (Not at all, Slightly, Somewhat, Moderately, and Extremely).

The study was conducted in August 2022, at a total cost of £250.74. Each text was annotated by three different annotators. The annotators were recruited using `https://www.prolific.co` with the following criteria:

- Age: Minimum 18 and Maximum 50.
- Nationality: UK, USA, IE.
- Place of most time spent before turning 18: United Kingdom, United States, Ireland.
- First language: English.
- Approval rate: Minimum approval rate .75.

## C.2 Appraisal Results

In the human evaluation in §5, we mainly focus on emotion evaluation. We now discuss briefly the results regarding appraisal variables.

The appraisal evaluation (Table 12) exhibits similar behavior to §4.3; the results for both automatic and human evaluation are similar (2nd and 4th row).

| Sec. | Statements |
|---|---|
| | **How much do these statements apply?** |
| Appraisal | The experiencer had to pay attention to the situation. |
| | The event was caused by the experiencer's own behavior. |
| | The experiencer was able to influence what was going on during the event. |
| | The situation was the result of outside influences over which nobody had control. |
| | The event was pleasant for the experiencer. |
| | The situation required her/him a great deal of energy. |
| | The experiencer anticipated the consequence of the event. |
| | **What do you think the writer of the text felt when experiencing this event?** |
| Emotion | Anger. |
| | Disgust. |
| | Fear. |
| | Guilt. |
| | Joy. |
| | Sadness. |
| | Shame. |
| | **How understandable is the text for you?** |
| Text quality | The text is fluent. |
| | The text has grammatical issues. |
| | The text is written by a native English speaker. |
| | The text is semantically coherent. |
| | What the text describes might have really happened. |
| | The text has been written by an artificial intelligence/machine. |
| | The text has been written by a human. |
| A.C. | Attention check. Please click "Moderately". |
| | The current question is an attention check, please select "Extremely". |

Table 13: Human evaluation survey

Therefore, it can be inferred that state-of-the-art classifiers are as good as humans, and that appraisal classification is a hard task. Even with *easy* texts (1st row) humans only achieve $78\%$ (while for emotions they achieve $100\%$). These results are aligned with Troiano et al. (2022).

# Leveraging Low-resource Parallel Data for Text Style Transfer

**Sourabrata Mukherjee** and **Ondřej Dušek**
Charles University, Faculty of Mathematics and Physics
Institute of Formal and Applied Linguistics
Prague, Czech Republic
{mukherjee,odusek}@ufal.mff.cuni.cz

## Abstract

Text style transfer (TST) involves transforming a text into a desired style while approximately preserving its content. The biggest challenge in TST in the general lack of parallel data. Many existing approaches rely on complex models using substantial non-parallel data, with mixed results. In this paper, we leverage a pretrained BART language model with minimal parallel data and incorporate low-resource methods such as hyperparameter tuning, data augmentation, and self-training, which have not been explored in TST. We further include novel style-based rewards in the training loss. Through extensive experiments in sentiment transfer, a sub-task of TST, we demonstrate that our simple yet effective approaches achieve well-balanced results, surpassing non-parallel approaches and highlighting the usefulness of parallel data even in small amounts.[1]

## 1 Introduction

Text style transfer (TST) aims to modify the style of a given text while preserving its underlying content (Shen et al., 2017; Prabhumoye et al., 2018; Li et al., 2018) (see Figure 1). The limited availability of parallel training data is a major obstacle in TST, as acquiring large-scale aligned datasets for specific style pairs is often impractical or unfeasible (Jin et al., 2022; Hu et al., 2022). The only TST study using parallel data and sequence-to-sequence learning known to us by Jhamtani et al. (2017) is a very specific application: converting modern English to Shakespeare's style, where extensive aligned paraphrases happen to exist for the purposes of literature research. Most recent TST research shifted to using non-parallel datasets and unsupervised learning (Hu et al., 2017; Zhao et al., 2018; Li et al., 2018). While it shows promising results, it does suffer a performance penalty and

---

[1]Our code and related details are available at: https://github.com/souro/low_tst.



Figure 1: An example of sentiment transfer as a TST task.

cannot avoid the data problem completely, as large quantities of non-parallel style-specific data are still hard to come by (Li et al., 2022b).

In this paper, we address the challenges of TST in low-resource scenarios by proposing methodologies that capitalize on minimal parallel data. Due to parallel data availability, we focus on sentiment transfer, a prominent sub-task within the realm of TST (Jin et al., 2022; Mukherjee et al., 2022; Luo et al., 2019a), in our experiments.[2] However, our model does not rely on a specific kind of textual styles and can be applied to TST in general.

In summary, our contributions are (1) building a TST system with low-resource parallel data, (2) applying multiple low-resource adaptation techniques, (3) and a novel style reward approach. This helps us achieve well-balanced results, surpassing previous non-parallel approaches on both automatic and human evaluation. Our experimental code is available on GitHub.[1]

## 2 Related Work

**TST with Parallel Data**   TST can be modeled as a sequence-to-sequence task and trained on pairs of texts with similar content but different styles. Here, Jhamtani et al. (2017) used a sequence-to-sequence model with a pointer network to translate modern English into Shakespearean English. However, this

---

[2]The task of sentiment transfer is related to sentence negation (Sarabi et al., 2019; Hosseini et al., 2021; Hossain and Blanco, 2022), but distinct from it, specifically aiming the scope of meaning change to sentiment only and going beyond using simple negation particles (cf. Table 3 in the Appendix).

approach to TST is inherently challenging due to the scarcity of parallel data (Hu et al., 2022).

**Non-Parallel Approaches to TST** Two main strategies were employed to avoid reliance on parallel data: (1) Straightforward text replacement, where style-specific phrases are explicitly identified and replaced (Li et al., 2018), (2) Implicit style-content disentanglement via latent representations through techniques such as backtranslation and autoencoding (Shen et al., 2017; Zhao et al., 2018; Fu et al., 2018; Prabhumoye et al., 2018; Hu et al., 2017), adversarial learning was shown to improve the results of both approaches (Lample et al., 2019; Dai et al., 2019; Li et al., 2019; Luo et al., 2019b). Despite a lot of progress, non-parallel approaches tend to produce mixed results and often require large amounts of non-parallel data, limiting their practical applicability (Li et al., 2022b).

## 3 Method

Our work sits between the parallel and non-parallel approaches, using parallel data but in very small amounts, in order to maximize performance while minimizing annotation costs. We build on transfer learning by finetuning a pretrained BART model on our task (Lewis et al., 2020). We further explore five techniques aimed at this low-resource scenario:

**Hyperparameter tuning:** As the effectiveness of Transformer models on low-resource data highly depends on hyperparameters (Araabi and Monz, 2020), we adapt our model, focusing on dropout regularization (Sennrich and Zhang, 2019) and label smoothing (Müller et al., 2019).

**Prompt-guided generation:** To align the style transfer finetuning with pre-training, we adopt using textual prompts, following Li and Liang (2021) and Li et al. (2022a). By adding prompts like *"POS:"* for positive sentences and *"NEG:"* for negative sentences, we provide explicit guidance to the decoder during fine-tuning.

**Data augmentation:** We use data augmentation by paraphrasing (see Section 4.2) to generate more training examples and improve data diversity (Shen et al., 2020; Qiu et al., 2020).

**Self-training:** To further expand our data, we use self-training, i.e., training on synthetic data generated by the model itself (He et al., 2020; Chai et al., 2022). To improve the quality of the synthetic data, we filter them using style classifier accuracy,

BLEU, and embedding similarity (cf. Section 5). We use a geometric mean of all three metrics as a sentence score, then choose a portion of the generated data with the top $k$ highest scores.

**Style reward:** To make our generator better focus on the target style accuracy, we incorporate rewards from a style classifier into the training loss. We use a simple reward $R$, which is $+1$ for instances where the generated output matches the target style, and $-1$ where it does not. We then modify the basic cross-entropy generation loss $\mathcal{L}_{\text{CE}}$ in the following way to get the overall loss $\mathcal{L}$:

$$\mathcal{L} = \alpha \cdot \text{norm}(R) + (1 - \alpha) \cdot \mathcal{L}_{\text{CE}} \qquad (1)$$

norm denotes normalization (zero mean, unit standard deviation), and $\alpha$ is a weight parameter.

## 4 Experiments

### 4.1 Dataset

We experiment on a small parallel sentiment transfer dataset of Yelp reviews by Li et al. (2018), comprising 500 positive-to-negative and 500 negative-to-positive sentences. The data was intended as an evaluation set only, but we repurpose it as a full low-resource set and split it into 400 examples for training, 100 for development, and 500 for testing. For self-training, we additionally use non-parallel sets of 2000+2000 positive and negative sentences from Li et al. (2018)'s development set.

### 4.2 Settings

We use BART-base (Lewis et al., 2020) from the HuggingFace library (Wolf et al., 2020).

**Hyperparameter tuning:** We ran three small-scale random searches for optimal values of individual parameters, resulting in the following changes from the defaults based on development set results: (1) We adjusted the learning rate (*LR*) $(5e-5 \rightarrow 1e-5)$ and *batch size* $(8 \rightarrow 3)$. (2) We increased the *Dropout* rate $(0.1 \rightarrow 0.15)$ and introduced additional attention and activation dropout (both $0.1$). (3) We introduced *L2* regularization with a value of $0.01$ and *label smoothing* with a value of $0.05$.

**Prompt-guided generation** does not have any specific settings; we only add the prompts on the input as described in Section 3.

**Data augmentation:** We used the following operations from the NLPAug library (Ma, 2019): substitute words with a *Spelling* mistake from a dictionary, *Insert* or *Substitute* words based on BERT embedding similarity, substitute words with a *Synonym* from WordNet, *Swap* or *Delete* words randomly, *Split* words into two tokens randomly. Additionally, we used *Back-translation* (Sennrich et al., 2016; Prabhumoye et al., 2018) via German using the online translation tool of Košarko et al. (2019).

We apply an augmentation to each training data example at random with a 50% probability (i.e., roughly 200 additional instances per augmentation type). We also consider an *"All"* setting where we include all augmented data.

**Self-training:** We generated parallel synthetic data of various sizes up to 2k examples. We further applied our filtering via automatic metrics (see Section 3) to choose the best 1k out of 2k examples.

**Style reward** We train a simple BERT-based (Devlin et al., 2019) sentiment classifier for this experiment, only using the same limited training set as for the main task. Its accuracy on our test set is 95.8%. We use this classifier for the style rewards, with a $\alpha = 0.5$, i.e., even split between the base cross-entropy loss and the style rewards.

### 4.3 External baselines

We compare our approaches to well-performing systems for sentiment transfer using large non-parallel datasets.[3] Our goal is to demonstrate the effectiveness of leveraging low-resource parallel data. We compare to Shen et al. (2017)'s cross-aligned autoencoder with style-specific decoders, Prabhumoye et al. (2018)'s system based on back-translation via French, and Li et al. (2018)'s text-replacement-based approach.

We also compare to state-of-the-art instruction-finetuned large language models: ChatGPT[4] and HuggingFace Chat.[5] We prompt them with a task specification and 10 randomly chosen examples from the training set. We only report results for ChatGPT, as HuggingFace Chat did not adhere to the given task, and its outputs were not parsable with our evaluation scripts.

---

[3]We faced difficulties when attempting to run some other recent approaches on our data (Xiao et al., 2021; Lee, 2020).

[4]https://openai.com, model gpt-3.5-turbo.

[5]https://huggingface.co/chat/, model OpenAssistant/oasst-sft-6-llama-30b (Köpf et al., 2023).

## 5 Evaluation & Results

We evaluate three main dimensions: style transfer accuracy, content preservation, and fluency.

We measure sentiment accuracy using Distil-BERT (Sanh et al., 2019) finetuned for sentiment analysis on the SST-2 dataset (Socher et al., 2013).[6] Following prior work (Jin et al., 2022; Hu et al., 2022), we evaluate content preservation using BLEU score (Papineni et al., 2002) and embedding similarity (Rahutomo et al., 2012) against the input sentences. We use Sentence-BERT (Reimers and Gurevych, 2019) and cosine similarity for the embedding similarity. We use GPT-2's (Radford et al., 2019) perplexity to estimate fluency.

We also run a small-scale in-house human evaluation on a random sample of 100 sentences from the test set (50 for each direction – positive-to-negative and negative-to-positive). Outputs are rated on a 5-point Likert scale for style transfer accuracy, content preservation, and fluency.

### 5.1 Automatic Metrics Results

Table 1 shows automatic metrics results. Our base BART model (experiment 01) performs decently in all metrics, but style accuracy is further improved via hyperparameter tuning (02-04), with a slight drop in BLEU score. Adding prompts (05) further increases style accuracy and makes up for the content similarity drop.

Data augmentation (06-14) leads to further improvements, especially for replacing *Synonym*s from WordNet (09), random word *Deletion* (10), and *Back-translation* (11). The best performance is achieved using *All* (14) data augmentation types (which also means a larger number of augmented examples). Augmentation generally leads to a style accuracy increase; perplexity rises, but BLEU and embedding similarity is preserved, indicative of less frequent expressions, but not much change in content.

Self-training with synthetic data (15-20) maintained the performance across the board with a slight improvement in BLEU score, but synthesizing too many examples does not lead to further improvements (18-19), likely due to an imbalance between original and synthetic data. The best results are achieved using 1k synthesized instances filtered using automatic metrics (20).

---

[6]https://huggingface.co/distilbert-base-uncased-finetuned-sst-2-english

| ID | Models | ACC | BLEU | CS | PPL |
|----|--------|-----|------|-----|-----|
| | | **Baseline** | | | |
| 01 | BART-base | 55.4 ± 2.6 | 33.8 ± 0.2 | 65.5 ± 0.9 | 127.7 ± 2.4 |
| | | **Hyperparameter tuning** | | | |
| 02 | 01 + LR & batch size | 61.7 ± 3.1 | 33.1 ± 0.2 | 67.6 ± 1.4 | 126.4 ± 1.6 |
| 03 | 02 + Dropout | 61.1 ± 2.7 | 33.3 ± 0.3 | 67.4 ± 1.3 | 126.1 ± 1.2 |
| 04 | 03 + L2 & label smoothing | 61.6 ± 3.1 | 33.2 ± 0.3 | 67.6 ± 1.4 | 126.9 ± 1.4 |
| | | **Prompt-guided generation** | | | |
| 05 | 04 + Prompt | 67.7 ± 2.6 | 33.3 ± 0.3 | 70.1 ± 1.0 | 126.7 ± 1.8 |
| | | **Data augmentation** | | | |
| 06 | 05 + Spelling | 71.1 ± 2.5 | 33.6 ± 0.4 | 70.0 ± 1.2 | 132.2 ± 2.2 |
| 07 | 05 + Insert | 71.6 ± 2.4 | 33.1 ± 0.4 | 70.8 ± 1.4 | 131.5 ± 0.9 |
| 08 | 05 + Substitute | 70.9 ± 3.5 | 33.2 ± 0.6 | 69.9 ± 1.2 | 131.9 ± 1.3 |
| 09 | 05 + Synonym | 71.5 ± 2.7 | 33.5 ± 0.5 | 71.2 ± 2.1 | 131.9 ± 0.9 |
| 10 | 05 + Delete | 72.0 ± 1.9 | 33.0 ± 0.5 | 70.7 ± 1.8 | 132.6 ± 0.8 |
| 11 | 05 + Back-translation | 72.7 ± 2.5 | 32.9 ± 0.7 | 70.6 ± 1.3 | 132.7 ± 1.6 |
| 12 | 05 + Swap | 71.1 ± 3.3 | 33.5 ± 0.1 | 70.1 ± 1.0 | 131.9 ± 1.4 |
| 13 | 05 + Split | 70.8 ± 4.5 | 33.5 ± 0.4 | 70.5 ± 1.4 | 133.5 ± 0.7 |
| 14 | 05 + All | 74.2 ± 3.2 | 33.2 ± 0.7 | 70.6 ± 2.7 | 132.5 ± 1.5 |
| | | **Self-training** | | | |
| 15 | 05 + 250 | 68.4 ± 2.5 | 33.4 ± 0.2 | 69.4 ± 1.5 | 132.5 ± 0.4 |
| 16 | 05 + 500 | 70.5 ± 5.0 | 33.6 ± 0.5 | 71.4 ± 2.3 | 132.3 ± 2.2 |
| 17 | 05 + 1k | 71.5 ± 4.8 | 34.1 ± 0.4 | 70.5 ± 2.7 | 131.0 ± 2.8 |
| 18 | 05 + 1.5k | 70.1 ± 5.0 | 34.2 ± 0.2 | 70.8 ± 2.8 | 132.4 ± 1.2 |
| 19 | 05 + 2k | 70.0 ± 4.6 | 34.3 ± 0.2 | 70.2 ± 2.2 | 132.4 ± 1.6 |
| 20 | 05 + 1k filtered | 72.6 ± 4.4 | 34.2 ± 0.4 | 71.5 ± 2.3 | 132.7 ± 1.3 |
| | | **Style reward** | | | |
| 21 | 14 + reward | 78.8 ± 2.7 | 33.1 ± 0.7 | 72.4 ± 2.4 | 132.8 ± 1.5 |
| 22 | 20 + reward | 78.4 ± 2.9 | 33.9 ± 0.7 | 72.2 ± 1.9 | 132.6 ± 1.2 |
| | | **External baselines** | | | |
| 23 | Shen et al. | 64.4 | 6.7 | 46.0 | 338.5 |
| 24 | Li et al. | 71.9 | 11.6 | 55.3 | 366.6 |
| 25 | Prabhumoye et al. | 72.4 | 3.0 | 41.7 | 318.8 |
| 26 | ChatGPT | 95.4 | 19.4 | 61.4 | 115.3 |

Table 1: Automatic evaluation results. We measure the sentiment classifier accuracy (ACC), BLEU score, Content Similarity (CS), and Fluency (PPL), see Section 5. The model names follow a format of experiment ID + Model name, indicating that the current model is built upon a base model from that particular ID. All our models' scores are averages of five runs with different random initializations, with standard deviations shown after "±".

| Models | Style | Content | Fluency |
|--------|-------|---------|---------|
| Li et al. | 2.36 | 1.57 | 1.58 |
| ChatGPT | 4.48 | 2.75 | 4.49 |
| Ours | 3.98 | 3.96 | 4.45 |

Table 2: Human evaluation of 100 randomly selected outputs on style transfer accuracy (Style), Content Preservation (Content), and Fluency (see Section 5).

Using style rewards and combining them with data augmentation (21) or self-training (22) brings further improved style accuracy, with other metrics staying approximately the same. Since both experiments 21 and 22 perform very similarly, we choose

22 as the best model for further evaluation because the self-training approach does not require additional tools, unlike the data augmentation toolkit needed for 21.

Compared to unsupervised approaches (23-25), our experiments show similar or better style accuracy while maintaining content preservation and fluency, both of which are very low for unsupervised systems. ChatGPT (26) excelled in style transfer accuracy and fluency, but also lacked in content preservation. Table 3 (see Appendix A) shows a few illustrative examples, comparing our chosen best model (22) with external baselines.

## 5.2 Human Evaluation

For the human evaluation, we compared our chosen model (experiment 22) with Li et al. (2018)'s work (24) and ChatGPT (26), chosen for their best automatic metrics results of the external models. The results in Table 2 largely confirm the automatic metrics results – the unsupervised system shows relatively poor performance, and while ChatGPT excels in hitting the target style, our approach is best on content preservation.

## 6 Conclusion

We showed that leveraging minimal parallel data in text style transfer can yield a good balance of style transfer accuracy, content preservation, and fluency. Standard low-resource approaches (hyperparameter tuning, data augmentation, self-training), improve results, while further improvement is achieved by using style classifier rewards. In future research, we plan to extend the range of low-resource techniques used and apply our approach to other style transfer tasks.

## Acknowledgments

## References

Ali Araabi and Christof Monz. 2020. Optimizing transformer for low-resource neural machine translation. In *Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020*, pages 3429–3435.

Junyi Chai, Reid Pryzant, Victor Ye Dong, Konstantin Golobokov, Chenguang Zhu, and Yi Liu. 2022. FAST: improving controllability for text generation with feedback aware self-training. *CoRR*, abs/2210.03167.

Ning Dai, Jianze Liang, Xipeng Qiu, and Xuanjing Huang. 2019. Style transformer: Unpaired text style transfer without disentangled latent representation. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Volume 1: Long Papers*, pages 5997–6007, Florence, Italy.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, MN, USA.

Zhenxin Fu, Xiaoye Tan, Nanyun Peng, Dongyan Zhao, and Rui Yan. 2018. Style transfer in text: Exploration and evaluation. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pages 663–670.

Junxian He, Jiatao Gu, Jiajun Shen, and Marc'Aurelio Ranzato. 2020. Revisiting self-training for neural sequence generation. In *8th International Conference on Learning Representations, ICLR 2020*, Addis Ababa, Ethiopia.

Md Mosharaf Hossain and Eduardo Blanco. 2022. Leveraging affirmative interpretations from negation improves natural language understanding. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022*, pages 5833–5847, Abu Dhabi.

Arian Hosseini, Siva Reddy, Dzmitry Bahdanau, R. Devon Hjelm, Alessandro Sordoni, and Aaron C. Courville. 2021. Understanding by understanding not: Modeling negation in language models. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021*, pages 1301–1312, Online.

Zhiqiang Hu, Roy Ka-Wei Lee, Charu C. Aggarwal, and Aston Zhang. 2022. Text style transfer: A review and experimental evaluation. *SIGKDD Explor.*, 24(1):14–45.

Zhiting Hu, Zichao Yang, Xiaodan Liang, Ruslan Salakhutdinov, and Eric P. Xing. 2017. Toward controlled generation of text. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 1587–1596, Sydney, NSW, Australia.

Harsh Jhamtani, Varun Gangal, Eduard Hovy, and Eric Nyberg. 2017. Shakespearizing modern language using copy-enriched sequence to sequence models. In *Proceedings of the Workshop on Stylistic Variation*, pages 10–19, Copenhagen, Denmark.

Di Jin, Zhijing Jin, Zhiting Hu, Olga Vechtomova, and Rada Mihalcea. 2022. Deep learning for text style transfer: A survey. *Comput. Linguistics*, 48(1):155–205.

Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi-Rui Tam, Keith Stevens, Abdullah Barhoum, Nguyen Minh Duc, Oliver Stanley, Richárd Nagyfi, Shahul ES, Sameer Suri, David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and Alexander Mattick. 2023. OpenAssistant conversations - democratizing large language model alignment. *CoRR*, abs/2304.07327.

Ondřej Košarko, Dušan Variš, and Martin Popel. 2019. LINDAT translation service.

Guillaume Lample, Sandeep Subramanian, Eric Michael Smith, Ludovic Denoyer, Marc'Aurelio Ranzato, and Y-Lan Boureau. 2019. Multiple-attribute text rewriting. In *7th International Conference on Learning Representations, ICLR 2019*, New Orleans, LA, USA, May 6-9, 2019.

Joosung Lee. 2020. Stable style transformer: Delete and generate approach with encoder-decoder for text style transfer. In *Proceedings of the 13th International Conference on Natural Language Generation, INLG 2020*, pages 195–204, Dublin, Ireland.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020*, pages 7871–7880, Online.

Dianqi Li, Yizhe Zhang, Zhe Gan, Yu Cheng, Chris Brockett, Bill Dolan, and Ming-Ting Sun. 2019. Domain adaptive text style transfer. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pages 3302–3311.

Juncen Li, Robin Jia, He He, and Percy Liang. 2018. Delete, retrieve, generate: a simple approach to sentiment and style transfer. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, Volume 1 (Long Papers)*, pages 1865–1874, New Orleans, Louisiana, USA.

Junyi Li, Tianyi Tang, Jian-Yun Nie, Ji-Rong Wen, and Xin Zhao. 2022a. Learning to transfer prompts for text generation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2022*, pages 3506–3518, Seattle, WA, United States.

Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers)*, pages 4582–4597, Virtual Event.

Xiangyang Li, Xiang Long, Yu Xia, and Sujian Li. 2022b. Low resource style transfer via domain adaptive meta learning. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2022*, pages 3014–3026, s, WA, United States.

Fuli Luo, Peng Li, Pengcheng Yang, Jie Zhou, Yutong Tan, Baobao Chang, Zhifang Sui, and Xu Sun. 2019a. Towards fine-grained text sentiment transfer. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Volume 1: Long Papers*, pages 2013–2022, Florence, Italy.

Fuli Luo, Peng Li, Jie Zhou, Pengcheng Yang, Baobao Chang, Xu Sun, and Zhifang Sui. 2019b. A dual reinforcement learning framework for unsupervised text style transfer. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI 2019*, pages 5116–5122, Macao.

Edward Ma. 2019. NLP augmentation. https://github.com/makcedward/nlpaug.

Sourabrata Mukherjee, Zdenek Kasner, and Ondrej Dusek. 2022. Balancing the style-content trade-off in sentiment transfer using polarity-aware denoising. In *Text, Speech, and Dialogue - 25th International Conference, TSD 2022*, volume 13502 of *Lecture Notes in Computer Science*, pages 172–186.

Rafael Müller, Simon Kornblith, and Geoffrey E. Hinton. 2019. When does label smoothing help? In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019*, pages 4696–4705, Vancouver, BC, Canada.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.

Shrimai Prabhumoye, Yulia Tsvetkov, Ruslan Salakhutdinov, and Alan W. Black. 2018. Style transfer through back-translation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Volume 1: Long Papers*, pages 866–876, Melbourne, Australia.

Siyuan Qiu, Binxia Xu, Jie Zhang, Yafang Wang, Xiaoyu Shen, Gerard de Melo, Chong Long, and Xiaolong Li. 2020. EasyAug: An automatic textual data augmentation platform for classification tasks. In *Companion of The 2020 Web Conference 2020*, pages 249–252, Taipei, Taiwan.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI Blog*.

Faisal Rahutomo, Teruaki Kitasuka, and Masayoshi Aritsugi. 2012. Semantic cosine similarity. In *The 7th international student conference on advanced science and technology ICAST*, volume 4, page 1.

Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using siamese BERT-Networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019*, pages 3980–3990, Hong Kong.

Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *CoRR*, abs/1910.01108.

Zahra Sarabi, Erin Killian, Eduardo Blanco, and Alexis Palmer. 2019. A corpus of negations and their underlying positive interpretations. In *Proceedings of the Eighth Joint Conference on Lexical and Computational Semantics, *SEM@NAACL-HLT 2019*, pages 158–167, Minneapolis, MN, USA.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, Volume 1: Long Papers*, Berlin, Germany.

Rico Sennrich and Biao Zhang. 2019. Revisiting low-resource neural machine translation: A case study. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Volume 1: Long Papers*, pages 211–221, Florence, Italy.

Dinghan Shen, Mingzhi Zheng, Yelong Shen, Yanru Qu, and Weizhu Chen. 2020. A simple but tough-to-beat data augmentation approach for natural language understanding and generation. *CoRR*, abs/2009.13818.

Tianxiao Shen, Tao Lei, Regina Barzilay, and Tommi S. Jaakkola. 2017. Style transfer from non-parallel text by cross-alignment. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017*, pages 6830–6841, Long Beach, CA, USA.

Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Y. Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing, EMNLP 2013*, pages 1631–1642, Seattle, Washington, USA.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Fei Xiao, Liang Pang, Yanyan Lan, Yan Wang, Huawei Shen, and Xueqi Cheng. 2021. Transductive learning for unsupervised text style transfer. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021*, pages 2510–2521, Virtual Event / Punta Cana, Dominican Republic.

Junbo Jake Zhao, Yoon Kim, Kelly Zhang, Alexander M. Rush, and Yann LeCun. 2018. Adversarially regularized autoencoders. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 5897–5906, Stockholm, Sweden.

# A  Appendix: Sample Outputs

| | Negative → Positive | Positive → Negative |
|---|---|---|
| **Source** | terrible menu, high prices, bad customer service . | it 's a much better option than the club scene . |
| **Gold** | nice menu , good prices , great service - for both dinner and breakfast ! | i would rather go to the club than here . |
| Shen et al. | fantastic selection of great customer ! | it 's a good experience for the whole airport i would . |
| Li et al. | no nonsense in service . | it 's a much better than the club scene . |
| Prabhumoye et al. | bad customer service with the food of this location . | she did n't go back with this place . |
| ChatGPT | marvelous entertainment, budget-friendly choices, exceptional atmosphere. | absolutely disastrous, it's worse than the late-night traffic. |
| Ours (exp. 22) | great menu, high prices, great customer service . | it's a terrible alternative to the club scene . |
| **Source** | the bad news that my vision had deteriorated made the visit even worse . | all of my clothes are returned in sparkling condition ! |
| **Gold** | the good news that my vision had improved made the visit even better . | all of my clothes are returned in terrible condition ! |
| Shen et al. | the good thing i have the whole nails made my whole gem ! | all of my car here are nothing in any room . |
| Li et al. | the problem was the red deal by handles the night my questions did n't . | all of my clothes are returned in my condition ! |
| Prabhumoye et al. | the worst time i have ever had to get a disappointment . | all of the food is not very good in all . |
| ChatGPT | the remarkable revelation of my surprise birthday party plans made the visit even more special. | The condition of all my belongings is extremely terrible! |
| Ours (exp. 22) | the good news that my vision had improved made the visit even better . | all of my clothes are returned in terrible condition ! |
| **Source** | it's located in a slum scottsdale area and isn't accomodating. | my father has decided to upgrade my mothers engagement ring this xmas . |
| **Gold** | it 's located in a great part of scottsdale and was really accommodating . | my father has decided not to upgrade my mothers engagement ring this Christmas. |
| Shen et al. | cute shop in a sunday area and desert ! | my son did to have my whole card to celebrate my appointment off . |
| Li et al. | no bueno in the north nonsense and not acknowledged a word or anything . | my father has decided to upgrade paint now . |
| Prabhumoye et al. | minutes later for the food and not worth the food . | my husband ordered me to get the worst service in the food . |
| ChatGPT | this place is family-owned, but it could greatly benefit from improving their staff. | my father has decided to downgrade my mother's engagement ring. |
| Ours (exp. 22) | it's located in a slum scottsdale area and is accomodating. | my father has decided not to upgrade my mothers engagement ring this xms. |

Table 3: Example output comparison on samples from the test set. Sentiment marker words are colored. Note that our model balances well between style transfer accuracy and content preservation, better than others.

# Reverse-Engineering Decoding Strategies Given Blackbox Access to a Language Generation System

**Daphne Ippolito***
dei@google.com

**Nicholas Carlini***
ncarlini@google.com

**Katherine Lee***
katherinelee@google.com

**Milad Nasr***
miladnasr@google.com

**Yun William Yu†**
ywyu@math.toronto.edu

## Abstract

Neural language models are increasingly deployed into APIs and websites that allow a user to pass in a prompt and receive generated text. Many of these systems do not reveal generation parameters. In this paper, we present methods to reverse-engineer the decoding method used to generate text (i.e., top-$k$ or nucleus sampling). Our ability to discover which decoding strategy was used has implications for detecting generated text. Additionally, the process of discovering the decoding strategy can reveal biases caused by selecting decoding settings which severely truncate a model's predicted distributions. We perform our attack on several families of open-source language models, as well as on production systems (e.g., ChatGPT).

## 1 Introduction

Language models are increasingly being incorporated into web applications and other user-facing tools.[1] These applications typically do not provide direct access to the underlying language model or the decoding configuration used for generation. In this paper, we show how even in this blackbox setting, it is possible to identify the decoding strategy employed for generation. We consider the case where one only has access to a system that inputs a prompt and outputs a generated response. We present algorithms for distinguishing the two most popular decoding strategies, top-$k$ and nucleus sampling (a.k.a. top-$p$), and estimating their respective hyperparameters ($k$ and $p$).

The choice of decoding strategy—the algorithm used to sample text from a language model—has a profound impact on the randomness of generated text, introducing biases toward some word choices. For example, when OpenAI's ChatGPT,[2] a chatbot built with large language models, is repeatedly passed a prompt

asking it to report the outcome of rolling a twenty-sided die, we found that it only returns 14 of the 20 options, even though all should be equally likely.

Prior work has shown that knowing the decoding method makes it easier to detect whether a writing sample was generated by a language model or else was human-written (Ippolito et al., 2020). As generated text proliferates on the web, in student homework, and elsewhere, this disambiguation is becoming increasingly important.

Concurrent work to ours by Naseh et al. (2023) has developed similar strategies for detecting decoding strategy from a blackbox API: however, they focus more on identifying hybrid decoding strategies (including beam search), whereas we focus more on prompt engineering to produce close-to-uniform token distributions that reduce the number of queries needed. Our proposed methods complement but are not comparable to those of Tay et al. (2020). Their method trains classifiers that input a generated text sequence and output a prediction for the decoding strategy used to generate it. In contrast, our method interacts with an API and does not require any data or ML training.

## 2 Background

Neural language models are not inherently generative. A causal language model $f_\theta$ takes as input a sequence of tokens $x_1,...,x_{t-1}$ and outputs a score for each possible next token $x_t$, computing the a likelihood score for each token in the vocabulary, which can be transformed into a probability distribution by applying a softmax such that $\text{Prob}(x_t|x_1,...,x_{t-1}) \sim f_\theta(x_1,...,x_{t-1})$.

A **decoding method** takes this probability distribution as input and samples a particular token to output. The simplest algorithm is *argmax decoding* (also known as 'greedy decoding'), where the most likely next token is outputted. Argmax is rarely used in practice because (1) only one generation can be produced for any given prompt, and (2) generations with argmax tend to be repetitive and low-quality.

---

*Google Deepmind,†University of Toronto

[1]E.g., see https://gpt3demo.com/ for a list of such apps.

[2]https://openai.com/blog/chatgpt/

396

Most commonly used decoding algorithms are based on random sampling: a token is chosen with probability proportional to the likelihood assigned to it by the model. Whereas argmax sampling has too little randomness, purely random sampling over the full distribution can have too much, leading to text that is too erratic and prone to errors. Thus, it is common to modify the distribution to reduce entropy before sampling from it.

In this short paper, we focus on two popular strategies researchers have developed for decoding: top-$k$ sampling (Fan et al., 2018) and top-$p$ sampling (Holtzman et al., 2019) (also known as nucleus sampling). Top-$k$ sampling involves the implementer picking a fixed hyperparemter $k$ then only ever sampling from the $k$ most likely items by assigning all other items a score of 0 before applying the softmax. Top-$p$ sampling involves the implementer picking a fixed hyperparamter $p$. Then at each step $t$ of generation, a $k_t$ is selected such that the $k_t$ most likely vocabulary items cover $p$ proportion of the total probability mass in the distribution. More precisely, let the notation $x^{(l)}$ refer to the $l$th most likely token in the distribution predicted at step $t$. We set $k_t$ to the first value for which $\sum_{l=1}^{k_t} \text{Prob}(x_t = x^l | x_1, ..., x_{t-1}) \geq p$. Then, the distribution is truncated to the $k_t$ most likely tokens, as described above for top-$k$.

Other common methods like beam search and temperature annealing are omitted in the interest of space (cf. Zarrieß et al. (2021) and Wiher et al. (2022)). Temperature annealing simply modifies the probability distributions of the output tokens, so the methods in this manuscript can be easily generalized (and indeed were in the concurrent work of Naseh et al. (2023)). Beam search is a bit more complicated, as tokens are not chosen independently of previous tokens; instead, multiple candidate token paths are retained. As such, it would be necessary to generate more than a single word for each prompt, which is the primary interrogative tool we use here.

## 3 Method

### 3.1 Threat Model

We assume black-box, query-only access to the system $\text{Gen}: m \mapsto r$ which takes as input a prompt string $m$ and outputs a textual response $r$; without loss of generality, we assume that the response $r$ is exactly one token long. The adversary can input arbitrary prompts and observe the output response. In most of our experiments, we assume Gen passes $m$ into the language model without any modification, then generates a con-

tinuation using an unknown decoding strategy. However, in some cases, such as for ChatGPT, the system might modify the provided prompt, $m$, such as by prepending few-shot examples, before passing it to the language model. Still, we assume that the causal language model can be repeatedly queried by a fixed prompt $m'$, even if modified from the original $m$.

The adversary's attack objective is to determine the decoding strategy employed by Gen, outputting either `topk` or `topp`, as well as the value for either $p$ or $k$.

### 3.2 Intuition for Method

We begin with the intuition of our attack. Suppose we were given a prompt $m$, such that the output of $\text{Gen}(m)$ is equally likely to be any item from a set of vocabulary items $\mathcal{V}_m \subseteq \mathcal{V}$. For example, the prompt "`List of capital English letters, chosen uniformly at random:`" ought to result in the model emitting each of the 26 letters of the alphabet with equal probability. However, suppose that when we repeatedly prompt the model in this way, it only ever emits 10 different letters. What could cause this?

One explanation could be that our prompt does not actually induce a uniform probability distribution over each of the 26 letters, and in fact that the model assigns (nearly) zero probability mass to the 11th most likely token. Suppose we know for a fact the prompt does induce a near-uniform distribution on all publicly-available language models: then the more likely explanation would be that the sampling algorithm itself truncated this distribution—either with top-$k$ or top-$p$ sampling. By measuring what fraction of the words we would expect to get generated actually do get generated for prompts with known output distributions, we can estimate values for $k$ and $p$ and distinguish between these two techniques.

### 3.3 Estimating $k$

Suppose, for a given prompt $m$, we call $\text{Gen}(m)$, $n$ number of times, each time keeping just the first token of the output. We can trivially lower bound $k$ by observing the number of unique items in a set of responses. As $n$ approaches $\infty$, all $k$ allowed responses will be observed. To achieve a compute-efficient attack, our goal is to estimate $k$ while keeping $n$ as small as possible. Appendix A gives theoretical accuracy/runtime estimates for this approach by posing it as the coupon collector problem (Pólya, 1930).

In practice, we use Algorithm 1 (see Appendix), which repeatedly estimates a lower bound for $k$ using two different prompts $m_1$ and $m_2$ for increasing

numbers of trials until (1) the two estimates match and (2) the $x^{(k)}$ token appears at least twice in both generations (to prevent spurious matching).

## 3.4 Estimating $p$

In this paper, we set a goal of determining $p$ to within $0.05$ of the true value. We can upper bound $p$ by constructing a prompt that yields a known, computable distribution over a set of vocab items $\mathcal{V}_m$. Then to attack a system, we repeatedly sample with the prompt, and count how many of those items are generated. Let's call this value $k$.[3] We estimate $p$ as the sum of the probabilities of the $k$ most likely tokens in the known distribution over items in $\mathcal{V}_m$. Because our guessed distributions are not perfect, instead of relying on just one prompt for our estimate, we instead average over two upper bounds of $p$ derived from two different prompts. Although our experiments here use only two prompts, increased precision is achievable by using additional prompts. The detailed algorithm can be found in the Appendix.

## 3.5 Distinguishing Top-$k$ from Top-$p$

To distinguish between top-$k$ and top-$p$, we need only reject the hypothesis that top-$k$ is used. It turns out that we can simply reuse Algorithm 1 because we already built in a measure of concordance in the $k$ predictions. If the two prompts used as input to Algorithm 1 continue to yield very different predictions of $k$ no matter how many samples are taken, we can reject the hypothesis of top-$k$ being used. For rejecting top-$k$, we found it useful to start with two prompts with radically different distributions; it suffices to choose prompts that with very differently sized $\mathcal{V}_m$, such as ADVERBS and MONTHS.

Although we did not explore it in this short paper, we could in theory also reject top-$p$ by looking at how closely the $p$ estimates from different prompts match. This may prove helpful if we wish to determine that neither top-$p$ or top-$k$ is being used, but is unnecessary for simply disambiguating the two.

**Prompt Selection** In addition to the distributional properties described above, we also need our prompts to have the property that the first space-separated word in the output of $\text{Gen}(m)$ is in-fact a word in the vocabulary. Since we often do not know which vocabulary was used by the model we are attacking, we choose prompts which yield distributions over words

---

[3]This is a slight abuse of notation since we used $k$ earlier for top-$k$, but in both cases, this value corresponds precisely to the number of unique tokens seen.

Table 1: Prompts used for top-$k$ and (top) top-$p$ (bottom) estimation on open-source models. The first two prompts include randomly selected exemplars (shown in blue). For MONTHS, Ramadan is included as the 13th month.

| Name | Prompt | $|\mathcal{V}_m|$ |
|---|---|---|
| NOUNS | List of nouns chosen completely randomly: dog, slash, altar | 8,432 |
| ADVERBS | List of adverbs chosen completely randomly: formally, blatantly, sadly | 504 |
| MONTHS | She came to visit in the month of | 13 |
| DATES | The accident occurred on March | 31 |



Figure 1: Error in top-$k$ estimation for 500 GPT-2 XL- and 500 bloom-3b-based systems. Error for both models is very low in common settings ($k < 100$).

Table 2: Performance at $k$ estimation over 100 systems with $k$ values randomly chosen between 1 to 500.

| Model | Acc | Acc$\pm$5 | Avg Error |
|---|---|---|---|
| GPT-2 Base | 28% | 76% | 1.3 |
| GPT-2 XL | 44% | 80% | 0.9 |
| BLOOM-3B | 0% | 71% | 2.3 |
| pythia-2.7b | 22% | 81% | 1.1 |

which are likely to be tokens in all models trained on webtext[4]. Table 1 shows all the prompts used in all experiments except for those we used on ChatGPT (which had to be longer), and Appendix C gives more details on prompt selection (including for ChatGPT).

## 4 Experiments

We conduct experiments on four language models where we can set the decoding strategy: GPT-2 Base and XL (Radford et al., 2019), GPT-3 Davinci (Brown et al., 2020), BLOOM 3B (Scao et al., 2022), and Pythia 2.7B (Biderman et al., 2023).

### 4.1 Predicting $k$

We used two prompt templates for estimating top-$k$: NOUNS and ADVERBS (see Table 1), each with 16 randomly selected exemplars. We build an evaluation

---

[4]Other prompts may be needed for attacking code models.

Figure 2: Performance at $p$ estimation over 100 systems with $p$ values ranging from 0.0 to 1.0, when the known distributions are computed using GPT-2 base (orange), and GPT-3 Davinci likelihoods (blue). Using Davinci as the known model leads to a better attack on GPT-3 models, but a worse one on all other models. RMSE in Table 6.
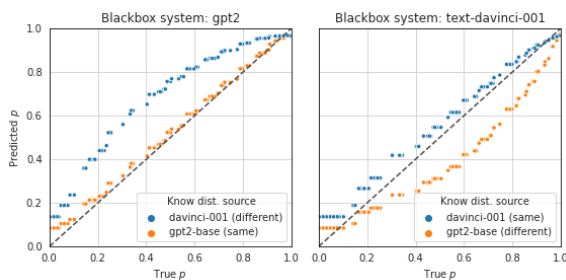


Figure 3: Mismatch between the known distributions and the distributions of the language model underlying the blackbox system lead to increased error. The matched estimates still have a slight systematic upward bias because we use the upper bounds for $p$ in our algorithm.

set of 100 systems, each with a $k$ selected uniform randomly to be between 1 and 500. Table 2 shows the accuracy of our approach on this evaluation set.[5] We see that while our method is not so great at guessing $k$ perfectly, on average its guesses are between 0.9 and 2.3 off (depending on the underlying model).

In Figure 1, we plot accuracy as a function of true $k$ for GPT-2 XL. This plot reveals that our method is especially effective at predicting $k$ for $k < 300$, and accuracy deteriorates for higher $k$. The vast majority of applications use $k$ well within this range, and it is simple to adjust for larger $k$ by increasing the max number of iterations parameter.

## 4.2 Predicting $p$

We build an evaluation set of 100 systems, each with a randomly assigned $p$ between 0 and 1. Table 1 shows the prompts used for top-$p$ estimation: MONTHS and DATES. For each prompt, we need to compute a known distribution over the next word. We experi-

---

[5]GPT-3 is omitted because the API does not expose top-$k$.



Figure 4: We plot the ratio between the $k$ estimated using ADVERBS and using MONTHS, for systems actually using top-$p$. Except when $p$ is near its extreme values, the prompt with larger $\mathcal{V}_m$ results in more generated words.

ment with using both GPT-2 Base and GPT-3 Davinci for this. For GPT-2, we compute the distribution directly; for GPT-3, we estimate it by running 1,000 trials with full random sampling. Figure 2 shows our method's accuracy at predicting within 0.05 of the true $p$ value. Figure 3 shows two limitations: (1) our estimates are worse when there is significant mismatch between Gen's distributions and our known distributions; and (2) the minimum $p$ our method can predict is $\left(\mathrm{Prob}(x_1^{(1)}) + \mathrm{Prob}(x_2^{(1)})\right)/2$, reducing accuracy for low $p$ values. Further research is needed into the design of prompts which induce consistent distributions over many families of language models.

## 4.3 Distinguishing Top-$k$ and Top-$p$

To evaluate our ability to distinguish between top-$k$ and top-$p$, we conduct the following experiment. We take 10 systems with $p$ values ranging from 0.0 to 1.0 and find the chance we misclassify each system as using top-$k$. Figure 4 reports the results of this experiment. We see that it would be fairly easy to distinguish the two methods by thresholding the ratio of the $k$ values returned by the two prompts. Note that at the extreme values of $k$ and $p$, the method are indistinguishable. Top-$k$ with $k{=}1$ and top-$p$ with $p{=}0.0$ are both the same as argmax; top-$k$ wirh $k{=}|\mathcal{V}|$ and top-$p$ with $p{=}1.0$ are both the same as full random sampling.

## 4.4 ChatGPT

We cannot repeat the exact same experiments with ChatGPT because (1) it does not use the exact prompt passed to the UI as the language model input, instead preprocessing it into a conversation format; and (2) the rate limiting prevents us from easily conducting many trials. We instead employ four conversational-format prompts (see Table A5). For the known distribution, we try out empirical

| Prompt | $n$ | Est. $p$ |
|---|---|---|
| MONTHSCHAT | 200 | 0.84 |
| DATESCHAT | 125 | 0.74 |
| D20CHAT | 115 | 0.79 |
| D100CHAT | 500 | 0.86 |

Table 3: The values of $p$ estimated for ChatGPT using different prompts, where $n$ is the number of samples taken.

distributions from five different versions of GPT-3, and take the one with lowest total variation distance from ChatGPT's output distribution. Table 3 shows the $p$ estimated using each prompt. Averaging these, we get $p$=0.81. This estimate could be further narrowed down by incorporating more prompts, though of course we cannot validate this number due to opacity of the ChatGPT proprietary system.

### 4.5 Room for Improvement

All of the estimates reported in this section could be improved with additional queries to the model For both $p$ and $k$ estimation, we average over the estimates from just two prompts, but using more prompts would lead to better estimates. In addition, to improve top-$k$ estimation for larger $k$, one can increase the minimum number of times the least frequent items needs to be seen before the sampling loop terminates; in this paper, we set that value to 2. Our methods could also be further improved by in-depth investigation of prompts which consistently produce close-to-uniform distributions across different families of language models.

Finally, while our methods do not currently address the case where temperature annealing is used in conjunction with top-$k$ or top-$p$, extending them to support this setting should be straightforward. Temperature followed by top-$k$ is still top-$k$, and should be detectable via our methods. Temperature followed by top-$p$ is trickier, because we no longer have a known distribution. However, this combination can be detected by comparing the empirical distribution against a set of known distributions for common models; if the distribution does not match any of them, then we can conclude that either it is not using any known model, or that other distribution shaping such as temperature has been applied.

### 5 Limitations

Our method is limited to identifying when top-$p$ sampling or top-$k$ sampling is used. We do not attempt to detect other decoding strategies which other systems might use. Additionally, there is no guarantee that a system would use a single decoding strategy—it is

possible that different prompts may trigger different decoding strategies, or that A/B testing results in different users seeing different decoding strategies. Our ChatGPT results were computed by two different authors on separate OpenAI accounts. Also, we have no guarantees that the decoding strategy is not changed over time. Some of our ChatGPT results were computed using the December 15, 2022 release while others were computed using the January 9, 2023 release.

Additionally, the biases in distributions that we see here could have other underlying reasons; for example, changes in the data can result in very different true distributions.

Furthermore, under the hood, an API might not be generating a new random generation each time an identical prompt is passed in. Either random seeds might be getting re-used, or generations could be retrieved from a cache. In both cases, the generations might look like argmax sampling. It's also conceivable that certain combinations of fixed models could look like top-$k/p$. For example, if a query is randomly routed to one of a series of $s$ servers, each serving a different model, we might interpret the decoding strategy to be top-$k$ even if each server is using argmax. In these cases, an approach more like that proposed by Tay et al. (2020), where classification of decoding strategy is made based on a long generated sequence (rather than single token system predictions, as in our approach), might be more effective.

For top-$p$ (though not top-$k$), we require access to an underlying distribution that approximates the model used. This is not an issue for open source models or models with API access that allows specifying the decoding strategy, but it does limit the applicability of our method to newer proprietary models. It may be possible to empirically determine distributions for carefully engineered prompts, but future work is needed for reverse engineering fully closed models.

### 6 Conclusion

Our attack shows that with even a little work, it is possible to reverse-engineer common decoding strategies. Although we have focused here only on top-$p$ and top-$k$ sampling, these approaches generalize readily to other common methodologies when the output probability distributions are well-approximated. Along with other recent work on reverse-engineering other parts of a language generation system (Zhang and Ippolito, 2023), it seems is infeasible to hide inference implementation details given black-box access to the system.

# References

Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. 2023. Pythia: A suite for analyzing large language models across training and scaling. *arXiv preprint arXiv:2304.01373*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 889–898.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration. In *International Conference on Learning Representations*.

Daphne Ippolito, Daniel Duckworth, Chris Callison-Burch, and Douglas Eck. 2020. Automatic detection of generated text is easiest when humans are fooled. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1808–1822.

Ali Naseh, Kalpesh Krishna, Mohit Iyyer, and Amir Houmansadr. 2023. On the risks of stealing the decoding algorithms of language models. *arXiv preprint arXiv:2303.04729*.

George Pólya. 1930. Eine wahrscheinlichkeitsaufgabe in der kundenwerbung. *Zeitschrift Angewandte Mathematik und Mechanik*, 10(1):96–97.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.

Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*.

Yi Tay, Dara Bahri, Che Zheng, Clifford Brunk, Donald Metzler, and Andrew Tomkins. 2020. Reverse engineering configurations of neural text generation models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 275–279, Online. Association for Computational Linguistics.

Gian Wiher, Clara Meister, and Ryan Cotterell. 2022. On decoding strategies for neural text generators. *Transactions of the Association for Computational Linguistics*, 10:997–1012.

Sina Zarrieß, Henrik Voigt, and Simeon Schüz. 2021. Decoding methods in neural language generation: a survey. *Information*, 12(9):355.

Yiming Zhang and Daphne Ippolito. 2023. Prompts should not be seen as secrets: Systematically measuring prompt extraction attack success. *arXiv preprint arXiv:2307.06865*.

## A Algorithm for Estimating $k$

---

**Algorithm 1** Algorithm for estimating $k$.

---

Given a system `Gen:m↦r` that takes an input prompt $m$ and outputs a single response token $r$, let $m_1$ and $m_2$ be two prompts that with high probability return responses from large ($\gg k$) sets of *different* sizes—e.g. $m_1$ returns random nouns and $m_2$ returns random adverbs.

**function** ESTIMATEK(Prompts $m_1,m_2$;`Gen:m↦r`)
    samples1 ← [], samples2 ← []
    **while** max number iterations not reached **do**                     // we set max iterations=32
        **for** {1..100} **do**                                  // at most 3200 samples generated
            samples1.insert(`Gen`($m_1$))
            samples2.insert(`Gen`($m_2$))
        **end for**
        $k_1$ ← # unique items in samples1
        $k_2$ ← # unique items in samples2
        minSamples ← (samples1$^{(k_1)}$ > 1 and
            samples2$^{(k_2)}$ > 1)                    // Boolean testing all items appear twice
        **if** $k_1 = k_2$ and minSamples **then**
            **break**
        **end if**
    **end while**
    **return** $\lfloor (k_1+k_2)/2 \rfloor$                 // guesses average if convergence not reached
**end function**

---

As we mentioned in the main paper text, suppose, for a given prompt $m$, we call `Gen`($m$), $n$ number of times, each time keeping just the first token of the output. We can trivially lower bound $k$ by observing the number of unique items in a set of responses. As $n$ approaches $\infty$, all $k$ allowed responses will be observed. Since this is infeasible, the adversary's goal is to estimate $k$ while keeping $n$ as small as possible.

It is easy to see that the ideal prompt $m$ is one that gives responses that are perfectly uniform over the entire vocabulary $\mathcal{V}$. In the uniform case, we are left with the standard coupon collector problem (Pólya, 1930). We would recover the exact value of $k$ with probability at least $1-\frac{1}{k}$ by setting $n > 2k\log k$. Unfortunately, such a prompt is exceedingly difficult to engineer (see Appendix D).

It turns out we can do almost as well without needing full uniformity. The key building block for our attack is the construction of an $m$ that distributes substantial probability mass onto a subset $\mathcal{V}_m \subseteq \mathcal{V}$ of the token space. We require that for any $k < |\mathcal{V}_m|$, we have $\text{Prob}\big(\texttt{Gen}(m) = x^{(k)}\big) \geq \frac{1}{ck}$, for some small constant $c$. Put in plain language, we want to ensure that for any number of tokens $k$ the distribution might be truncated at, the least likely token that can be generated is no more than $c$ times less likely to appear than if the distribution were truly uniform. If $n \geq 2ck\log(ck)$, then with probability at least $1-\frac{1}{ck}$, our prediction is exactly correct. This result is far from tight, but follows easily from coupon collector on a uniformly random set of size $ck$.

In practice, we use Algorithm 1, which repeatedly estimates a lower bound for $k$ using two different prompts $m_1$ and $m_2$ for increasing numbers of trials until (1) the two estimates match and (2) the $x^{(k)}$ token appears at least twice in both generations (to prevent spurious matching). In such a case, the expected number of trials $n$ is approximately bounded above by $2ck\log(ck)$ via coupon collector[6].

## B Algorithm for Estimating $p$

Our goal is to estimate $p$ to within a factor of $\epsilon$. This would be trivial to do if we could construct a prompt $m$ that is uniform over a subset $\mathcal{V}_m \subseteq \mathcal{V}$ of size at least $\frac{1}{\epsilon}$. Then estimating $p$ would be equivalent to estimating top-$k$ for $k \approx \frac{p}{\epsilon}$ because each unique token seen implies a probability mass of $\epsilon$.

It is impossible to design a prompt which yields a perfectly uniform distribution. However, although uniformity is desirable, for top-$p$ estimation, it is more important that the distribution of `Gen`($m$) is known, i.e., we have access to the underlying language model $f_\theta$. If $k$ distinct tokens appear in the $p$-truncated distribution,

---

[6]Aside: the constant 2 that appears in the expected number of trials is due to requiring that the $k$th most likely token appears at least twice. However, it is unrelated to the constant 2 that appears in the bound in the previous paragraph, which is chosen to ensure the $\frac{1}{ck}$ failure probability.

---

**Algorithm 2** Algorithm for estimating $p$.

---

Consider a language model $f_\theta : (m,r) \to \mathbb{R}$ that scores a prompt/response pair and a system `Gen`$: m \mapsto r$ that takes an input prompt $m$ and outputs a single response token $r$ using $f_\theta$ and top-$p$ sampling. Let $m_1$ and $m_2$ be two prompts that return responses from known distributions over relatively small sets ($|\mathcal{V}_m|$ around 10-40)—e.g. $m_1$ returns random months and $m_2$ returns random dates within the month of March.

**function** ESTIMATEP(Prompts $m_1, m_2$; `Gen`$:m \mapsto r, f_\theta$)

    $p_1 \leftarrow$ HELPER($m_1$, `Gen`, $f_\theta$)

    $p_2 \leftarrow$ HELPER($m_2$, `Gen`, $f_\theta$)

    **return** $(p_1 + p_2)/2$

**end function**

**function** HELPER(Prompt $m$; `Gen`, known LM $f_\theta$ )

    baseProbs $\leftarrow []$                                             // Will store known probability distribution

    **for** $v \in \mathcal{V}_m$ **do**                               // $\mathcal{V}_m$ is the subset of tokens we consider

        baseProbs.insert($P_{f_\theta}(r = v | m)$)           // Probabilities using full random sampling

    **end for**

    Sort baseProbs from largest to smallest.

    baseProbs.insert($\sum_{v \in \mathcal{V} - \mathcal{V}_m} \text{Prob}_{f_\theta}(r = v | m)$)       // Summed probabilities of all out-of-set tokens

    samples $\leftarrow []$

    **for** $\{1..N\}$ **do**

        samples.insert(`Gen`($m$))

    **end for**

    $l \leftarrow$ num unique items in samples

    $p \leftarrow \sum_{i=1}^{l} \text{baseProbs}[i]$

    **return** $p$

**end function**

---

then (using the same notation as above), we can bound $p$ as:

$$\sum_{l=1}^{k-1} \text{Prob}_{f_\theta}(x^{(l)}) < p \le \sum_{l=1}^{k} \text{Prob}_{f_\theta}(x^{(l)}).$$

Thus, given a known distribution, the top-$p$ reverse engineering problem reduces to top-$k$.

Even if we do not know exactly the underlying model for a blackbox system, we can construct prompts that appear to often return distributions close to a family of known distributions. Then the error in estimating $p$ is just determined by how far off our guess of distribution is from the true underlying one. Note that to ensure robustness against an imperfectly guessed distribution, we estimate $p$ using the sum of the $k$ largest in-vocabulary probabilities, rather than trying to actually match the probabilities for the unique items sampled. This turns out to be important when prompts including exemplars are used, as the exemplars often create a bias in the tokens returned, but the overall drop-off in probabilities of most to least likely tokens tends to be more consistent. However, for distribution matching, we use the actual distributions over tokens.

In this paper, we set a goal of determining $p$ to within $\epsilon = 0.05$ and construct two prompts with almost-known distributions over $k = 13$ and $k = 31$. Because our guessed distributions are not perfect, instead of relying on a single distribution to bound our estimate, we instead average over the two upper bounds of $p$ derived from the different prompts and return that as our guess. Additionally, as $k$ is small for both prompts, instead of using the stopping criterion of Algorithm 1, for each prompt, we always generate 3000 samples. This means that with very high likelihood, we correctly return all possible items from the prompt's vocabulary. Algorithm 2 gives our implementation.

Figure 5: For each prompt style (DIGITS or ABC), we prompted with either [1, 2, 4, 8, 16, 32] exemplars selected randomly (RS), or with a random permutation of all expected outcomes (P). We ran three trials for each number of exemplars and generated the next word 5,000 times per trial. The majority of next-word generated were within the vocabulary, however, sometimes they were not, in those cases, we discarded that output. Missing bars indicate that there weren't enough generated next-words that were in the vocabulary to compute entropy.

## C  Prompt Selection

This Appendix gives more details on selecting good prompts for the decoding strategy detection task.

## D  Challenges

We encountered many challenges in selecting appropriate prompts. Our initial aim was to find prompts that induced an as-close-to-uniform distribution over the next token as possible. In addition to the prompts decided on for our main experiments (Table 1), we tried prompts meant to elicit a uniform distribution over digits, letters, dice rolls, and alphanumeric characters. For some of these prompt styles, the main difficulty was in getting the language model to assign higher probability to the expected outputs for the prompt than to unexpected outputs. For example, a prompt designed to elicit random digits would result in "and" being a more likely next token than several of the digits. For other prompts, the distribution was not as random as we would have expected. If exemplars were involved, even if they were chosen completely randomly, the model would try to follow any patterns observed in the exemplars. For example, if a prompt containing randomly selected exemplars of digits happened to end with "2 4 6", then "8" would be by far the most likely next token. Our difficulty here conforms with prior work that has shown that language models have significant biases toward certain numbers and words, even in settings where there should not be such bias.

Table 4: Prompts showcasing the sensitivity of models to different exemplar choices. The exemplars, shown in blue, can be varied in order and count.

| Name | Prompt v2 | $|\mathcal{V}_p|$ |
|---|---|---|
| DIGITS | Digits: 4, 3, 2 | 10 |
| ABC | Letters: E, F, P | 26 |

Figure 5 shows the variance in output distributions for two exemplar-based prompts, DIGITS and ABC (Table 4), across different numbers of exemplars and different random selection of exemplars. The DIGITS prompt is expected to output digits [0-9] with equal likelihood, and the random letters prompt is expected to output the letters [A-Z] with equal likelihood. While with enough exemplars, the DIGITS prompt yielded consistently high entropy (i.e., close to uniform-random) distributions, the ABC prompt did not consistently improve with more exemplars. In the end, we decided to avoid these prompts, and others which were too dependent on choice of exemplars.

For prompts to be used in top-$p$ estimation, one additional challenge is that ideally the prompt should yield a similar distribution when inputted to all popular language models. As discussed in the paper, our estimates for $p$ are worse when there is a greater mismatch between the known distribution used for top-$p$ estimation and the true distribution of the language model underlying the blackbox system being attacked. Figure 6 shows

Figure 6: The likelihoods for the digits {1-31} given the prompt "The accident occurred on March" (left) and for {January thru December + Ramadan} given the prompt "She arrived in the month of" (right), ordered from most to least likely. The sum probability of the remaining items in $\mathcal{V}$ is shown in red.

Table 5: Prompts uses to estimate the $p$ value for ChatGPT.

| Name | Prompt | $|\mathcal{V}_m|$ |
|------|--------|-------|
| MONTHSCHAT | write one word for the rest of this sentence: "She came to visit in the month of" | 13 |
| DATESCHAT | write one word for the rest of this sentence: "The accident occured on March" | 31 |
| D20CHAT | write one number for the rest of this sentence: "I rolled a D20 and the outcome was" | 20 |
| D100CHAT | Could you roll me a D100? We're playing D&D. Answer with just the roll value and nothing else. | 100 |

the known distributions for the two prompts we used in top-$p$ estimation, across several different models. We see that some models have much spikier distributions than others. The best approach (and the one we used to attack ChatGPT) is to choose the known distribution to use for top-$p$ estimation by keeping around a database of distributions from a bunch of different models, and then comparing the output distribution from the blackbox system to each distribution in the database. We can then choose to estimate $p$ using the known distribution with the lowest relative entropy with the blackbox's one.

## D.1 Chosen Prompt Details

Here we describe the actual prompts used in our experiments. For the NOUNS and ADVERBS prompts, we assumed access to the GPT-2 vocabulary and used Spacy (`en_core_web_sm`) to identify all tokens in the vocabulary corresponding to nouns and adverbs. In all experiments with these prompts, we used 16 randomly selected exemplars from these lists. An example prompt for NOUNS is: "List of nouns chosen completely randomly: negativity diarrhea problems eloqu money aspect vertex fraternity stone breast skies pushes probabilities ink north creditor". In our experiments estimating top-$k$, for each system being evaluated, we varied the random seed, resulting in a slightly different prompt. We did this to avoid any systematic biases resulting from always using the same choice of exemplars. For the non-exemplar-based prompts, we did not assume vocabulary access and instead relied on the expectation that letters, digits, and common words are present in most model vocabularies.

As mentioned in the main paper, different prompts were needed to attack ChatGPT than for the experiments on open-source models because ChatGPT expects its inputs to be in a conversation format and it does not offer control over the number of words generated (without careful prompt design, it tends to return tens to hundreds of words). Table 5 gives the prompts used to attack ChatGPT.

## E  Scientific Artifacts

We use the following language models in our research:

- **BLOOM 3B**: This model was released by BigScience under the RAIL License v1.0 with the goal to "to enable public research on large language models" (Scao et al., 2022). It can be downloaded at `https://huggingface.co/bigscience/bloom`.

- **Pythia 2.7B**: This model was released by EleutherAI under the MIT License with the goal of enabling

research on "interpretability analysis and scaling laws" (Biderman et al., 2023). It can be downloaded at `https://github.com/EleutherAI/pythia`.

- **GPT-2 base and XL**: These models were released by OpenAI under the MIT license with the goal of fostering language model research (Radford et al., 2019). They can be downloaded at `https://huggingface.co/gpt2`.

- **ChatGPT and GPT-3 model family**: These models are only available via OpenAI's API or through OpenAI's web interface. Our experiments with them fall under OpenAI's research policy, found at `https://openai.com/api/policies/sharing-publication/#research-policy`.

We chose these models evaluate on because (1) we wanted to evaluate our method on a wide range of independently trained models using different paradigms and training dataset choices. For example, though we conduct all our experiments using English prompts, we can observe the impact of BLOOM being trained on multilingual data, in that for the MONTHS prompt, BLOOM puts significant probability-mass on non-English month names, which could affect our $p$ estimates for BLOOM models.

## F  Computational Resources

Preliminary experiments were run in Google Colab using a Pro membership, which gave access to one Tesla T4. Subsequent experiments were running on a Google Cloud machine with 8 Tesla V100s. No more than 100 hours were spent running computation on this machine, which has a cost of $17 per hour.

## G  Additional Results

Table 6 gives the numbers used in Figure 2 in the main paper, as well as the root mean-square error between the true and estimated $p$ values.

Table 6: Performance at $p$ estimation across 100 estimations with $p$ values random from 0 to 1. On the left, GPT-2 Base was used to compute known distributions, and on the right GPT-3 was used to compute the known distributions.

| Model | GPT-2 Base | | GPT-3 Davinci v1 | |
|---|---|---|---|---|
| | Acc±.05 | RMSE | Acc±.05 | RMSE |
| GPT-2 Base | 0.93 | 0.03 | 0.08 | 0.19 |
| GPT-2 XL | 0.82 | 0.04 | 0.07 | 0.21 |
| Davinci v1 | 0.23 | 0.14 | 0.51 | 0.06 |
| BLOOM-3B | 0.77 | 0.04 | 0.07 | 0.22 |
| pythia-2.7b | 0.88 | 0.03 | 0.07 | 0.20 |

# Controlling keywords and their positions in text generation

**Yuichi Sasazawa, Terufumi Morishita, Hiroaki Ozaki,**
**Osamu Imaichi, Yasuhiro Sogawa**
Hitachi, Ltd. Research and Development Group
{yuichi.sasazawa.bj, terufumi.morishita.wp, hiroaki.ozaki.yu,
osamu.imaichi.xc, yasuhiro.sogawa.tp}@hitachi.com

## Abstract

One of the challenges in text generation is to control text generation as intended by the user. Previous studies proposed specifying the keywords that should be included in the generated text. However, this approach is insufficient to generate text that reflect the user's intent. For example, placing an important keyword at the beginning of the text would help attract the reader's attention; however, existing methods do not enable such flexible control. In this paper, we tackle a novel task of controlling not only keywords but also the position of each keyword in the text generation. To this end, we propose a task-independent method that uses special tokens to control the relative position of keywords. Experimental results on summarization and story generation tasks show that the proposed method can control keywords and their positions. The experimental results also demonstrate that controlling the keyword positions can generate summary texts that are closer to the user's intent than baseline.

## 1 Introduction

One of the challenges in text generation is to generate text that is consistent with the user's intent. Many methods for specifying the keywords that should be included in the generated text to reflect the user's intent have been proposed. As for summarization, by providing the model with keywords that should be included in the summary, it is possible to generate summaries that focus on specific parts of the document (Fan et al., 2018; He et al., 2022; Dou et al., 2021). As for story generation, keywords are used to control the narrative storyline (Jain et al., 2017; Fan et al., 2019; Yao et al., 2019). As for other tasks, such as e-commerce generation, review generation, and question generation, keywords are also used to control text generation (Chan et al., 2019; Shao et al., 2021; Ni and McAuley, 2018; Chan et al., 2021; Zhang and Zhu,

2021). In addition, more-advanced methods that specify the order of keywords to be included in the generated text to control the rough storyline have been proposed (Su et al., 2021; Shao et al., 2021).

The above-described methods, however, cannot generate texts that reflect more fine-grained intentions. Specifically, the user may want to reflect the intended importance of each keyword in the generated text. An effective way to reflect the intended importance of each keyword is to adjust the position of keywords within the text. For example, important keywords such as topic words and eye-catching words can be placed at the beginning of the text to attract the reader's attention, while the keywords for supplementary information can be placed in the middle or later in the text. By controlling the specific position of each keyword according to its importance, it is possible to generate appropriate text for each situation. That is, controlling the specific position of keywords in the generated text is a challenge in terms of reflecting more-specific user intentions and generating texts that attract readers. However, as far as we know, no previous work has tackled this challenge.

In this paper, we tackle a novel task of controlling keywords and the position of each keyword in text generation. Inspired by previous work that controlled text attributes by using special tokens (Iwama and Kano, 2019; Lakew et al., 2019; Martin et al., 2020), we propose a task-independent method that uses special tokens to control text generation. Specifically, the position of the keyword is specified by providing the model with a special token that represents the target relative position of the keyword (0-10%, 10-20%, etc.) and length of target text (20-24 words, 25-29 words, etc.). We use relative positions (rather than absolute positions) because it is more practical to specify relative positions such as "at the beginning," "in the middle," or "at the end" of the target text. Moreover, length of
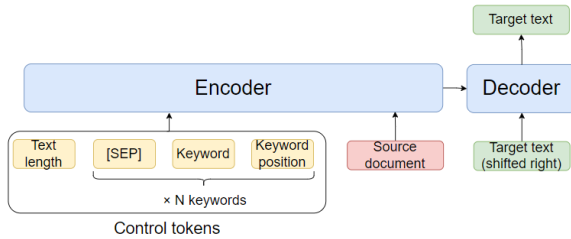
Figure 1: Overview of proposed method. The model is provided with control tokens: keywords in the target text, positions of each keyword, and target-text length to control text generation.

the target text is controlled because text length is considered to be one of the important factors that users want to control when considering where to place keywords. During training of the model, the model is provided with control tokens, including keywords randomly extracted from the target text, the positions of each keyword, and the length of the target text. The model is trained with cross-entropy loss in the same manner as conventional text generation; as a result, the model can learn the correspondence between the input control tokens and the target text.

The proposed "task-independent text-generation-control method" ("proposed method" hereafter) was comprehensively evaluated by applying it to summarization and story-generation tasks. The results of the evaluation show that the proposed method can control keywords and their positions in both tasks (Section 3.2). They also show that the proposed method can generate summary texts that are more similar to the gold summary than the baseline, indicating that text closer to the user's intent can be generated (Section 3.3). Case studies show that a model specifying keyword position can reflect the user's fine-grained intention (Section 3.4).

## 2 Method

### 2.1 Models

A BART model (Lewis et al., 2020) is used for the summarization task, and a GPT model (Radford et al., 2018) is used for the story-generation task. When the BART model is used, the source document is combined with the control tokens (i.e., keywords in the text to be generated, positions of each keyword, and length of the text to be generated) and given to the encoder as shown in Figure 1. When the GPT model is used, the control tokens are given to the decoder. As with regular text gen-

eration using BART and GPT models, the model is trained to maximize the conditional probabilities $p(y_i|y_{<i}, x)$ by using cross-entropy loss, where $y$ denotes the target text and $x$ denotes the input to the model, including the control tokens and the source document used in summarization task.

### 2.2 Control tokens

Inspired by existing methods that control text attributes by special tokens (Iwama and Kano, 2019; Lakew et al., 2019; Martin et al., 2020), the model is provided with the position of each keyword and text length as special tokens. For example, if the keyword phrase "two dogs" is located in the first 20-30% of the text and text length is in the range of 50-54 words, "[LENGTH50][SEP]two dogs[POSITION20]" will be given to the model as the control token. Here, [LENGTH50] and [POSITION20] are new tokens added to the vocabulary, and the corresponding word embedding is initialized randomly.

Note that control tokens that represent the oracle information of the target text are given to the model during both training and inference. This setting is appropriate because we aim to generate the intended text by providing additional information to the model. It is also possible that the model automatically determines keywords and their positions (i.e., control tokens are not given to the model), but that approach is left for future work.

Control tokens are extracted from the target text as follows. More details are given in Appendix A.3.

**Keywords** Keywords in this paper are not limited to important words in the target text; they can also be any phrase consisting of one to three consecutive words in the target text. For example, from the target text "Marcia was looking forward to trying hang gliding.", the phrases "Marsha", "was", "looking forward", "to trying", and "trying hang gliding" are first extracted as keyword candidates. However, frequent words with little meaning such as "was" and "to trying" are excluded from the keyword candidates, because they are considered unlikely to be given as keywords by the user. During training, a random number of phrases from the keyword candidates are given to the model as keywords. During inference, the user has the flexibility to give arbitrary keywords to the model. However, for the experiments conducted in this paper, we follow the same approach as during training: the keywords are randomly selected from the keyword

| Control | CNN/DM | | XSum | | ROCStories | |
|---|---|---|---|---|---|---|
| | Include | Pos | Include | Pos | Include | Pos |
| One keyword | | | | | | |
| w/o Control | 27.5 | 8.3 | 23.4 | 9.4 | 0.5 | 0.1 |
| Keyword | 71.3 | 18.7 | 86.4 | 28.7 | 53.0 | 14.3 |
| +Len | 72.7 | 20.4 | 85.8 | 30.8 | 50.9 | 13.5 |
| +Pos | 80.8 | 47.0 | **92.1** | 63.0 | 57.2 | 27.4 |
| +Pos+Len | **85.8** | **48.8** | 91.8 | **64.1** | **58.8** | **29.1** |
| Two keywords | | | | | | |
| Keyword | 52.4 | 5.1 | 74.1 | 14.1 | 22.9 | 1.6 |
| +Pos+Len | **75.9** | **28.6** | **85.9** | **46.4** | **31.1** | **7.9** |
| Three keywords | | | | | | |
| Keyword | 39.1 | 2.0 | 62.5 | 9.8 | 9.2 | 0.3 |
| +Pos+Len | **70.6** | **21.8** | **80.5** | **37.3** | **15.5** | **2.2** |

Table 1: Evaluation of the control of keywords and their positions in terms of (i) accuracy of generating text **Includ**ing all of the target keywords and (ii) accuracy of generating text in which all of the target keywords are placed in each target **Pos**ition.

candidates and given to the model.

**Keyword Position** The position of each keyword is expressed as a relative position. Specifically, the absolute position of the target keyword when counted from the beginning of the text is divided by the number of words in the text, and the quantized position in units of 10% are given to the model.

**Text Length** Number of words in the target text (quantized in 5 word units) is given to the model.

## 3 Experiment

### 3.1 Experiment setting

The proposed method was comprehensively evaluated by applying it to well-established summarization and story-generation tasks. These two tasks have different characteristics. As for summarization, the model extracts information from a source document and compresses it into a short text by using the given control tokens. As for story generation, the model generates text solely on the basis of the given control tokens. For summarization, we used the CNN/DailyMail (Hermann et al., 2015) and the XSum (Narayan et al., 2018) dataset and the BART$_{LARGE}$ model (400M parameters) (Lewis et al., 2020). For story generation, we used the ROCStories (Mostafazadeh et al., 2016) dataset and the GPT2 model (120M parameters) (Radford et al., 2018).

We extract candidate keywords from a target text by using the method described in Section 2.2. Dur-

ing training, no more than three keywords were randomly selected from the keyword candidates for each epoch and given to the model. During inference, one to three keywords randomly selected were given to the model in the experiment of Table 1, and one keyword randomly selected was given to the model in the experiment of Table 2 and Table 3.

In all experiments, training and inference were performed three times, and the mean score was reported. See Appendix A for more details on the experimental setup.

### 3.2 Evaluation of keyword-position control

Whether the given keywords are placed at given positions was evaluated first in terms of (i) the accuracy of generating text including all target keywords and (ii) the accuracy of generating text in which all target keywords are placed in each target position. As shown in Table 1, the proposed method using special tokens (+Pos and +Pos+Len) can generate text that includes the target keyword at the target position. Providing text-length information along with position information (+Pos+Len) improves the accuracy of keyword-position control, particularly in datasets with long text lengths (CNN/DM and ROCStories). In other words, combining relative position and length information enables the model to place the keywords in appropriate positions. The accuracy of the keyword inclusion is also improved when the keyword position is given. We suspect that the model was informed in advance of where the keywords should be placed; as a result, preventing the model from forgetting to place keywords in the text. It is clear that control accuracy is much lower in the case of story generation compared to summarization. This finding may be because the model is not given the source document and generates text from condition tokens only, so the model is more likely to generate the inappropriate context for keyword inclusion.

A more detailed evaluation is given in Table 2. For each target relative position of the keyword, the keyword position in the text was classified as (i) located in the target position, (ii) located at a positional deviation within 10%, (iii) located at a positional deviation greater than 10%, or (iv) not included in the text. It is clear from the results in the table that at all target positions, the accuracy of the keyword-position control is improved compared with that achieved using keyword-only

| Keyword position in the generated summary | Target-keyword position (relative position) | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | 0-10% | 10-20% | 20-30% | 30-40% | 40-50% | 50-60% | 60-70% | 70-80% | 80-90% | 90-100% |
| **Keyword only Control** | | | | | | | | | | |
| Correct position | 52.6 | 23.8 | 14.5 | 9.5 | 9.5 | 9.1 | 8.7 | 11.8 | 12.7 | 15.6 |
| **Keyword + Position + Length Control** | | | | | | | | | | |
| Correct position | **84.0** | **57.9** | **49.1** | **41.4** | **36.0** | **36.2** | **33.7** | **36.0** | **46.2** | **47.9** |
| Within 10% diff | 8.1 | 27.5 | 31.9 | 34.4 | 36.1 | 34.1 | 35.5 | 34.3 | 23.3 | 8.9 |
| Over 10% diff | 3.2 | 5.3 | 8.3 | 12.8 | 15.1 | 15.5 | 15.1 | 11.4 | 6.7 | 10.9 |
| Not included | 4.7 | 9.4 | 10.7 | 11.4 | 12.8 | 14.1 | 15.7 | 18.4 | 23.7 | 32.4 |

Table 2: Detailed evaluation of the control of the keyword and its position in the CNN/DM dataset. For each target relative position of the keyword, the keyword position in the text was classified as (i) located in the target position (**Correct position**), (ii) located at a positional deviation within 10% (**Within 10% diff**), (iii) located at a positional deviation greater than 10% (**Over 10% diff**), or (iv) not included in the text (**Not included**).

| Control | CNN/DM | | | XSum | | |
|---|---|---|---|---|---|---|
| | R1 | R2 | RL | R1 | R2 | RL |
| w/o Control | 43.6 | 20.6 | 40.5 | 44.3 | 21.1 | 36.5 |
| Keyword | 44.4 | 21.4 | 41.3 | 45.9 | 22.7 | 38.4 |
| +Len | 45.7 | 22.1 | 42.5 | 47.0 | 23.5 | 39.3 |
| +Pos | 44.9 | 21.9 | 41.8 | 46.7 | 23.6 | 40.2 |
| +Pos+Len | **46.4** | **22.8** | **43.2** | **47.8** | **24.5** | **41.2** |

Table 3: Summarization evaluation by ROUGE score. To reduce the effect on the ROUGE score due to giving target keywords, target keywords are excluded from both the target and generated summaries.

control, and that finding suggests the effectiveness of the proposed method. The results also show a high success rate of keyword inclusion and positional control near the beginning of the text, and a low success rate in the middle and at the end of the text. This may be because the closer to the end of the text, the more difficult it becomes for the model to generate text that contains the specified keywords while maintaining consistency with the context provided by the preceding words.

### 3.3 Evaluation of summary-content control

We show that controlling the text makes it easier for the user to generate the intended text in summarization. The results of the evaluation of summary-content control in summarization by ROUGE score (Lin, 2004) are shown in Table 3. Note that to reduce the effect on the ROUGE score due to giving target keywords, target keywords are excluded from both the target and generated summaries. It is clear from the results in the table that the score is improved by controlling keyword positions and text length, and that finding indicates that such control makes it easier to generate text that is close to the user's intended content.

### 3.4 Case study

To better understand how the proposed model behaved, representative examples of generated texts are shown in Table 4 and Table 5 [1]. In these examples, the keywords and their positions were controlled, although in some examples, the position of the keyword deviates slightly from the target position. It is clear from the table that by assigning different positions for the keywords, it was possible to generate several valid texts with different characteristics. For example, in the example in Table 4, placing the keyword "true miracle dog" at the 0% position generates a text that draws the reader's attention with an eye-catching keyword at the beginning of the text. In contrast, placing that keyword at the 90% position generates a narrative-style text that describes events in chronological order. It is also clear that even when multiple keywords are given, the order of the keywords can be adjusted by controlling the position of each keyword.

We also show some cases in which the proposed model produced errors. When a keyword position near the end of the text is specified, the instruction is often ignored, and the keyword is placed in a completely different position or not included in the text. As can be seen from the results in Table 2, the model tends to be poor at placing keywords at the back of the text.

When comparing the generated text of the summarization task with that of the story generation task, we observed that each of the specified keywords is usually used only once in the generated text of the summarization task, while each of the specified keywords is sometimes used multiple times in the generated text of the story generation

---

[1] A source document of summarization, gold texts, and additional examples of generated texts are given in Appendix D.

| Keyword & Position | Generated text |
|---|---|
| true miracle dog (0%) | "She's a **true miracle dog** and she deserves a good life," foster mother says. Theia was apparently hit by a car and buried in a field. Four days later, she was found emaciated and dirt-covered by a farm worker. A fundraising page has raised more than $10,000. |
| true miracle dog (40%) | A stray dog apparently hit by a car and buried in a field survives. The dog, named Theia, is "a **true miracle dog**," her foster owner says. She suffered a dislocated jaw, leg injuries and a caved-in sinus cavity. A fundraising page has raised more than $10,000. |
| true miracle dog (90%) | Stray dog apparently hit by a car and buried in a field survives. Four days later, she is found emaciated and dirt-covered on a farm. She has a dislocated jaw, leg injuries and a caved-in sinus cavity. A good Samaritan has set up a fundraising page for "the **true miracle dog**" |
| hit by (40%) bully breed mix (60%) | A stray dog in Washington State apparently died after being **hit by** a car. Four days later, she was found emaciated and dirt-covered in a field. Theia, a **bully breed mix**, has a dislocated jaw, leg injuries and a caved-in sinus cavity. |
| hit by (10%) bully breed mix (0%) | Theia, a **bully breed mix**, was apparently **hit by** a car and buried in a field. Four days later, she was found emaciated and dirt-covered by a farm worker. She suffered a dislocated jaw, leg injuries and a caved-in sinus cavity. A fundraising page has raised more than $10,000. |
| Keyword & Position | Generated text (Error cases) |
| buried in (80%) | Stray dog apparently hit by a car, apparently killed and **buried in** a field. Four days later, the dog manages to stagger to a nearby farm. She was found by a farm worker, who took her to a vet for help. Theia is now receiving treatment. |
| whacked (70%) | A stray dog in Washington State apparently died after being hit by a car. Four days later, she was found emaciated and dirt-covered by a farm worker. She suffered a dislocated jaw, leg injuries and a caved-in sinus cavity. A fundraising page has raised more than $10,000. |

Table 4: Examples of generated texts from the CNN/DM dataset

| Keyword & Position | Generated text |
|---|---|
| dealership (0%) | I went to the **dealership** to look at a car. I was very nervous about buying a car. I asked the salesman if I could try out a certain car. The salesman told me that I could try out a blue car. I drove the blue car home and loved it so much, I bought it. |
| dealership (30%) | I went to the **dealership** to buy a new car. I was very nervous about the car and the price. I went to the **dealership** and looked at the price. I decided to buy the car and I was very happy with it. I am glad I went to the **dealership** because it was a great deal. |
| dealership (80%) | I went to the **dealership** to look at a car. I was very nervous because I didn't know what I wanted. I asked the salesman if I could pick out a car. The salesman told me that I could get a new car if I paid $40,000. I drove to the **dealership** and bought a new car. |
| drive to work (50%) saved enough (70%) | I was driving to work one day when I saw a car in the road. I pulled over and asked if I could **drive to work**. The driver told me that he had **saved enough** money to buy a new car. I drove to work and paid him back. I drove to work and paid him back and he was very happy. |
| drive to work (20%) saved enough (0%) | I **saved enough** money to buy a new car. I went to the car dealership to test drive my new car. I drove the car for a few hours before I left. When I got home, I realized I had forgotten my wallet. I had to **drive to work** to get my wallet back, but I was happy. |
| Keyword & Position | Generated text (Error cases) |
| began saving money (90%) | I was in a hurry to get to work. I had to hurry because I didn't have my car keys. I looked everywhere for my keys, but couldn't find them. I finally found them under my coat, and I was relieved. I was able to grab my keys and walk to work without losing my car keys. |
| local (0%) Bob (40%) enough (60%) | **Bob** was a **local** handyman. He was hired to fix up a leaky roof on his home. **Bob** was very handy and did a good job at it. Unfortunately, the roof was too deep and the water would not come out. **Bob** had to call a **local** handyman to fix the leaky roof. |

Table 5: Examples of generated texts from the ROCStories dataset

task. This may be because the story generation task requires the model to generate text content conditionally only on the specified keywords, causing the model to become overly dependent on them.

# 4 Conclusion

A method for controlling keywords and the position of each keyword in generated text is proposed and evaluated experimentally by applying it to two tasks: summarization and story generation. The results of the evaluation show that the proposed method, which uses special tokens, can control the keyword positions in both tasks. They also show that the method can generate summary texts that are more similar to the gold summary than the baseline, and that finding indicates that text closer to the user's intent can be generated.

**Supplementary Materials Availability Statement**

**Source code**

- The source code is available at Github[2].

**Dataset**

- The CNN/DM dataset is available at Github[3].
- The XSum dataset is available at Hugging-Face[4].
- The ROCStories dataset is available at here[5].

# Acknowledgments

# References

Ying-Hong Chan, Ho-Lam Chung, and Yao-Chung Fan. 2021. Improving controllability of educational question generation by keyword provision. *CoRR*, abs/2112.01012.

Zhangming Chan, Xiuying Chen, Yongliang Wang, Juntao Li, Zhiqiang Zhang, Kun Gai, Dongyan Zhao, and Rui Yan. 2019. Stick to the facts: Learning towards a fidelity-oriented E-commerce product description generation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4959–4968, Hong Kong, China. Association for Computational Linguistics.

Zi-Yi Dou, Pengfei Liu, Hiroaki Hayashi, Zhengbao Jiang, and Graham Neubig. 2021. GSum: A general framework for guided neural abstractive summarization. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4830–4842, Online. Association for Computational Linguistics.

Angela Fan, David Grangier, and Michael Auli. 2018. Controllable abstractive summarization. In *Proceedings of the 2nd Workshop on Neural Machine Translation and Generation*, pages 45–54, Melbourne, Australia. Association for Computational Linguistics.

Angela Fan, Mike Lewis, and Yann Dauphin. 2019. Strategies for structuring story generation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2650–2660, Florence, Italy. Association for Computational Linguistics.

Junxian He, Wojciech Kryscinski, Bryan McCann, Nazneen Rajani, and Caiming Xiong. 2022. CTRL-sum: Towards generic controllable text summarization. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5879–5915, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. In *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc.

Kango Iwama and Yoshinobu Kano. 2019. Multiple news headlines generation using page metadata. In *Proceedings of the 12th International Conference on Natural Language Generation*, pages 101–105, Tokyo, Japan. Association for Computational Linguistics.

Parag Jain, Priyanka Agrawal, Abhijit Mishra, Mohak Sukhwani, Anirban Laha, and Karthik Sankaranarayanan. 2017. Story generation from sequence of independent short descriptions. *CoRR*, abs/1707.05501.

Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.

Surafel Melaku Lakew, Mattia Di Gangi, and Marcello Federico. 2019. Controlling the output length of neural machine translation. In *Proceedings of the 16th International Conference on Spoken Language Translation*, Hong Kong. Association for Computational Linguistics.

Jason Lee, Elman Mansimov, and Kyunghyun Cho. 2018. Deterministic non-autoregressive neural sequence modeling by iterative refinement. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1173–1182, Brussels, Belgium. Association for Computational Linguistics.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Louis Martin, Éric de la Clergerie, Benoît Sagot, and Antoine Bordes. 2020. Controllable sentence simplification. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 4689–4698, Marseille, France. European Language Resources Association.

Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Vanderwende, Pushmeet Kohli, and James Allen. 2016. A corpus and cloze evaluation for deeper understanding of commonsense stories. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 839–849, San Diego, California. Association for Computational Linguistics.

Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1797–1807, Brussels, Belgium. Association for Computational Linguistics.

Jianmo Ni and Julian McAuley. 2018. Personalized review generation by expanding phrases and attending on aspect-aware representations. In *Proceedings*

*of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 706–711, Melbourne, Australia. Association for Computational Linguistics.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2018. Language models are unsupervised multitask learners.

Huajie Shao, Jun Wang, Haohong Lin, Xuezhou Zhang, Aston Zhang, Heng Ji, and Tarek Abdelzaher. 2021. Controllable and diverse text generation in e-commerce. In *Proceedings of the Web Conference 2021*, WWW '21, page 2392–2401, New York, NY, USA. Association for Computing Machinery.

Yixuan Su, David Vandyke, Sihui Wang, Yimai Fang, and Nigel Collier. 2021. Plan-then-generate: Controlled data-to-text generation via planning. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 895–909, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. 2016. Rethinking the inception architecture for computer vision. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2818–2826.

Lili Yao, Nanyun Peng, Ralph Weischedel, Kevin Knight, Dongyan Zhao, and Rui Yan. 2019. Plan-and-write: Towards better automatic storytelling. In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence and Thirty-First Innovative Applications of Artificial Intelligence Conference and Ninth AAAI Symposium on Educational Advances in Artificial Intelligence*, AAAI'19/IAAI'19/EAAI'19. AAAI Press.

Wonjin Yoon, Yoon Sun Yeo, Minbyul Jeong, Bong-Jun Yi, and Jaewoo Kang. 2020. Learning by semantic similarity makes abstractive summarization better. *CoRR*, abs/2002.07767.

Zhiling Zhang and Kenny Zhu. 2021. Diverse and specific clarification question generation with keywords. In *Proceedings of the Web Conference 2021*, WWW '21, page 3501–3511, New York, NY, USA. Association for Computing Machinery.

Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. 2018. Texygen: A benchmarking platform for text generation models. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, SIGIR '18, page 1097–1100, New York, NY, USA. Association for Computing Machinery.

# Tackling Hallucinations in Neural Chart Summarization

**Saad Obaid ul Islam**[1,2,3]    **Iza Škrjanec**[2]    **Ondřej Dušek**[1]    **Vera Demberg**[2]

[1]Charles University, Faculty of Mathematics and Physics, Prague, Czechia
[2]Saarland University, Saarbrücken, Germany
[3]Julius-Maximilians-Universität Würzburg, Germany

`saadob12@gmail.com`, `{skrjanec,vera}@coli.uni-saarland.de`, `odusek@ufal.mff.cuni.cz`

## Abstract

Hallucinations in text generation occur when the system produces text that is not grounded in the input. In this work, we tackle the problem of hallucinations in neural chart summarization. Our analysis shows that the target side of chart summarization training datasets often contains additional information, leading to hallucinations. We propose a natural language inference (NLI) based method to preprocess the training data and show through human evaluation that our method significantly reduces hallucinations. We also found that shortening long-distance dependencies in the input sequence and adding chart-related information like title and legends improves the overall performance.

## 1 Introduction

The task of generating a summary to accompany a chart is an instance of data-to-text generation and has a long tradition in natural language generation (NLG) (Elzer et al., 2007; Ferres et al., 2007; Demir et al., 2012). Recent neural models for chart summarization (Obeid and Hoque, 2020; Hsu et al., 2021; Zhu et al., 2021; Kantharaj et al., 2022) carry the promise to be trainable from data and hence more versatile than approaches using manually constructed templates, and to produce more fluent text than previous statistical NLG systems. However, texts generated by state-of-the-art neural systems frequently include information which is not grounded in the input ("extrinsic hallucination"), or is even contradictory to it ("intrinsic hallucination"), see an example in Table 1.

Hallucinations in NLG (Koehn and Knowles, 2017; Raunak et al., 2021) have been a concern in neural models for various tasks (Huang et al., 2021; Lee et al., 2019; Rebuffel et al., 2022). We identify two reasons for hallucinations in chart summarization: (1) complexity and missing information in the input format of chart data; (2) presence of ungrounded information in chart summaries of the training data. Our contributions are as follows:



Road rage behavior among drivers in the U.S. as of 2015.

This statistic shows the road rage behavior of drivers in the United States as of 2015. <u>Four percent of the drivers said they have been on the receiving end of a rude gesture.</u> *The survey was conducted online and all the participants had a valid U.S. driving license.*

Table 1: This output example from the chart-to-text NLG system by Kantharaj et al. (2022) includes <u>intrinsic</u>, and *extrinsic* hallucinations.

- We demonstrate the importance of providing more context and reducing long-distance dependencies in the linearized input format.

- We propose an NLI cleaning step to remove ungrounded information in the training data.

Our experimental code and model output will be released on Github under an open license.[1]

## 2 Background and Related Work

### 2.1 Recent work in Chart Summarization

Several chart summarization datasets and models were developed recently. Obeid and Hoque (2020) created the Chart-to-Text data with English charts from statista.com (dubbed c2t-small in this paper). They model chart summarization as a data-to-text problem and adapt a transformer by Gong et al. (2019). Kantharaj et al. (2022) released an extended dataset crawled from the same platform, also called Chart-to-text (c2t-big in this paper).

---

[1]https://github.com/WorldHellow/Hallucinations-C2T

414

| Dataset | Training | Validation | Test | Total |
|---------|----------|------------|------|-------|
| c2t-small | 5,703 | 1,222 | 1,222 | 8,147 |
| c2t-big | 24,367 | 5,222 | 5,222 | 34,811 |

Table 2: Dataset sizes and splits: c2t-small by Obeid and Hoque (2020) and c2t-big by Kantharaj et al. (2022).

They finetune multiple pretrained models, such as BART (Lewis et al., 2020) and T5 (Raffel et al., 2020). Table 2 shows the statistics of both datasets.

## 2.2 Hallucinations in Data-to-text NLG

In NLG, hallucination or unfaithful output means generated text that is not grounded in the input. Ji et al. (2023) describe two types of hallucinations: **Intrinsic Hallucinations** refer to generated output that contradicts the source content, and **Extrinsic Hallucinations** refer to output that cannot be verified by the source. Ji et al. (2023) name three main causes for hallucinations: (1) source-reference divergence (reference text not supported by the input data), (2) modeling choices, and (3) decoding strategies. Efforts are made to mitigate hallucinations in NLG, particularly in data-to-text models. At the data level, clean and faithful datasets such as ToTTo (Parikh et al., 2020) and RotoWire-FG (Wang, 2019) were created with significant manual effort. For data preprocessing, Nie et al. (2019) use a domain-specific language understanding module to automatically remove noise from input meaning representations and reduce hallucinations. A more domain-general automatic approach was proposed by Dušek and Kasner (2020) with a transformer finetuned for NLI, but it was only applied to evaluate faithfulness of generated text, not to correct it. Our approach combines the latter two by applying NLI for automatic data cleaning.

## 3 Problem Identification

We manually analyzed two samples of 50 generated chart summaries of Obeid and Hoque (2020)'s transformer trained on c2t-small and Kantharaj et al. (2022)'s T5 model finetuned on c2t-big. For the c2t-small transformer, 22 summaries had intrinsic and 13 had extrinsic hallucinations. Some summaries were also incoherent and repetitive. For the c2t-big T5 model, 4 out of 50 summaries had intrinsic and 11 had extrinsic hallucinations.

Based on our analysis, we speculate that intrinsic and extrinsic hallucinations have different causes

Road rage behavior among drivers in the U.S. as of 2015 x-y labels situation - share of respondents x-y values On the receiving end of a rude gesture 53%, Yelled or used profanity 26%, Made a rude gesture 17%, Felt physically threatened 13%, Exited their vehicle to engage angrily 4%

Table 3: Example of our proposed linearization (See the chart in Table 1). After the chart title, the input contains $x$ and $y$ axis labels, followed by $x$-$y$ pairs, such that each $y$ value is adjacent to its corresponding $x$ value.

as identified by Ji et al. (2023). Intrinsic hallucinations may happen due to input sequence formatting (i.e., modeling choices), while extrinsic hallucinations are caused by the source-reference divergence in training data (missing input information and noise) Maynez et al. (2020). We address the input format in Section 4 by adding context and reducing long-distance dependencies, and we further address noise in references in Section 5.

## 4 Input Format Adjustment

### 4.1 Context and Distance in Input Formatting

The source data table must be linearized for input into a sequence-to-sequence model. Obeid and Hoque (2020)'s linearization includes $x$ and $y$ axis labels, values, and chart type, but lacks the chart title (see Table 8 in the Appendix for an example). We speculate that excluding the title results in extrinsic hallucination, i.e., generation of entities from parametric knowledge instead of the input data (Longpre et al., 2021).

Kantharaj et al. (2022)'s T5-based approach produced far better results than Obeid and Hoque (2020), but some hallucinations were still present. They format the data table by following the template: *title + y-values + x-values* (see Table 9 in the Appendix for an example). This format includes the title, but it lacks the $x$ and $y$ axis labels and the corresponding $x$-$y$ values are not adjacent. The distance between each $x$ and its corresponding $y$ value is large, and we speculate the model faces difficulty when learning pairwise relationships between $x$ and $y$, leading to intrinsic hallucination.

### 4.2 Proposed Input Formatting

Considering these input format problems, we hypothesise that reducing long-distance dependencies between $x$ and $y$ axis values in the linearized input data will alleviate intrinsic hallucinations; adding title and $x$ and $y$ axis labels should reduce extrinsic ones. We thus propose a linearized input with adja-

cent *x-y* pairs. The template we use is: *title + x-y labels + x-y values*. See Table 3 for an example.

### 4.3 Experimental Setup

We finetune T5 (Raffel et al., 2020) with our linearization proposed in Section 4.2, comparing to both original linearizations discussed in Section 4.1. We include ablated versions to check the effects of including the title, including axis labels, or using adjacent *x-y* pairs. As a prefix to T5's decoder input, we use *"C2T: "*. More training details are provided in Table 11 in the Appendix. We evaluate using BLEU (Post, 2018), ROUGE-2 (Lin, 2004), perplexity,[2] and NUBIA (Kane et al., 2020). NUBIA produces a score based on logical agreement, contradiction, neutrality, and semantic similarity.

### 4.4 Metrics Results

On the c2t-small data, *T5-S-OL* (our linearization) is compared to the original model of Obeid and Hoque (2020), a T5 finetuned using their linearization (*T5-S-O&HL*), and an ablation variant which uses their linearization and adds the chart title (*T5-S-O&HL + title*). Results in Table 4 show that our linearization improves almost all metrics. While the overall NUBIA score is lower, its most important elements (logical agreement, contradiction, semantic similarity) are improved (cf. Section 6.2). Obeid and Hoque (2020)'s input format produces many entity hallucinations. Including the chart title format improves performance substantially, which is expected as this provides crucial context for the model. Further small gains stem from less redundancy in our linearization.[3]

On the c2t-big data, *T5-B-OL* (our linearization) is compared to the original T5 model of Kantharaj et al. (2022) and an ablation using their linearization with added axis labels (*T5-B-K + axis labels*). Table 5 shows improvements on almost all metrics, with NUBIA not reflecting its individual elements' improvements, similar as above (cf. Section 6.2). Adding axis labels to Kantharaj et al. (2022)'s format is a very modest help, but using adjacent *x-y* pairs in our format yields a larger improvement.[4]

### 4.5 Manual Analysis

We manually analyzed 50 output samples from T5-S-OL, checking for hallucinations. To find intrinsic

hallucinations, we checked for any information in the summaries that would conflict with the input (x-y values, entities, or trends). For extrinsic hallucinations, we checked for the presence of any information that was not verifiable from the input data. We found no intrinsic hallucinations, but 18 summaries still had extrinsic hallucinations. Table 7 in the Appendix provides example outputs.

## 5 Cleaning Noisy References with NLI

### 5.1 Noise in Training Summaries

Since source-reference divergence can also cause hallucinations, we analyzed the reference side of the same sample of 50 instances from the c2t-small dataset as in Section 3 to look for text that is not grounded in the source chart. 20 out of 50 summaries contained ungrounded information. While this ungrounded information makes the summaries more interesting, it cannot be verified from the chart and hence counts as extrinsic hallucination. We also analyzed references in the c2t-big dataset and found a similar pattern, which is expected since both datasets come from the same source.

### 5.2 Influence on Generation

To show that ungrounded information in training data influence system outputs, we run an experiment on the Autochart dataset (Zhu et al., 2021), which is handcrafted and thus guaranteed not to contain hallucinations. We introduce synthetic ungrounded text at random places in Autochart summaries using vanilla GPT-2 (Radford et al., 2019) generation prompted by preceding summary text, thus creating a noisy Autochart version.[5] The summary of the chart is segmented, and a random sentence from the summary is used as a prompt for GPT-2 to produce an ungrounded sentence. The generated text is then inserted at a random location in the segmented summary, creating a new summary with ungrounded information. We analyze 50 outputs from a T5 model finetuned on both original (T5-AC-orig) and noisy (T5-AC-noisy) versions. While we found no hallucinations in outputs of T5-AC-orig, we identified 27 extrinsic hallucinations in T5-AC-noisy's outputs. T5-AC-noisy also produced repetitive summaries (43 out of 50).

### 5.3 Improving Faithfulness using NLI

To alleviate extrinsic hallucinations caused by the training data, we propose using NLI, taking inspira-

---

[2]https://huggingface.co/docs/transformers/perplexity

[3]Axis labels and chart type are not repeated in our format, compare Table 8 and 10 in the Appendix.

[4]More on ablations in Appendix A.2.

[5]See Section A.1 in the Appendix for details.

| Model | BL↑ | RG-2↑ | PPL↓ | L↑ | C↓ | Neu | SS↑ | N↑ |
|---|---|---|---|---|---|---|---|---|
| Transformer by Obeid and Hoque (2020) | 18.5 | - | - | - | - | - | - | - |
| T5-S-O&HL | 26.1 | 33.5 | 7.4 | 5.5 | 67.8 | 26.5 | 3.0/5 | 35.4 |
| T5-S-O&HL + title | 31.0 | 44.0 | 16.0 | 24.8 | 23.5 | 51.5 | 3.0/5 | 60.6 |
| T5-S-OL | 33.9 | 44.8 | 7.5 | 33.2 | 22.3 | 44.4 | 3.5/5 | 46.9 |
| T5-S-OL-NLI | 34.2 | 43.7 | 7.1 | 33.1 | 10.2 | 56.5 | 3.5/5 | 44.5 |

Table 4: Evaluation results on c2t-small for input format improvements (Section 4) and NLI filtering (Section 5): BLEU-4 (BL), ROUGE-2 (RG-2), Perplexity (PPL), Logical Agreement (L), Contradiction (C), Neu (Neutrality), Semantic Similarity (SS) and the NUBIA (N) score.

| Model | BL↑ | RG-2↑ | PPL↓ | L↑ | C↓ | Neu | SS↑ | N↑ |
|---|---|---|---|---|---|---|---|---|
| T5 by Kantharaj et al. (2022) | 37.0 | 50.5 | 10.0 | 34.5 | 22.9 | 42.5 | 3.6/5 | 53.5 |
| T5-B-K + axis labels | 37.6 | 50.5 | 8.2 | 33.0 | 23.9 | 42.9 | 3.6/5 | 51.4 |
| T5-B-OL | 39.8 | 55.0 | 8.2 | 39.3 | 21.3 | 39.3 | 3.6/5 | 55.6 |
| T5-B-OL-NLI | 42.2 | 50.7 | 8.2 | 40.3 | 15.1 | 44.5 | 3.6/5 | 53.5 |

Table 5: Evaluation results for comparing linearization methods on c2t-big (see Table 4 for metrics).



Figure 1: Summary cleaning process using zero-shot BART-NLI.

| Model | VC↑ | OIP↓ | Inf↑ | Coh↑ | Flu↑ |
|---|---|---|---|---|---|
| T5-S-OL | 56.00% | 38.00% | 3.80/5 | 3.81/5 | 3.88/5 |
| T5-S-OL-NLI | 76.00% | 17.00% | 3.60/5 | 3.91/5 | 3.96/5 |
| $p$-value | 1.5e-5 | 1.0e-5 | 0.914 | 0.333 | 0.465 |

Table 6: Human evaluation results (see Section 5.5). We used a $\chi^2$ test for VC and OIP, and one-way ANOVA for Inf, Coh, and Flu.

tion from Pang et al. (2021) and Dušek and Kasner (2020). We use NLI as a preprocessing tool: any sentences in a summary that are not entailed in the linearized data will be discarded. We use the BART-MNLI pretrained model[6] for this. It is based on Yin et al. (2019)'s study on zero-shot text classification as an entailment problem and trained to produce an entailment score on a scale of 0-100 (with no specific neutral or contradiction labels).

Figure 1 shows a diagram of the overall cleaning process. We apply the filtering step to the entire c2t-small dataset by segmenting the summaries, passing each sentence as hypothesis and linearized data as premise. If the sentence gets an entailment score above the threshold of 0.3, we keep it, otherwise we discard the sentence.[7] Out of the 100 filtered summaries, 69 were correctly filtered, and 31 were incorrectly filtered. We also applied the above process to c2t-big dataset.

## 5.4 Finetuning T5 on Filtered Data

To evaluate the filtering, we compare two pairs of finetuned T5 models on both c2t-small and c2t-big, filtered (T5-B/S-OL-NLI) and non-filtered (T5-B/S-OL). All use our proposed linearization method. We use data splits described in Section 2.

Tables 4 and 5 show the results of T5 finetuned on the filtered data and the baselines. The automatic metrics do not show any clear trend for T5-S-OL-NLI: BLEU improves but ROUGE drops; NUBIA logical agreement stays unchanged, contradictions decrease but neutral statements (also interpretable as hallucinations) increase. For T5-B-OL-NLI, we get minor improvements in BLEU and NUBIA logical agreement, but again a drop in ROUGE. This does not indicate improvements; however, upon manual analysis of 50 summaries for the same charts as in Section 4.5, we found that only 4 outputs from T5-S-OL-NLI still contained hallucinations (compared to 18 for T5-S-OL).

---

[6] https://huggingface.co/facebook/bart-large-mnli
[7] The threshold was found empirically (see Appendix A.1.3).

## 5.5 Human Evaluation

We conduct a detailed human evaluation, comparing T5-S-OL (see Section 4.1) and T5-S-OL-NLI (trained in Section 5.4). We evaluate the following: (1) **Value Correctness (VC)**: Numbers/values in the summary are from the chart, (2) **Outside Information Presence (OIP)**: The summary contains information not grounded in the chart, (3) **Informativeness (Inf)**: The summary conveys a lot of information about the chart, (4) **Coherence (Coh)**: The summary content is orderly and logically consistent, and (5) **Fluency (Flu)**: The text is grammatically correct and is not repetitive.

We used Prolific to recruit English native speakers from the UK[8] and Google Forms to conduct the survey. For each model, 50 samples were used and split into 5 experiments with 10 samples each. Each sample was annotated by 5 participants, in total 25 participants completed the survey. Table 6 shows the result. Considering faithfulness (VC and OIP), the T5-S-OL-NLI model trained on filtered data performs significantly better than the baseline T5-S-OL, showing that our method of alleviating hallucinations via cleaning training summaries through NLI is effective.

## 6 Discussions

### 6.1 Ungrounded Information in Training Data

In Section 5.1, we reported on ungrounded information in training data and showed in Section 5.2 that this leads to hallucinations in generated outputs. While there are good reasons for ungrounded information in human-written summaries (e.g., providing additional detail/background or linking to other events; cf. Thomson and Reiter, 2020), using such data to train an end-to-end model that does not distinguish between describing the chart and providing additional information is not appropriate and leads to inaccurate outputs, which is unsuitable for real-world scenarios (cf. Maynez et al., 2020; Xu et al., 2021).

### 6.2 Metrics' Shortcomings in Assessing Hallucinations

None of the automatic metrics we used (see Section 4.3) measure hallucinations explicitly. BLEU and ROUGE are reference-based and prone to biases stemming from ungrounded information in references (see Section 6.1), such as assigning

higher scores to hallucinated outputs. NUBIA is a trained aggregate metric of several components (NLI-based logical agreement, contradiction and neutrality, plus semantic similarity and fluency). We found that this aggregation can lead to non-transparent or confusing results: even if individual components are clearly improved, NUBIA may drop (see Section 4.4). Therefore, we found logical agreement and semantic similarity scores to be the most useful constituents of NUBIA for indicating the presence of hallucinations in generated texts.

## 7 Conclusions

We show that reducing long-distance dependencies and providing more context on the model's input results in fewer intrinsic hallucinations, and demonstrate that extrinsic hallucinations are a result of ungrounded information in the training summaries. Furthermore, we show through human evaluation that employing NLI to filter training summaries results in a significant drop in hallucinations.

## Limitations

The main limitation of our work is that we were unable to eliminate the extrinsic hallucinations completely. In Section 5.3, we mentioned that 31 out of 100 summaries were not filtered correctly, meaning that these summaries were left with ungrounded information, which resulted in 4 out of 50 generated summaries with extrinsic hallucinations. BART-NLI is developed for linguistic input and we employ it to infer from non-linguistic input. The second limitation is that generated summaries are shorter on average. T5-S-OL-NLI generated 28% of the summaries of just a single sentence. We expected this problem because our method of filtering only removes sentences and does not replace them with statements entailed in the data. The final limitation is that our model is only limited to producing summaries in the English language as it is trained on English summaries.

---

[8] https://www.prolific.co/, the hourly pay rate was 9.5 GBP.

## Ethics Statement

The human evaluation study was approved by the ethics committee of the respective national professional linguistic association. All the annotators were from the United Kingdom and each annotator was paid according to the hourly minimum wage in the United Kingdom, i.e. 9.5 GBP. The annotators were paid immediately after the results were analyzed. We only collected Prolific IDs of the users and they were deleted after the analysis of the data.

## References

Steven Bird, Ewan Klein, and Edward Loper. 2009. *Natural language processing with Python: analyzing text with the natural language toolkit*. O'Reilly Media, Inc.

Seniz Demir, Sandra Carberry, and Kathleen F. McCoy. 2012. Summarizing information graphics textually. *Computational Linguistics*, 38(3):527–574.

Ondřej Dušek and Zdeněk Kasner. 2020. Evaluating semantic accuracy of data-to-text generation with natural language inference. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 131–137, Dublin, Ireland. Association for Computational Linguistics.

Stephanie Elzer, Edward Schwartz, Sandra Carberry, Daniel Chester, Seniz Demir, and Peng Wu. 2007. A browser extension for providing visually impaired users access to the content of bar charts on the web. In *Proceedings of the Third International Conference on Web Information Systems and Technologies - Web Interfaces and Applications*, pages 59–66.

Leo Ferres, Petro Verkhogliad, Gitte Lindgaard, Louis Boucher, Antoine Chretien, and Martin Lachance. 2007. Improving accessibility to statistical graphs: The IGraph-Lite system. In *Proceedings of the 9th International ACM SIGACCESS Conference on Computers and Accessibility*, Assets '07, page 67–74, New York, NY, USA. Association for Computing Machinery.

Li Gong, Josep Crego, and Jean Senellart. 2019. Enhanced transformer model for data-to-text generation. In *Proceedings of the 3rd Workshop on Neural Generation and Translation*, pages 148–156, Hong Kong. Association for Computational Linguistics.

Ting-Yao Hsu, C Lee Giles, and Ting-Hao Huang. 2021. SciCap: Generating captions for scientific figures. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3258–3264, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Yichong Huang, Xiachong Feng, Xiaocheng Feng, and Bing Qin. 2021. The factual inconsistency problem in abstractive text summarization: A survey. *arXiv preprint arXiv:2104.14839*.

Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *ACM Comput. Surv.*, 55(12).

Hassan Kane, Muhammed Yusuf Kocyigit, Ali Abdalla, Pelkins Ajanoh, and Mohamed Coulibali. 2020. NUBIA: NeUral based interchangeability assessor for text generation. In *Proceedings of the 1st Workshop on Evaluating NLG Evaluation*, pages 28–37, Online (Dublin, Ireland). Association for Computational Linguistics.

Shankar Kantharaj, Rixie Tiffany Leong, Xiang Lin, Ahmed Masry, Megh Thakkar, Enamul Hoque, and Shafiq Joty. 2022. Chart-to-text: A large-scale benchmark for chart summarization. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4005–4023, Dublin, Ireland. Association for Computational Linguistics.

Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. In *Proceedings of the 3rd International Conference on Learning Representations (ICLR 2015)*. ICLR.

Philipp Koehn and Rebecca Knowles. 2017. Six challenges for neural machine translation. In *Proceedings of the First Workshop on Neural Machine Translation*, pages 28–39, Vancouver. Association for Computational Linguistics.

Katherine Lee, Orhan Firat, Ashish Agarwal, Clara Fannjiang, and David Sussillo. 2019. Hallucinations in neural machine translation.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Rensis Likert. 1932. A technique for the measurement of attitudes. *Archives of Psychology, 22 140, 55*.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Shayne Longpre, Kartik Perisetla, Anthony Chen, Nikhil Ramesh, Chris DuBois, and Sameer Singh. 2021. Entity-based knowledge conflicts in question answering. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7052–7063, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1906–1919, Online. Association for Computational Linguistics.

Feng Nie, Jin-Ge Yao, Jinpeng Wang, Rong Pan, and Chin-Yew Lin. 2019. A simple recipe towards reducing hallucination in neural surface realisation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2673–2679, Florence, Italy. Association for Computational Linguistics.

Jason Obeid and Enamul Hoque. 2020. Chart-to-text: Generating natural language descriptions for charts by adapting the transformer model. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 138–147, Dublin, Ireland. Association for Computational Linguistics.

Richard Yuanzhe Pang, Adam Lelkes, Vinh Tran, and Cong Yu. 2021. AgreeSum: Agreement-oriented multi-document summarization. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3377–3391, Online. Association for Computational Linguistics.

Ankur Parikh, Xuezhi Wang, Sebastian Gehrmann, Manaal Faruqui, Bhuwan Dhingra, Diyi Yang, and Dipanjan Das. 2020. ToTTo: A controlled table-to-text generation dataset. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1173–1186, Online. Association for Computational Linguistics.

Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.

Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67.

Vikas Raunak, Arul Menezes, and Marcin Junczys-Dowmunt. 2021. The curious case of hallucinations in neural machine translation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1172–1183, Online. Association for Computational Linguistics.

Clément Rebuffel, Marco Roberti, Laure Soulier, Geoffrey Scoutheeten, Rossella Cancelliere, and Patrick Gallinari. 2022. Controlling hallucinations at word level in data-to-text generation. *Data Mining and Knowledge Discovery*, 36(1):318–354.

Craig Thomson and Ehud Reiter. 2020. A Gold Standard Methodology for Evaluating Accuracy in Data-To-Text Systems. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 158–168, Dublin, Ireland. Association for Computational Linguistics.

Hongmin Wang. 2019. Revisiting challenges in data-to-text generation with fact grounding. In *Proceedings of the 12th International Conference on Natural Language Generation*, pages 311–322, Tokyo, Japan. Association for Computational Linguistics.

Xinnuo Xu, Ondřej Dušek, Shashi Narayan, Verena Rieser, and Ioannis Konstas. 2021. MiRANews: Dataset and Benchmarks for Multi-Resource-Assisted News Summarization. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 1541–1552, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Wenpeng Yin, Jamaal Hay, and Dan Roth. 2019. Benchmarking zero-shot text classification: Datasets, evaluation and entailment approach. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3914–3923, Hong Kong, China. Association for Computational Linguistics.

Jiawen Zhu, Jinye Ran, Roy Ka-Wei Lee, Zhi Li, and Kenny Choo. 2021. AutoChart: A dataset for chart-to-text generation task. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021)*, pages 1636–1644, Held Online. INCOMA Ltd.

# A Appendix

## A.1 Experiment using Autochart

### A.1.1 Splits

The authors of autochart did not create any data splits. The total size of the dataset is 23,543 chart-summary pairs. From the dataset, we use 10,593 and split it with the ratio of 70:15:15.

### A.1.2 GPT-2 Noise Generation

To inject noise in the summaries, we first segment the summary using NLTK (Bird et al., 2009) sentence tokenizer. After segmenting the summary, we randomly pick a sentence and give it as a prompt to the GPT-2 model. For GPT-2 generation, we use greedy search. The generated sentence is then inserted at a random location in the segmented summary list, and then all the sentences are combined.

### A.1.3 Threshold Determination

We analyzed a random sample of 100 filtered summaries and found that the mean entailment score of the entailed sentences was 89, while the mean entailment score given to the non-entailed sentences was 8.7. This means that the model is sure when assigning the score, and making minor adjustments to the threshold would not lead to significant improvements.

### A.2 Ablations

We conducted ablation experiments to investigate where the improvement comes from in the linearization. We computed two ablations: (1) Obeid and Hoque (2020)'s linearization with added chart title (T5-S-O&HL+title), and (2) Kantharaj et al. (2022)'s linearization with added axis labels (T5-B-K + axis labels). Results are shown in Tables 4 and 5. Both ablation show improvements over the original linearizations, but do not match our scores, showing that all our format changes are helpful. For (1) resulting scores are between our linearization and Obeid and Hoque (2020)'s linearization, and for (2) resulting scores are very close to the original linearization results, however, they are slightly worse than ours. The most significant enhancement is due to (1), which is understandable as the model's performance improves with increased context. Likewise, the inclusion of x-y labels leads to an enhancement over the initial outcomes. However, the performance is further boosted when the input data values are combined.

### A.3 Human Evaluation Survey Details

### A.3.1 Consent Form

Each user was asked to sign the consent form based on the following text: This study is being conducted as part of ongoing research at [——]. If you have any questions or comments about the study, please contact us on Prolific. You must be at least 18 years old to participate. Your participation in this research is voluntary. There are no risks or benefits to participating in this study. In the next section we will ask for your Prolific ID. All data will be anonymized prior to analysis and Prolific IDs will not be published.

### A.3.2 Survey Description

Dear Participants, you will be evaluating summaries of charts. Choose the summary that has Value Correctness and Outside Information Presence. Rate the informativeness, coherence, and fluency of the summaries given the chart.

**Value Correctness:** Numbers/figures/values in the summary are from the chart. Here you determine which of the summaries are accurate.

**Outside Information:** Information that is not from the chart at all. Here you determine which of the summaries have information not taken from the chart.

**Informativeness:** The summary conveys a lot of information about the chart. 1 being the least informative and 5 being the most informative.

**Coherence:** The information included in the summary is orderly and logically consistent. Here you rate the coherence of the summary. 1 being the least coherent and 5 being the most coherent.

**Fluency:** Summary is grammatically correct and does not contain any repetitions. Here you rate the fluency of the summary. 1 being the least fluent and 5 being the most fluent.

### A.3.3 Evaluation of Measured Properties

Value correctness gives us a binary scores, meaning, either the summary has correct values or not. Similarly for outside information presence, we also get binary scores. For, informativeness, coherence, and fluency, we get scores out of 5-point Likert scale (Likert, 1932), 5 being the highest score, and 1 being the lowest score.

| T5-S-OL | T5-B-OL | T5-S-OL-NLI |
|---|---|---|
| This statistic shows road rage behavior among drivers in the United States. During the survey, 53 percent of respondents stated they had been on the receiving end of a rude gesture. *All the participants in this survey had a valid U.S. driving license.* | This statistic represents the road rage behavior among drivers in the United States as of April 2015. During the survey, 13 percent of respondents stated that they felt physically threatened by another driver to engage in angrily with another driver. *The survey was conducted online and all the participants had a valid U.S. driving license.* | This statistic shows the road rage behavior among drivers in the United States. 53 percent of respondents said they had been on the receiving end of a rude gesture and 26 percent of the respondents said they have yelled or used profanity at another driver. |

Table 7: Generated summaries from three different models for the chart in Table 1. The summaries from T5-S-OL and T5-B-OL contain *extrinsic hallucinations*.



Table 8: Linearized input format used by Obeid and Hoque (2020). Example from c2t-small dataset.



Table 9: Linearized input format used by Kantharaj et al. (2022). Example from c2t-big dataset.

| Sales volume of beer in Prince Edward Island ( P.E.I ) from FY 2012 to FY 2019 , by product type ( in million liters ) | Sales volume of beer in Prince Edward Island ( P.E.I ) from FY 2012 to FY 2019 , by product type ( in million liters ) labels Year - Packaged - Draught values 2019 8.62 1.13 , 2018 8.65 1.1 , 2017 8.19 0.98 , 2016 8.48 0.91 , 2015 8.39 0.83 , 2014 8.47 0.74 , 2013 8.84 0.65 , 2012 8.79 0.64 |

Table 10: The proposed input format that we use for both the c2t-small and c2t-big dataset following the formula $title$ + $x$-$y$ labels + $x$-$y$ values. Individual parts are separated by a comma.

| Model Version | Model Repository |
|---|---|
| Pre-trained T5-base (Raffel et al., 2020) | https://huggingface.co/t5-base |
| **Parameter** | **Value** |
| Maximum input length | 1024 |
| Maximum target length | 512 |
| Truncation | True |
| Padding | max_length |
| batch size | 2 |
| Optimizer | Weighted Adam (Kingma and Ba, 2015) |
| Learning rate | 3e-4 |
| Weight decay | 0.01 |
| Training epochs and hours for T5-S-O&HL | 6 epochs, 11 hours |
| Training epochs and hours for T5-S-OL | 6 epochs, 11 hours |
| Training epochs and hours for T5-AC-orig | 8 epochs, 6 hours |
| Training epochs and hours for T5-AC-noisy | 8 epochs, 6 hours |
| Training epochs for T5-S-OL-NLI | 6 epochs, 11 hours |
| Training epochs for T5-B-OL | 12 epochs, 37 hours |
| Training runs for all the models | Single run |
| Beam size | 4 |
| GPU | Tesla T4 16 GB |

Table 11: Hyper-parameters used and training details of our experiments.

# Learning Disentangled Meaning and Style Representations for Positive Text Reframing

**Sheng Xu**   **Fumiyo Fukumoto**   **Jiyi Li**   **Kentaro Go**   **Yoshimi Suzuki**
University of Yamanashi
{g22dts03,fukumoto,jyli,go,ysuzuki}@yamanashi.ac.jp

## Abstract

The positive text reframing (PTR) task which generates a text giving a positive perspective with preserving the sense of the input text, has attracted considerable attention as one of the NLP applications. Due to the significant representation capability of the pre-trained language model (PLM), a beneficial baseline can be easily obtained by just fine-tuning the PLM. However, how to interpret a diversity of contexts to give a positive perspective is still an open problem. Especially, it is more serious when the size of the training data is limited. In this paper, we present a PTR framework, that learns representations where the meaning and style of text are disentangled. The method utilizes pseudo-positive reframing datasets which are generated with two augmentation strategies. A simple but effective multi-task learning-based model is applied to fuse the generation capabilities from these datasets. Experimental results on Positive Psychology Frames (PPF) dataset, show that our approach outperforms the baselines, BART by five and T5 by six evaluation metrics. Our source codes and data are available online.

## 1 Introduction

Text style transfer (TST) has been a long history from the early works, i.e., the eariler attempts are the frame language-based systems (McDonald and Pustejovsky, 1985) and schema-based Natural Language Generation (Hovy, 1987) in the 1980s, and more recent attempts such as CTPM (contrastive transfer pattern mining) (Han et al., 2023) and TST BT (Text Style Transfer Back Translation) (Wei et al., 2023). The goal is to change the text style, such as formality, and politeness with preserving the sense of the input text. With a recent surge of interest in deep learning (DL) techniques, positive text reframing (PTR) has been explored as one of the sub-fields in the TST study. Likewise, human-annotated data such as Positive Psychology Frames
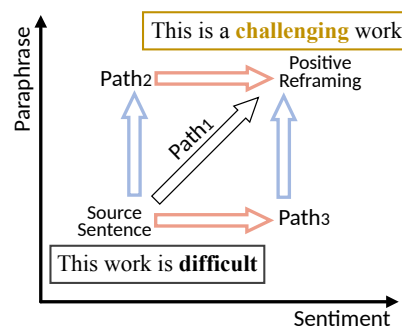


Figure 1: Disentangling Meaning and Style

(PPF) has been created for this task (Ziems et al., 2022).

One major approach for the TST task is to apply supervised learning for parallel data. Xu et al. (Xu et al., 2019) and Zhang et al. (Zhang et al., 2020) attempted multi-task learning for parallel data. To mitigate the small size of the parallel data, Rao (Rao and Tetreault, 2018) presented data augmentation strategies. Another attempt is to utilize a non-parallel dataset and train a model in an un-supervised manner (Shen et al., 2017; Fu et al., 2018). John et al. proposed a method that disentangles content- and style-related features and makes the decoder generate ideal output by using the disentangled features (John et al., 2019). Lai et al. designed two types of rewards for target style and content based on reinforcement learning (Lai et al., 2021). Many of these methods attained significant progress on the TST task while they still fail to handle the fine-grained transfer, i.e., disentangle style from content with preserving the meaning of the input that is required for the PTR task.

The main challenge in the PTR task is how to control diversity and the extent of style transfer. The concept of our PTR can be illustrated in Figure 1. The straightforward fine-tuning of PLM, proposed by (Ziems et al., 2022), is shown in the path, $Path_1$. We regard this strategy as our baseline
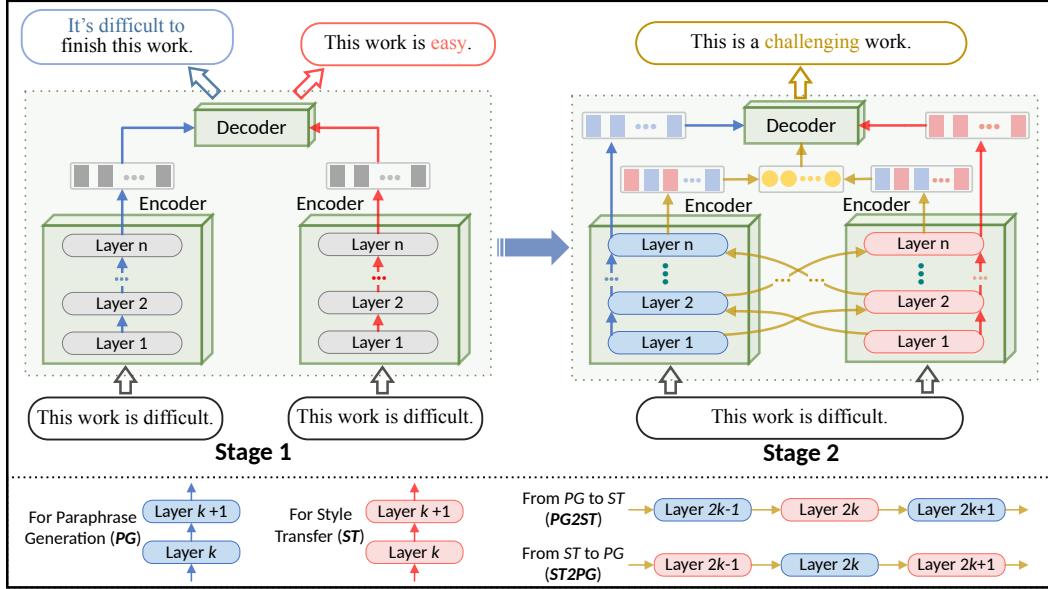
Figure 2: The model architecture and data flow: The architecture consists of two fine-tuning stages, **Stage 1** and **2**, and four data flows, $PG$, $ST$, $PG2ST$, and $ST2PG$.

which requires the model to directly learn the capability of paraphrase generation (PG) with diversity, and sentiment transfer (ST) with a positive perspective for the input. However, it is challenging for the model to directly capture all of the complicated features at once. We thus divide this path into two components to make the problem easier i.e., PG and ST, which are marked with blue and pink colors. Specifically, there are two paths $Path_2$ (from PG to ST) and $Path_3$ (from ST to PG) to obtain the target sentence. The method leverages two pseudo-datasets, paraphrase pairs with sentiment polarities, and sentiment pairs with paraphrases to disentangle meaning and style and transfer the source text into a diverse and positive target sentence. The contributions of this paper can be summarized: (1) we propose a simple but effective reframing model for the PTR task, (2) we propose two data augmentation strategies for generating pseudo-positive reframing datasets, and (3) The experimental results show that our approach improves the performance compared with the baseline on PPF dataset.

## 2 Methodology

### 2.1 Creating Pseudo Data as Prior Knowledge

#### (1) Selecting Annotation Pairs

We choose Microsoft Common Objects in COntext (MSCOCO) which are widely used to learn the paraphrase generation model. We call the data $D_{pg}$. Shen et al. modified the huge Yelp reviews dataset for sentence-level sentiment analysis (Shen et al., 2017). We utilized it to learn the sentiment transfer model. We divided it into two sets, $S_{neg}$ and $S_{pos}$ consisting of sentences with negative and positive sentiment labels, respectively. We created pairs for $\forall s_i \in S_{neg}$, and $\forall s'_i \in S_{pos}$. To reduce the computation cost, for a given $s_i$, we randomly chose the number of $0.05 \times |S_{pos}|$ samples from the set $S_{pos}$. We thus obtained a set $D_{st}$ consisting of $0.05 \times |S_{pos}| \times |S_{neg}|$ sentence pairs.

#### (2) Filtering and Creating Two Pseudo Datasets

To create pseudo datasets from two datasets, $D_{pq}$ and $D_{st}$, each sentence of a pair extracted from $D_{pq}$ should be different polarity from each other. Similarly, each sentence of a pair from $D_{st}$ should be a similar meaning. To this end, a semantic similarity classifier $F_{sem}$ and a sentiment classifier $F_{senti}$ are trained by leveraging BERT (Devlin et al., 2019). We utilized Semantic Textual Similarity Benchmark (STSB) (Cer et al., 2017) and TweetEval Sentiment (TE-sentiment) (Barbieri et al., 2020) for training $F_{sem}$, and $F_{senti}$, respectively. The semantic similarity score obtained by $F_{sem}$ ranges from 0 to 5.0. The higher the score value, the more semantically similar the two sentences are. We chose sentence pairs from the set $D_{st}$ whose similarity score $\delta$ is larger than a certain threshold value and obtained pseudo set $D'_{st}$. Likewise, we chose only two types of sentence pairs labeled with the fine-grained sentiment classifier $F_{senti}$, i.e., (Negative,

425

| Dataset | Train | Validation | Test |
|---------|-------|------------|------|
| PPF | 6,679 | 835 | 835 |
| $D'_{pg}$ | 15,181 | 134 | 1,899 |
| $D'_{st}$ | 14,807 | 139 | 215 |
| STSB | 5,749 | 1,500 | 1,379 |
| TE-sentiment | 45,615 | 2,000 | 12,284 |

Table 1: The statistics of dataset

Neutral) and (Neutral, Positive) from the set $D_{pg}$, resulting in pseudo set $D'_{pg}$.

## 2.2 Fusion Strategies

We recall that the straightforward fine-tuning of PLM illustrated in the path $Path_1$ of Figure 1 requires directly capturing all of the complicated features at once. We thus divide this path into two relative steps: paraphrase generation, and sentiment transfer. The model architecture and data flow are illustrated in Figure 2. It consists of two fine-tuning stages and four data flows. More specifically, in stage 1, the PLM encoder is copied and each encoder is fine-tuned for PG and ST, respectively. We utilize the multi-task learning algorithm proposed by Liu et al. (Liu et al., 2019) to fine-tune the PLM on two pseudo datasets, $D'_{pg}$ and $D'_{st}$. It can balance the $PG$ and $ST$. After processing stage 1, the same model is further fine-tuned on PPF dataset following four variants of data flows, $PG$, $ST$, $PG2ST$, and $ST2PG$. In stage 2, when the data flow is $PG$ or $ST$, each independent encoder is utilized to fine-tune the model, while $PG2ST$ or $ST2PG$, both encoders are utilized. Let $E_{pg} = [l_{p_1}, \cdots, l_{p_n}]$, and $E_{st} = [l_{s_1}, \cdots, l_{s_n}]$ be the encoder for $PG$, and $ST$, respectively. Here, $l_{p_i} \in E_{pg}$ and $l_{s_i} \in E_{st}$ are the $i$-th block layer in the encoder ("Layer $i$" of blue, and pink color in Figure 2, respectively). The encoder by $PG2ST$ and $ST2PG$ flows are shown as $E_{pg2st} = [l_{p_1}, l_{s_2}, \cdots, l_{p_{n-1}}, l_{s_n}]$, and $E_{st2pg} = [l_{s_1}, l_{p_2}, \cdots, l_{s_{n-1}}, l_{p_n}]$, respectively.

## 3 Experiment

## 3.1 Experimental Setting

We chose BART (Lewis et al., 2020) and T5 (Raffel et al., 2020) pre-trained model as the PLM in our method (Lewis et al., 2020) since Ziems et al. (Ziems et al., 2022) reported that they provided the best quality of positive reframes among other PLMs such as GPT-2 (Radford et al., 2019)

and CopyNMT (See et al., 2017). We utilized the version "facebook/bart-base", and "t5-base" on Hugging Face[1] as the backbones. The statistics of datasets are summarised in Table 1. Semantic Textual Similarity Benchmark (STSB) (Cer et al., 2017) and TweetEval Sentiment (TE-sentiment) (Barbieri et al., 2020) are used to train the classifiers, $F_{sem}$, and $F_{senti}$, respectively.

We utilized the PPF dataset[2] to evaluate our method. It consists of 8,349 sentence pairs with manual annotation. The same BART trained in stage 1 is further trained on the PPF training set. The semantic similarity value $\delta$ is set to 3. We tuned the hyperparameters as follows: the batch size is 4, 8, 16, 32, the number of epochs is from 2 to 5, the number of layers $n$ is 12, and the value of the learning rate is from 1e-5 to 1e-4. The procedure of tuning hyperparameters is automatically conducted by the "Ray Tune"[3] library.

For a fair comparison with the baseline by (Ziems et al., 2022), we used the eight metrics, which are (1) ROUGE-1, -2, -LCS (longest common subsequence) (Lin, 2004), BLUE (Papineni et al., 2002) and BERT-Score (Zhang et al., 2019) referring to the gold reference for assessing the performance on content preservation, (2) The $\Delta$TextBlob value (Loria, 2018) for assessing positivity transfer effectiveness, and (3) The Average Length and Perplexity (Yang et al., 2018), followed by (Jin et al., 2022) for measuring the fluency of the output sentences.

## 3.2 Results

Table 2 shows the results on the PPF test dataset. We can see from Table 2 that the results obtained by our approach improve the performance compared with the baseline with the BART model except for BScore and $\Delta$TB. Similarly, our results are better than the baseline with the T5 model except for Avg.Len. This shows that our approach contributes to giving a positive perspective while preserving the original contents. Our variants show that the BART is more effective than T5 by five metrics, R-1, 2, LCS, BLUE, and Avg.Len. However, the variants with T5 are more robust as they work well on content preservation (BScore), positivity transfer ($\Delta$TB), and fluency (PPL).

The performance on the baseline by Avg.Len is more affected by the PLM model than our models

---

| | Method | R-1 | R-2 | R-LCS | BLEU | BScore | ΔTB | Avg.Len | PPL |
|---|---|---|---|---|---|---|---|---|---|
| | (Ziems et al., 2022) | 27.7 | 10.8 | 24.3 | 10.3 | **89.3** | **0.23** | 24.4 | - |
| | ST (ours) | 32.5 | 13.4 | 26.6 | 10.1 | 88.4 | 0.22 | **26.9** | **24.6** |
| BART | PG (ours) | 32.8 | **13.7** | **27.1** | 10.6 | 88.3 | 0.17 | 26.8 | 26.6 |
| | PG2ST (ours) | 32.6 | 13.5 | 26.9 | 10.3 | 88.4 | 0.19 | 26.7 | 24.8 |
| | ST2PG (ours) | **32.9** | 13.6 | **27.1** | **10.9** | 88.4 | 0.20 | 26.6 | 25.6 |
| | (Ziems et al., 2022) | 27.4 | 9.8 | 23.8 | 8.7 | 88.7 | 0.38 | **35.3** | - |
| | ST (ours) | **31.1** | 11.2 | 25.4 | **8.9** | 88.7 | **0.39** | 24.3 | **14.0** |
| T5 | PG (ours) | 30.8 | 11.2 | **25.5** | 8.7 | 88.7 | 0.33 | 23.5 | 15.4 |
| | PG2ST (ours) | **31.1** | 11.2 | **25.5** | **8.9** | 88.7 | 0.35 | 23.4 | 14.5 |
| | ST2PG (ours) | 30.8 | **11.3** | 25.5 | 8.8 | 88.7 | 0.33 | 23.0 | 15.1 |

Table 2: Main results Against the baseline (Ziems et al., 2022) on PPF dataset. ST and PG are the results obtained by only applying stage 1. R-1, R-2, and R-L refer to ROUGE-1, 2, and LCS. BSocre indicates BERT-Score and Avg.Len shows the Average length. The bold font indicates the best result obtained by each backbone.

as there is a significant difference (35.3-24.4) between T5 and BART baselines. Overall, $PG2ST$ and $ST2PG$ except for Avg.Len of T5, preserve the balance between the meaning of the contents and positivity as these results have medium scores between $ST$ and $PG$.

Note that in the BART backbone, the results by the $PG$ strategy are best on all ROUGE scores, while the $ST$ strategy can perform best on average length and perplexity. The reason could be that for $PG$, the encoder is fine-tuned on $D'_{pg}$ which is obtained from paraphrase generation data during the first stage in Figure 2. In contrast, the encoder used by $ST$ is fine-tuned on $D'_{st}$ whose source is sentiment data. Therefore, the model can perform better in terms of preserving the semantic features and sentiment transfer in $PG$, and $ST$, respectively. The ST2PG could balance the functions of $ST$ and $PG$ and obtain the best result on the BLEU score. Why the $PG2ST$ can not perform similarly to $ST2PG$ is still unknown and needs further investigation as future work.

For the backbone of T5, although the results obtained by our four strategies are better than those of the baseline except for the average length, the best performances by each metric are varied on all of these four variants without clear rules. We also need further investigation to make the reasons clear. To conclude our results, our two steps of fine-tuning combined with each strategy can provide different advantages on semantic preserving, sentiment transfer, and the balance of these two.

Table 3 illustrates example sentences obtained by one of the variants of our approach, $ST2PG$, and the baseline with BART. As shown in the se-

quences highlighted in blue and pink, the output sentences generated by our model express more positively than compared with the baseline, properly preserving the meaning of the given input. For instance, in sentence 1, "hope" is a more positive expression and the rest part keeps the meaning and topic of the input sentence. In contrast, the output of the baseline is duplicated with the source input and appends a short sequence including three pronouns "it" which causes ambiguity, although the "it's worth" is positive. In Sentence 2, the baseline blindly copies the counterpart from the input, and just changes the word, "Its" to "It's". The output does not give positive perspectives, while it preserves the meaning of the original sentence. In contrast, "is going through a rough patch" which is generated by our model $ST2PG$ shows that it keeps the bad situation of the stock market, and proposes a positive wish by using the sentence, "I'm sure it will get better".

Likewise, in sentences 3 and 4, the output from our $ST2PG$ can balance the content preserving and sentiment transfer, while the baseline only captures either content preserving or sentiment transfer. For the special case, in sentence 5, "I'm sure I can get through it" in the baseline is positive, while its meaning is too abstract and the context is not suitable. In contrast, the output of our $ST2PG$, "I'll feel better in the morning", is very similar to "get up fresh in the morning." in the ground truth.

## 4 Conclusion

In this paper, we proposed a method for PTR that learns disentangled meaning and style representations. The method utilizes pseudo-positive refram-

| | |
|---|---|
| Sentence 1 | Why is it so hard to get one day of alone time, of peace and quiet??? |
| Ground truth | I could really do with some peace and quiet for a day. |
| $ST2PG$ | I hope I get one day of alone time, of peace and quiet . |
| Baseline | It's hard to get one day of alone time, of peace and quiet , but it's worth it when you get it. |
| Sentence 2 | Its like Friday righ now with the stock market like wtf is going on Lmaoo. |
| Ground truth | It may be a confusing time, but I know that I can remain calm in the storm. |
| $ST2PG$ | The stock market is going through a rough patch right now, but I'm sure it will get better . |
| Baseline | It's like Friday righ now with the stock market like wtf is going on Lmaoo. |
| Sentence 3 | What's the opposite of sleeping like a baby? Because that's me as of late tired. |
| Ground truth | I have not been sleeping like a baby lately, but I will find a way to not be so tired. |
| $ST2PG$ | I'm tired , but I'm sure I'll get better soon . |
| Baseline | I hope I get a good night's sleep so that I can get into a better frame of mind tomorrow. |
| Sentence 4 | I WISH YOU WOULD JUST CHILL OUT! I am working as hard as I can. |
| Ground truth | I am doing the best I can rifht now. |
| $ST2PG$ | I am working as hard as I can , but I know I can do it . |
| Baseline | I am working as hard as I can to make you happy. |
| Sentence 5 | the reason why I'm wide awake right now is unexplainable and is gonna hurt me bad in the morning. |
| Ground truth | I don't know why I am not able to sleep, but I hope to get up fresh in the morning. |
| $ST2PG$ | I'm wide awake right now , but I'm sure I'll feel better in the morning . |
| Baseline | The reason why I'm wide awake right now is unexplainable and is going to hurt me bad in the morning , but I'm sure I can get through it. |

Table 3: Examples from the test set of PPF, their ground truth, and generated positive reframe by $ST2PG$ and Baseline with BART. The sequences marked with blue preserve the original contents, and those of pink give a positive perspective.

ing datasets, and a simple multi-task learning-based model is learned to fuse the generation capabilities from these datasets. The experimental results showed that our approach is comparable with the baseline, especially, since we found that it is effective for generating fluent sentences. Future work will include: (1) exploring more effective augmentation strategies by leveraging a huge number of unlabeled datasets, (2) evaluating our method by using other PPF datasets, and (3) applying it to other TST tasks such as formality and politeness.

## Acknowledgement

## References

Francesco Barbieri, Jose Camacho-Collados, Luis Espinosa Anke, and Leonardo Neves. 2020. TweetEval: Unified benchmark and comparative evaluation for tweet classification. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1644–1650.

Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. 2017. SemEval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 1–14.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of

deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.

Zhenxin Fu, Xiaoye Tan, Nanyun Peng, Dongyan Zhao, and Rui Yan. 2018. Style transfer in text: Exploration and evaluation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1):663–670.

Jingxuan Han, Quan Wang, Licheng Zhang, Weidong Chen, Yan Song, and Zhendong Mao. 2023. Text style transfer with contrastive transfer pattern mining. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7914–7927. Association for Computational Linguistics.

Eduard Hovy. 1987. Generating natural language under pragmatic constraints. *Journal of Pragmatics*, 11(6):689–719.

Di Jin, Zhijing Jin, Zhiting Hu, Olga Vechtomova, and Rada Mihalcea. 2022. Deep learning for text style transfer: A survey. *Computational Linguistics*, 48(1):155–205.

Vineet John, Lili Mou, Hareesh Bahuleyan, and Olga Vechtomova. 2019. Disentangled representation learning for non-parallel text style transfer. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 424–434.

Huiyuan Lai, Antonio Toral, and Malvina Nissim. 2021. Thank you BART! rewarding pre-trained models improves formality style transfer. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 484–494.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81.

Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. 2019. Multi-task deep neural networks for natural language understanding. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4487–4496.

Steven Loria. 2018. textblob documentation. *Release 0.16*, 2.

David D. McDonald and James D. Pustejovsky. 1985. A computational theory of prose style for natural language generation. In *Second Conference of the European Chapter of the Association for Computational Linguistics*, pages 187–193.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. In *OpenAI blog 1(8):9*.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67.

Sudha Rao and Joel Tetreault. 2018. Dear sir or madam, may I introduce the GYAFC dataset: Corpus, benchmarks and metrics for formality style transfer. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 129–140.

Abigail See, Peter J. Liu, and Christopher D. Manning. 2017. Get to the point: Summarization with pointer-generator networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1073–1083. Association for Computational Linguistics.

Tianxiao Shen, Tao Lei, Regina Barzilay, and Tommi Jaakkola. 2017. Style transfer from non-parallel text by cross-alignment. In *Advances in Neural Information Processing Systems*, volume 30.

Daimeng Wei, Zhanglin Wu, Hengchao Shang, Zongyao Li, Minghan Wang, Jiaxin Guo, Xiaoyu Chen, Zhengzhe Yu, and Hao Yang. 2023. Text style transfer back-translation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7944–7959. Association for Computational Linguistics.

Ruochen Xu, Tao Ge, and Furu Wei. 2019. Formality style transfer with hybrid textual annotations. *arXiv preprint arXiv:1903.06353*.

Zichao Yang, Zhiting Hu, Chris Dyer, Eric P Xing, and Taylor Berg-Kirkpatrick. 2018. Unsupervised text style transfer using language models as discriminators. In *Advances in Neural Information Processing Systems*, volume 31.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.

Yi Zhang, Tao Ge, and Xu Sun. 2020. Parallel data augmentation for formality style transfer. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3221–3228.

Caleb Ziems, Minzhi Li, Anthony Zhang, and Diyi Yang. 2022. Inducing positive perspectives with text reframing. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3682–3700.

# Generating clickbait spoilers with an ensemble of large language models

**Mateusz Woźny**[1]  and  **Mateusz Lango**[1,2]
[1] Poznan University of Technology, Faculty of Computing and Telecommunications, Poznan, Poland
[2] Charles University, Faculty of Mathematics and Physics, Prague, Czech Republic
`mateusz.wozny@student.put.edu.pl`, `mlango@cs.put.edu.pl`

## Abstract

Clickbait posts are a widespread problem in the webspace. The generation of spoilers, i.e. short texts that neutralize clickbait by providing information that satisfies the curiosity induced by it, is one of the proposed solutions to the problem. Current state-of-the-art methods are based on passage retrieval or question answering approaches and are limited to generating spoilers only in the form of a phrase or a passage. In this work, we propose an ensemble of fine-tuned large language models for clickbait spoiler generation. Our approach is not limited to phrase or passage spoilers, but is also able to generate multipart spoilers that refer to several non-consecutive parts of text. Experimental evaluation demonstrates that the proposed ensemble model outperforms the baselines in terms of BLEU, METEOR and BERTScore metrics.

## 1 Introduction and related works

Catchy headlines or social media posts designed to entice users to click, known as clickbait, are widespread on the internet. Although they often increase website traffic and generate revenue, they usually fall short of readers' expectations, wasting their time and causing disappointment (Molyneux and Coddington, 2020).

To deal with this problem, Rubin et al. (2015) proposed the clickbait detection task, which received some research attention (Potthast et al., 2016; Chakraborty et al., 2016). More recently, Hagen et al. (2022a) found that clickbaits can be neutralized by providing short texts that clarify what the reader can expect from the linked article, often making the clickbait uninteresting. They developed corpora of clickbait spoilers and classified them into three types: phrase spoilers (containing a single word or a short phrase), passage spoilers (a few sentences at most), and multi-part spoilers

| type | clickbait | spoiler |
|---|---|---|
| phrase | You're missing this major way to save money | promotional code |
| passage | Scientists unearth big surprise near celebrated pyramids | remains of a bustling port and barracks for sailors or troops |
| multi | This is what RE-ALLY happens when you don't brush your teeth | Bad breath, Coronary heart disease, Bleeding gums, (...) |

Table 1: Abbreviated examples of spoilers and clickbaits from Webis-Clickbait-22 corpus (Hagen et al., 2022a)

(containing many non-consecutive phrases and/or passages). See the examples in Tab. 1.

Hagen et al. (2022a) also experimented with 20 approaches for clickbait spoiler generation, which were based on passage retrieval or extractive question answering algorithms. However, all of these methods were only evaluated on phrase and passage spoilers, as they are not suitable for generating multi-part spoilers and their generation requires a specialised approach.

In this work, we demonstrate that all three types of spoilers can be effectively generated by means of conditional language generation with large language models. We put forward a simple yet effective proposal of an ensemble of LLMs that selects the final spoiler by exploiting learning-to-rank techniques. Finally, we verify the performance of the proposed approach and investigate the possibility of combining it with previously developed methods that provide phrase and passage spoilers.

## 2 Ensemble of LLMs for clickbait spoiler generation

The task of clickbait spoiler generation is defined as follows. For a given clickbait text $c$, the content of the linked article $a$ and the requested spoiler type $t$, generate a textual spoiler $s$ whose goal is to make the clickbait $c$ uninteresting for the user by providing the additional information from the referred article $a$. The possible types of spoilers $t$ are phrase, passage, and multi-part.

In this paper, we propose an ensemble of language models for clickbait spoiler generation. The proposed approach consists of three steps: converting the text of a clickbait $c$ into a question $q$, generating candidate spoilers from various prompted large language models, and finally selecting the final spoiler by a trained scoring model.

### 2.1 Converting clickbaits to questions

Clickbaits usually take the form of declarative or exclamatory sentences. In contrast, question answering, which is one of the most related tasks according to related works, naturally deals with problems structured as interrogative sentences. Due to their popularity, QA datasets are often used as a part of LLMs' (pre)training sets, enabling better knowledge transfer for these tasks. Therefore, to better exploit knowledge acquired by LLMs during pretraining, we convert each clickbait into a question before passing it for further processing.

The conversion is made in a zero-shot fashion using the recent Vicuna language model (Chiang et al., 2023). For each clickbait $c$, we construct the following prompt: "Below is a sentence from which write a question.\n Sentence: $c$ \n Question:", where \n is the sign of a new line. The resulting question $q$ is generated by initializing the language model with the prompt and completing the text with the greedy search algorithm until the sign of a new line is generated.

### 2.2 Generating spoilers with LLMs

The next step of our approach is to use a set of different pretrained language models to produce a diversified set of spoiler candidates.

Each component of our ensemble is fine-tuned on the standard language modeling task using an adapter-based approach LoRA (Hu et al., 2021). Such transfer learning approaches allow parameter-efficient fine-tuning by leaving all the pretrained weights unchanged and modifying the model op-

erations by adding shallow, trainable feed-forward networks between the transformer layers. The results of these feed-forward networks are incorporated into the transformer architecture by adding their output to the output of successive transformer layers. Such fine-tuning approaches have proven to be well-suited for relatively small supervised datasets like ours (Houlsby et al., 2019).

In order to create training corpora for the language modeling task, each training example was converted to a textual form by filling in the hand-designed prompt template: "Below is a question paired with a context for which you should generate an answer. Write an answer with type $t$ that appropriately completes the question.\n Question: $q$ \n Context: $a$ \n Answer: $s$\n ". During training, cross-entropy loss was optimized, i.e.

$$-\sum_{i=1}^{n} \log P(w_i|w_1, w_2, ..., w_{i-1})$$

where $w_i$ is the $i$-th token of the filled template and $n$ is its length. During testing, standard greedy decoding was used to retrieve the clickbait spoiler.

### 2.3 Selecting the final spoiler with scoring model

After generating several clickbait spoiler candidates, the final step of the approach is to select the most appropriate spoiler by running a trainable model that evaluates them. This problem can be viewed as a learning-to-rank problem (Liu, 2009), where our goal is to construct the ranking of spoiler candidates and later select the best candidate i.e. the spoiler at the top of the ranking. We experiment with two popular learning-to-rank approaches: 1) a pointwise approach, which assigns a score to each candidate and later uses it to sort the list of candidates. 2) a pairwise approach, which compares all pairs of candidates and decides which spoiler from the pair is more suitable.

#### 2.3.1 Pointwise approach

To evaluate each candidate, we develop a regressor that tries to predict the value of BLEU score for each spoiler. As a regressor fine-tuned DeBERTa model (He et al., 2021) with one linear layer on top of CLS token is used. The input to the model consists of a question (clickbait) $q$, candidate spoiler $s$, and article $a$, separated by the sign of a new line and concatenated into one input text. The output of the model is the predicted BLEU score.

The training data for the regressor was generated by running all the LLMs used in the ensemble on the training data and evaluating their BLEU score against the available gold standard. During training, the classical sum of squared errors was optimized.

### 2.3.2 Pairwise approach

The second method considered for selecting the best spoiler among the candidates is the pairwise approach. This approach relies on a classifier that, for a given pair of spoiler candidates, decides which of them is more suitable. More specifically, the BERT-based classifier receives the same input as in the pointwise approach, but with two spoiler candidates $s_1$ and $s_2$. The output of the binary classification model is 1 if $s_1$ is better than $s_2$ in terms of BLEU, and 0 otherwise.

The classifier is trained on the generated data as follows. First, all LLMs used in the ensemble were run on the training set, generating a collection of spoiler candidates for each clickbait. Later, all possible pairs from each collection were considered and converted into binary classification instances by comparing the BLEU scores of the candidates. The pairs containing spoilers with BLEU= 0 or pairs containing identical spoilers were removed from the training data. During training, the standard cross-entropy loss was optimised.

## 3 Experiments

We have performed computational experiments aimed at verifying the effectiveness of ensembling LLMs with pointwise and pairwise rankers, and comparing its effectiveness with the previous SOTA methods based on question answering. In addition, we also investigate the possibility of combining the previous QA approaches, which are best suited for passage and phrase spoiler types, with the proposed approach for multi-part spoilers.

### 3.1 Experimental setup

We experiment with an ensemble of three fine-tuned LLMs, which were constructed from two pretrained models: LLaMA (Touvron et al., 2023) and Vicuna (Chiang et al., 2023). Both of these models are open-source and were fine-tuned using the prompt described in Sec. 2.2. However, we observed additional improvements with Vicuna model while using specially tailored prompts for each spoiler type separately (see appendix), therefore we also report the results of this fine-tuned model and use it as a part of the ensemble.

The obtained results were compared against the performance of two extractive QA approaches, which on top of pretrained encoder perform begin/end span classification[1]. These approaches are based on RoBERTa (Liu et al., 2019) and De-BERTa (He et al., 2021) models since among 20 approaches compared on clickbait spoiler generation task by Hagen et al. (2022a) these two were the most effective ones. Note, that the results reported for these approaches in this work are significantly lower than therein, since we report the performance over all three types of spoilers, including multipart.

The ensemble approach with the pointwise ranker used DeBERTa-based regressor, which obtained MSE of 0.384 on test set. Similarly, the classifier used in the pairwise approach achieved balanced accuracy of 90,8% on the test data.

Following earlier works, we evaluated the approaches with three metrics: BLEU (Papineni et al., 2002), METEOR (Banerjee and Lavie, 2005), and BERT Score (Zhang et al., 2020). All experiments were performed on a single Nvidia A100 GPU.

Some additional experiment details and results can be found in the online appendix[2].

### 3.2 Comparing the proposed approaches with related works

The results of QA-based approaches, ensemble models as well as individual fine-tuned LLMs can be found in Table 2. The best-performing approach according to BLEU, METEOR, and BERT Score F1 is the proposed ensemble with a pointwise ranker. This ensemble provides the improvement of approx. 2 percentage points in terms of BLEU and METEOR over the best of its components i.e. Vicuna model with specific prompts for each spoiler type. The second-best approach was the ensemble with pairwise ranker which offered limited improvement over the individual LLMs.

Overall, each of the approaches using LLMs obtained better results than previous state-of-the-art approaches based on extractive question answering. The only metric on which the QA-based approaches (RoBERTa and DeBERTa) stand out is BERTScore Precision. Still, BERTScore Recall is higher for LLM and ensemble approaches, making them more effective in terms of F1 measure, which combines both precision and recall.

---

[1]This is the default fine-tuning approach for QA-task of BERT (Devlin et al., 2019), more details therein.

[2]https://www.cs.put.poznan.pl/mlango/publications/inlg23.pdf

| Model | BLEU | METEOR | BERT Score | | |
| | | | Precision | Recall | F1 |
|---|---|---|---|---|---|
| RoBERTa | 31,78 | 0,387 | 0,904 | 0,883 | 0,893 |
| DeBERTa | 32,20 | 0,398 | 0,907 | 0,884 | 0,894 |
| LLaMA 13B | 37,70 | 0,474 | 0,895 | 0,901 | 0,897 |
| Vicuna 13B | 38,80 | 0,481 | 0,898 | 0,903 | 0,900 |
| Vicuna 13B with type-based prompts | 40,02 | 0,492 | 0,899 | 0,905 | 0,901 |
| Ensemble with pairwise ranker | 40,76 | 0,500 | 0,901 | 0,907 | 0,904 |
| Ensemble with pointwise ranker | 42,13 | 0,517 | 0,902 | 0,909 | 0,905 |

Table 2: The experimental results of previous state-of-the-art QA-based methods compared with our ensembling approaches and LLMs. All the metrics are computed on Webis Clickbait 22 corpora (Hagen et al., 2022b).

| Model | BLEU | METEOR | BERT Score | | |
| | | | Precision | Recall | F1 |
|---|---|---|---|---|---|
| DeBERTa | 32,20 | 0,398 | 0,907 | 0,884 | 0,894 |
| DeBERTa trained on questions | 37,82 | 0,451 | 0,913 | 0,895 | 0,903 |
| Baseline ensemble | 42,28 | 0,506 | 0,910 | 0,908 | 0,909 |
| Ensemble with pairwise ranker | 43,57 | 0,520 | 0,912 | 0,911 | 0,911 |
| Ensemble with pointwise ranker | 44,45 | 0,532 | 0,911 | 0,913 | 0,911 |

Table 3: The experimental results of the ensemble that combines previously proposed methods for passage and phrase spoiler types (DeBERTa) with the newly proposed approaches for clickbait spoiler generation.

## 3.3 Combining previous SOTA models with the proposed ones

In the final experiment, we decided to verify whether it is possible to obtain even better results by combining our approaches with QA-based models previously designed for clickbait spoilers of phrase and passage types only.

As a QA model, we use the fine-tuned DeBERTa model, since it gave the best results both in our experiments from the previous section and in the experiments of Hagen et al. (2022a). As we mentioned in Sec. 2.1, our approaches generate spoilers for clickbaits converted into interrogative sentences in order to facilitate better knowledge transfer from the pre-trained models. We also fine-tuned DeBERTa on such preprocessed data and found improved performance. Therefore, this model (later denoted DeBERTa-q or "DeBERTa trained on questions") was used for ensembling.

The operation of ensembles reported in this section slightly differs from what was described in Sec 2.3 to better account for QA-based approaches' suitability for passage and phrase spoiler types. If the generated spoiler is of mutli-part type, the list of candidate spoilers is generated as previously, i.e. from three fine-tuned LLMs. However, if the generated spoiler is of phrase/passage type then only spoiler candidates from DeBERTa-q and Vicuna

with customised prompts are considered. As previously, the selection of the best spoiler among the candidates is performed by a ranker.

The results of these approaches are reported in Table 3. It can be seen that DeBERTa-q achieves significantly better results than DeBERTa for all metrics considered. However, the ensemble with LLMs provides further significant improvements.

As a form of sanity check for our ensemble approach that uses a ranker to select the best spoiler, we have also implemented a trivial ensemble (baseline ensemble) that uses DeBERTa-q for all passage and phrase spoilers and Vicuna with type-based prompts for multi-part spoilers. Although such a form of ensembling is also advantageous in that the performance obtained is better than that of the individual models, using any variant of the proposed ensemble with a ranker still improves the results. For example, for the METEOR measure, the improvement is almost 2% for the pairwise approach and 3% for the pointwise approach.

Table 4 presents a more detailed analysis of the performance of the spoiler generation methods, i.e. the results are reported separately for each spoiler type. Although the newly proposed ensemble approaches achieve better performance for each spoiler type, the most significant improvement is observed for the most difficult multi-part spoilers,

| Model | spoiler type | BLEU | METEOR | BERT Score Recall | F1 | Precision |
|---|---|---|---|---|---|---|
| DeBERTa | phrase | 56,00 | 0,569 | 0,934 | 0,932 | 0,931 |
| | passage | 20,10 | 0,304 | 0,858 | 0,869 | 0,883 |
| | multi-part | 2,07 | 0,204 | 0,822 | 0,860 | 0,902 |
| DeBERTa trained on questions | phrase | 62,50 | 0,627 | 0,946 | 0,942 | 0,939 |
| | passage | 21,60 | 0,319 | 0,861 | 0,872 | 0,885 |
| | multi-part | 4,10 | 0,224 | 0,829 | 0,862 | 0,900 |
| Ensemble with pairwise ranker | phrase | 63,30 | 0,617 | 0,945 | 0,943 | 0,942 |
| | passage | 30,00 | 0,416 | 0,882 | 0,887 | 0,894 |
| | multi-part | 27,30 | 0,525 | 0,894 | 0,886 | 0,879 |
| Ensemble with pointwise ranker | phrase | 65,30 | 0,645 | 0,950 | 0,946 | 0,942 |
| | passage | 30,30 | 0,424 | 0,884 | 0,887 | 0,892 |
| | multi-part | 28,60 | 0,522 | 0,894 | 0,887 | 0,881 |

Table 4: The performance of clickbait spoiler generation models for each spoiler type.

| Clickbait | Reference | DeBERTa | Vicuna |
|---|---|---|---|
| Agency might plant a garden on the moon. | NASA | *NASA* | lunar sunlight |
| 10 habits of incredibly happy people | 1 they slow down to appreciate life's little pleasure 2 they exercise 3 they spend money on other people (...) | they have a growth mindset | *1 they slow down to appreciate life's little pleasure 2 they exercise 3 they spend money on other people (...)* |

Table 5: Two examples of spoilers generated by different methods. The spoilers in italics were selected by pointwise ranker (The ensemble also includes the LLaMA model, the output of which is not shown due to page limits.).

which are the focus of our paper. For example, the best ensemble model achieves an improvement of over 25 in terms of BLEU score over the previous DeBERTa model.

Two examples of generated spoilers are presented in Table 5 (more examples in the appendix). It can be observed that the previous SOTA approach (DeBERTa) fails to extract the correct spoiler of the multi-part type, but Vicuna model generates it correctly. In contrast, DeBERTa extracts the correct phrase spoiler and Vicuna fails to do so. In both cases, the pointwise ranker was able to indicate the correct spoiler.

# 4 Summary

In this paper, we have shown that using fine-tuned LLMs can be a simple, yet effective way of dealing with clickbait spoiler generation for all three spoiler types considered, i.e. including multi-part spoilers. We also demonstrated that ensembling several such models with a ranker that selects the most suitable spoiler leads to improved results, es-

pecially when using a pointwise ranker. Finally, we show that combining state-of-the-art approaches for phrase and passage spoiler types based on question answering with the newly proposed ones based on LLMs leads to further improvements.

# References

Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with im-

---

[3] https://github.com/mateusz-wozny/spoiler_generation

proved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.

Abhijnan Chakraborty, Bhargavi Paranjape, Sourya Kakarla, and Niloy Ganguly. 2016. Stop clickbait: Detecting and preventing clickbaits in online news media. *CoRR*, abs/1610.09786.

Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Matthias Hagen, Maik Fröbe, Artur Jurk, and Martin Potthast. 2022a. Clickbait Spoiling via Question Answering and Passage Retrieval. In *60th Annual Meeting of the Association for Computational Linguistics (ACL 2022)*, pages 7025–7036. Association for Computational Linguistics.

Matthias Hagen, Maik Fröbe, Artur Jurk, and Martin Potthast. 2022b. Webis clickbait spoiling corpus 2022.

Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. Deberta: Decoding-enhanced bert with disentangled attention. In *International Conference on Learning Representations*.

Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for NLP. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 2790–2799. PMLR.

Edward Hu, Yelong Shen, Phil Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models.

Tie-Yan Liu. 2009. Learning to rank for information retrieval. *Found. Trends Inf. Retr.*, 3(3):225–331.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized BERT pretraining approach. *CoRR*, abs/1907.11692.

Logan Molyneux and Mark Coddington. 2020. Aggregation, clickbait and their effect on perceptions of journalistic credibility and quality. *Journalism Practice*, 14(4):429–446.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. pages 311–318.

Martin Potthast, Sebastian Köpsel, Benno Stein, and Matthias Hagen. 2016. Clickbait detection. In *European Conference on Information Retrieval*.

Victoria Rubin, Nadia Conroy, and Yimin Chen. 2015. Towards news verification: Deception detection methods for news discourse.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

# Reducing named entity hallucination risk to ensure faithful summary generation

**Eunice Akani**[1,2] and **Benoit Favre**[1] and **Frederic Bechet**[1] and **Romain Gemignani**[2]

[1] Aix-Marseille Univ, CNRS, LIS, Marseille, France

[2] Enedis, Marseille, France

`firstname.lastname@lis-lab.fr`

## Abstract

The faithfulness of abstractive text summarization at the named entities level is the focus of this study. We propose to add a new criterion to the summary selection method based on the "*risk*" of generating entities that do not belong to the source document. This method is based on the assumption that Out-Of-Document entities are more likely to be *hallucinations*. This assumption was verified by a manual annotation of the entities occurring in a set of generated summaries on the CNN/DM corpus. This study showed that only 29% of the entities outside the source document were inferrable by the annotators, leading to 71% of hallucinations among OOD entities. We test our selection method on the CNN/DM corpus and show that it significantly reduces the hallucination risk on named entities while maintaining competitive results with respect to automatic evaluation metrics like ROUGE.

## 1 Introduction

Abstractive text summarization methods aim at generating shorter versions of documents, possibly containing new words with respect to original documents. Recent pre-trained language models (Devlin et al., 2019; Lewis et al., 2020; Zhang et al., 2019; Raffel et al., 2020) allow to obtain fluent generated text. Despite their remarkable performance, those models tend to generate information that cannot be inferred from the source document. According to a study by Cao et al. (2018), 30% of the summaries generated by various systems have inconsistencies, qualified as "hallucination" by (Maynez et al., 2020). Current metrics used to assess automatic text summarization systems, such as ROUGE (Lin, 2004) and BERTScore (Zhang* et al., 2020), do not account for these issues. These hallucinations led to several studies on evaluating the faithfulness of generated summaries and generating more faithful texts (Li et al., 2022; Ji et al., 2023). For example Durmus et al., 2020; Deutsch

et al., 2021 proposed a QA-based metric. They produced questions from the generated summary and check that the answer in the document matches the one in the summary. Natural Language Inference (NLI) have also been proposed to evaluate the factuality of a generated summary by checking if it entailed the source document (Falke et al., 2019; Maynez et al., 2020; Laban et al., 2022). Most recent studies (Luo et al., 2023) used ChatGPT[1] as NLI system to evaluate the factual inconsistency of generated summary. Chen et al., 2021 suggested using contrast candidate generation and selection as a post-processing method to avoid hallucination. They create candidate summaries and replace named entities with ones found in the document to summarize. Nan et al., 2021 introduced an entity-based metric to compare the entities in the gold summary to the generated summary. They also proposed a method based on training data filtering and multi-task learning to avoid hallucinations.

This paper studies hallucinations at the entity level, which is a crucial level (Chen et al., 2021) in news-related corpora such as CNN/DM (Hermann et al., 2015) or XSum (Narayan et al., 2018). We aim to reduce the "*risk*" of generating nonfactual summaries by reducing the number of possibly hallucinated entities using a new criterion called *Named Entity Hallucination Risk* (*NEHR*). This criterion stands on the precision-source of Nan et al., 2021 for summary selection. Based on the assumption that an Out-Of-Document (OOD) entities have more chance to be hallucinations, we empirically verify that the entities in the summary are occurring in the source document.

Our contributions are as follows :

- We create summaries using sampling methods and we propose a selection criterion called *Named Entity Hallucination Risk* (NEHR) to minimize factual errors.

---

[1] https://openai.com/blog/chatgpt

437

- We empirically examine the relevance of NEHR as a selection criterion.

- We evaluate summaries chosen with NEHR and the model's performance on two benchmark corpora CNN/DM and XSum.

## 2   Named Entity Hallucination Risk

We propose to assess the risk of using entities and quantities incorrectly, leading to hallucinations in generated summaries. In the following, we call an *entity* a word string belonging to a named entity category such as person, location, organisation and quantities. An entity can be considered as an *in-document* entity if its corresponding word string can be found in the source document (*exact-match in-document entities)*, or if it relates to an alternative form (*e.g. New York/Big Apple*). If an entity is not mentioned in the source document, it is considered as an Out-Of-Document (OOD) entity. OOD entites can be either acceptable if they can be inferred from the document (*e.g. New York → USA*) or can be considered as *hallucinations*.

We call *hallucinations* the occurrence in generated summaries of entities that cannot be inferred from the source document. Hallucinations can be obvious errors, entities that have nothing to do with the source document context, or entities that could be acceptable by looking at external sources but which were considered as dubious by human annotators: accepting them would require some form of world knowledge beyond general understanding, not directly available from the document. Let's point out that an entity can be an in-document entity or an acceptable OOD entity but still be incorrect in its context of occurrence in a generated summary.

This study aims to reduce the risk of hallucinating entities thanks to a selection criterion called *NEHR* for *Named Entity hallucination Risk* inspired by the *precision-source* metrics proposed by Nan et al., 2021. A simple selection system can be built from our NEHR criterion, which selects less risky hypotheses from a large sample of summaries.

**NEHR definition**   The NEHR criterion is based on the assumption that in-document entities, and moreover exact-match in-document entities are less prone to be incorrectly used in generated summaries than OOD entities. From this assumption, a summary-level hallucination risk, called NEHR,

can be computed as the percentage of entities in that summary not in the source document. For a document $d$ and a summary $s$ we have:

$$NEHR(d, s) = (1 - \frac{|entities \in d \wedge s|}{|entities \in s|}) \times 100 \quad (1)$$

For detecting named entities and quantities in $d$ and $s$ we rely on an automatic NER system as in (Chen et al., 2021). However, we have no direct way to check if an entity considered as *risky* (*i.e.* not in the source document) is correct or not. Therefore, we rely on human annotations to verify that this criterion is effectively correlated with hallucination errors in generated summaries.

**Assessing the relevance of the NEHR criterion** The following steps are involved in the empirical study designed to test the relevance of our risk criterion:

1. Select a document/summary corpus $C$, train several summarization generation models on the training instances of $C$; generate a set $S_d$ of alternative summaries for all documents $d$ belonging to the test partition of $C$ with the different models and compute the ROUGE and NEHR scores for all summaries $s \in S_d$.

2. For each document $d \in C$, select the maximum ROUGE hypothesis $\hat{s}_d$:

$$\hat{s}_d = \operatorname*{argmax}_{s \in S_d} ROUGE(s, s_{ref})$$

3. Run a NER on each summary $e(\hat{s}_d)$.

4. Manually inspect all entities $e$ detected in summaries $\hat{s}_d$ to classify them according to 2 dimensions: inside/outside the source document; correct/incorrect use of $e$ in $\hat{s}_d$.

This study can tell us whether our risk criterion is indeed correlated with the factuality of the summary generated, and whether incorrectly used entities are more frequent outside the document than inside. The motivation for using the max ROUGE summary $\hat{s}_d$ for each document $d$ is to abstract from a given system by using an upper bound of what current state-of-the-art systems can produced. The following section describes the empirical study we did on the CNN/DM corpus.

## 3   Empirical study on the NEHR criterion

**Dataset and model**   We used the test set of the CNN/DM (Hermann et al., 2015), a popular abstractive dataset for the automatic text summary

task that comes from news articles from the CNN and the Daily Mail websites; and BART (Lewis et al., 2020), a transformer encoder-encoder based model that can be fine-tuned to perform automatic text summarization. BART's pre-training consists in applying an arbitrary noising function to corrupt the text deliberately, and training a model to reconstruct the original text from this corrupted version. We initialised the model with pretrained weights from the hugging face library (Wolf et al., 2020) resulting from the fine-tuning of BART-large on CNN/DM[2].

**Sampling a diverse population of summaries** We generated multiple summaries using 4 different sampling methods to select the next token from a language model: *beam search* that keeps the *num_beams* highest probability paths at each step; *Temperature Sampling* that consists in re-scaling *logits* before applying the softmax; *Top-K Sampling* (Fan et al., 2018) that only keeps the $K$ most likely next words and redistributes the probability among those $K$ words and *Top-P Sampling* (Holtzman et al., 2019) that consists in, given a probability $p$, taking the smallest possible set of next words whose cumulative probability exceeds a given probability mass and redistributes the probability among them. By using several sampling parameters, we are able to generate a set of 77 summaries for each document to summarize. Refer to appendix A for sampling parameters.

**Named entity extraction** In this study we use FLERT (Schweter and Akbik, 2020)[3] to extract Named Entities from documents and summaries. FLERT is a NER system that yields 90.93% of F1-score on OntoNotes, a large NER annotated dataset. By extracting NEs on the source documents and on the set of summaries generated by our sampling method on each source document, we are able to estimate the NEHR value for each summary.

Our first evaluation consists of computing the ROUGE metrics on the summary sets generated by our sampling method in order to check that each set contain a large diversity of summaries corresponding to a large range of ROUGE values. We computed ROUGE and NEHR for different summaries produced by the sampling strategy (Table 1). The variations of ROUGE show how the summaries

| summary set | ROUGE (R-1/R-2/R-L) | NEHR |
|---|---|---|
| ROUGE max | 57.45 / 32.59 /41.63 | 4.6 |
| ROUGE min | 30.04 / 09.33 /19.47 | 6.0 |

Table 1: Maximal and minimal values of ROUGE on the summary set produced by our sampling method on CNN/DM. The NEHR value for the min/max summaries is also reported.

| | in-doc | | out-doc | |
|---|---|---|---|---|
| Entity dist (%). | 79.7 | | 20.3 | |
| Type | *exact.* | *var.* | *inf.* | *hall.* |
| Type Dist. (%) | 62.8 | 37.2 | 28.8 | 71.2 |
| % correct | 90 | 90 | 88 | - |

Table 2: % of correctly used entities for each subset of in-document and out of document entities. Dist. is the percentage of entities belonging to each type for entities inside (*exact match, variation*) and outside (*inferable, hallucination*) the document.

generated cover a large space. Large variations of ROUGE do not translate into large variations of NEHR values for the summaries corresponding to the min and max values of the ROUGE criterion, highlighting the fact that ROUGE might not be a good indicator of summary faithfulness.

**Are Out-Of-Document entities mostly hallucinations?** We manually analyzed the generated summaries to identify if our assumption that OOD entities were likely to be hallucinations was correct. We randomly selected 50 generated summaries from the test set of CNN/DM with the following constraint: each generated summary must contain at least one exact-match in-document entities and at least one that had no match (either variation of in-document entities or OOD entities). We used these constraints in order to oversample in-document and OOD entities. By using the *exact-match* constraint we were able to select automatically the in-document summaries. In each summary, we manually inspected the same number of entities with exact matches to entities in the source document as those with no matches, resulting in 145 entities with exact matches and 145 with no matches. Three annotators were asked to label each entity as correctly or incorrectly used w.r.t. the following definition: *the entity is used in the correct context according to the document*. The entities inside and outside the document were divided into two types: *exact match* or *variation* for in-document entities and *inferable* or *hallucination* for OOD entities. While *exact match* refers to entities that match exactly those in the source document, *variation* refers

to entities in the document that are written with misspellings or using alternative forms, such as a date written differently in the document and in the summary. *Inferable* refers to entities outside the document whose correct use and veracity can be inferred as presented in section 2. We collect the annotations and report the results based on a majority vote among the annotators. The inter-annotators agreement Cohen kappa (Cohen, 1960) of at least 0.63 was obtained for each pair (refer to appendix B for kappa). According to the annotation obtained, most of the annotated entities belong to the set of in-document entities (80%); only 20% belong to OOD entities. In table 2, we report the % of correctly used entities for each subset of in-document and OOD entities. For in-document entities, 90% are labelled as "*correct*", and there are no differences between exact-match and alternative forms of entities. When dealing with OOD entities, 71% of them were considered as hallucinations by our annotators. This confirms our assumption that in-document entities are a good source of information for computing NEHR. It is interesting to see that the set of 29% inferable entities have almost the same correctness (88% *v.s.* 90%) than in-document entities. So the drop in correctness between in-document and out-of-document entities exclusively comes from hallucinations, which represent about 71% of the out-of-document entities. Therefore, by explicitly minimizing the risk of having out-of-document entities, we reduce the risk of hallucinations and this can lead to an increase in summary faithfulness. In the next section we apply the criterion to select summaries at inference time and check the impact on summary quality of explicitly minimizing this criterion.

## 4 Experiments on summary selection

We evaluate the impact of using our NEHR criterion for selecting a summary at inference time among a possible set of summaries and compare the results obtained in terms of ROUGE and human evaluation with three baseline selection methods: summary with the highest score given by the model among the summaries generated; the 1-best with $beam = 4$ and the summary with the best entailment score compare to the source document as proposed in (Maynez et al., 2020). We propose a criterion based on both NEHR and model scores: First, the population of possible summaries is reduced to those with the lowest NEHR value, then

the summary with the highest model score is selected from that subgroup. Let $H$ be the set of summaries sampled from the model, $V$ set of summaries with minimum risk, $P(\cdot|model)$ the probability given by the model to a summary, and $\hat{s}$ the final system output:

$$V = \left\{ x \in H | risk(x) = \min_{s' \in H} NEHR(s') \right\} \quad (2)$$

$$\hat{s} = \operatorname*{argmax}_{s \in V} P(s|model) \quad (3)$$

In section 3, we saw that 37.2% of the entities in the document were variations of other entities belonging to the source document. Thus, using heuristics described in appendix C, we identified entities that were possible variations of in-document entities to ensure that the OOD entities were not mistakenly considered to be alternative forms of in-document entities by our automatic system .

**Automatic evaluation** We report into table 3, the ROUGE (Lin, 2004), NEHR (see 1) and the percentage of summary with at least one OOD entity computed on CNN/DM. We also report the results of the same experiment made on XSum (Narayan et al., 2018), a more abstractive text summarization dataset than CNN/DM where the reference summaries contain a large number of OOD ngrams: 36% new 1-grams and 83% new 2-grams in XSum whereas CNN/DM has 17% and 54% respectively (Narayan et al., 2018).
The table shows that if BART-Large with beam=4 yields the best automatic evaluation scores, it results in a higher value of NEHR for both dataset. Our approach minimizes that number while maintaining almost equivalent ROUGE for CNN/DM. For XSum, there is a big loss in term of ROUGE. But our proposed method reduce the NEHR as well as the percentage of summary with at least one entity that could be an hallucination. This may be due to the fact that XSum's references summaries are not faithful to the document. In order to check if our selection criterion did not negatively impact the subjective quality of the chosen summary, we performed a manual evaluation on a subsample of the XSum test set.

**Human evaluation** We evaluated 10 XSum test examples selected randomly. This evaluation consists in annotating the faithfulness of the different entities with respect to the document. Two annotators were asked to label each entity as true if it was

| | | ROUGE-1 ↑ | ROUGE-2 ↑ | ROUGE-L ↑ | NEHR ↓ | %HallSum ↓ |
|---|---|---|---|---|---|---|
| CNN/DM | BEAM 4 | **43.74** | **20.84** | **30.44** | 0.5 | 3.86 |
| | BEST PROBA | 41.99 | 18.96 | 28.01 | 2.6 | 20.57 |
| | ENTAILMENT | 43.61 | 19.69 | 29.26 | 1.62 | 12.92 |
| | MIN NEHR + VAR (OUR) | 42.19 | 19.12 | 28.24 | **0.003** | **0.035** |
| XSUM | BEAM 4 | **45.32** | **22.20** | **37.10** | 27.67 | 52.48 |
| | BEST PROBA | 40.26 | 16.79 | 31.29 | 31.05 | 61.24 |
| | ENTAILMENT | 40.92 | 17.14 | 31.96 | 27.08 | 54.98 |
| | MIN NEHR + VAR (OUR) | 40.16 | 16.54 | 31.31 | **6.92** | **21.49** |

Table 3: Evaluation on CNNDM and XSum. BEST PROBA – the summary with the best score among all generated summaries given by the model. MIN NEHR – our proposed summary selection method after using variation heuristics. %HallSum – the percentage of summary with at least one entity out of the document. ↑ higher is better, ↓ lower is better.

| | Annotator. 1 | Annotator. 2 | Adjudicated |
|---|---|---|---|
| BEAM 4 | 49.33 | 61.16 | 48.67 |
| OUR | 59.76 | 61.07 | 61.67 |

Table 4: Average percentage of well-used entities for each system annotated. In adjudicated results, the annotators had to agree on a label for each entity.

used in the right context with respect to the document. An inter-annotator agreement kappa of about 0.38 was obtained. The disagreement among annotators could be attributed to the challenges they faced while annotating entities in sports-related news items where a lot of external knowledge was required to check the correctness of an OOD entity. We held a post-annotation adjudication phase to get annotators to agree on a label for each entity. We report, in Table 4, the percentage of correct entities for each annotator and the adjudication. Looking at the adjudication, 61% of the entities of our method has been tag as correct while 48% for the beam one. That means our method seems to increase the number of correctly used entities by reducing the hallucination risk.

## 5 Conclusion

We propose in this study a new Named Entity Hallucination Risk criterion for summary selection. Using human evaluation on oracle summaries generated by Bart on CNN/DM, we showed that in-document entities are mostly correct, but this performance drops when considering entities outside documents because of *hallucinations*. We observed empirically that our proposed summary selection method did not significantly impact the ROUGE score for CNN/DM while significantly reducing the hallucination risk. On the highly abstractive dataset XSum, our method was able to drop dramatically the hallucination risk but with a significant drop in ROUGE. Human evaluation of the generated sum-

maries selected from XSum using NEHR showed that the occurring entities were more often correct with respect to those obtained without our selection criteria.

## Acknowledgements

## References

Ziqiang Cao, Furu Wei, Wenjie Li, and Sujian Li. 2018. Faithful to the original: Fact aware neural abstractive summarization. *ArXiv*, abs/1711.04434.

Sihao Chen, Fan Zhang, Kazoo Sone, and Dan Roth. 2021. Improving faithfulness in abstractive summarization with contrast candidate generation and selection. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5935–5941, Online. Association for Computational Linguistics.

Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1):37–46.

Daniel Deutsch, Tania Bedrax-Weiss, and Dan Roth. 2021. Towards question-answering as an automatic metric for evaluating the content quality of a summary. *Transactions of the Association for Computational Linguistics*, 9:774–789.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Esin Durmus, He He, and Mona Diab. 2020. FEQA: A question answering evaluation framework for faithfulness assessment in abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5055–5070, Online. Association for Computational Linguistics.

Tobias Falke, Leonardo F. R. Ribeiro, Prasetya Ajie Utama, Ido Dagan, and Iryna Gurevych. 2019. Ranking generated summaries by correctness: An interesting but challenging application for natural language inference. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2214–2220, Florence, Italy. Association for Computational Linguistics.

Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 889–898, Melbourne, Australia. Association for Computational Linguistics.

Karl Moritz Hermann, Tomáš Kočiský, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration.

Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. 55(12).

Philippe Laban, Tobias Schnabel, Paul N. Bennett, and Marti A. Hearst. 2022. SummaC: Re-visiting NLI-based models for inconsistency detection in summarization. *Transactions of the Association for Computational Linguistics*, 10:163–177.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

Wei Li, Wenhao Wu, Moye Chen, Jiachen Liu, Xinyan Xiao, and Hua Wu. 2022. Faithfulness in natural language generation: A systematic survey of analysis, evaluation and optimization methods. *ArXiv*, abs/2203.05227.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Zheheng Luo, Qianqian Xie, and Sophia Ananiadou. 2023. Chatgpt as a factual inconsistency evaluator for text summarization.

Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1906–1919, Online. Association for Computational Linguistics.

Feng Nan, Ramesh Nallapati, Zhiguo Wang, Cicero Nogueira dos Santos, Henghui Zhu, Dejiao Zhang, Kathleen McKeown, and Bing Xiang. 2021. Entity-level factual consistency of abstractive text summarization. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2727–2733, Online. Association for Computational Linguistics.

Shashi Narayan, Shay B. Cohen, and Mirella Lapata. 2018. Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1797–1807, Brussels, Belgium. Association for Computational Linguistics.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *arXiv:1910.10683 [cs, stat]*. ArXiv: 1910.10683.

Stefan Schweter and Alan Akbik. 2020. Flert: Document-level features for named entity recognition.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.

Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter J. Liu. 2019. Pegasus: Pre-training with extracted gap-sentences for abstractive summarization.

Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

# Building a dual dataset of text- and image-grounded conversations and summarisation in Gàidhlig (Scottish Gaelic)

**David M. Howcroft*     Will Lamb†     Anna Groundwater‡     Dimitra Gkatzia***

*Edinburgh Napier University   {D.Howcroft,D.Gkatzia}@napier.ac.uk
†University of Edinburgh   W.Lamb@ed.ac.uk
‡National Museum of Scotland   a.groundwater@nms.ac.uk

## Abstract

Gàidhlig (Scottish Gaelic; gd) is spoken by about 57k people in Scotland,[1] but remains an under-resourced language with respect to natural language processing in general and natural language generation (NLG) in particular. To address this gap, we developed the first datasets for Scottish Gaelic NLG, collecting both conversational and summarisation data in a single setting. Our task setup involves *dialogues* between a pair of proficient speakers discussing museum exhibits, grounding the conversation in images and texts. Then, each interlocutor summarises the dialogue resulting in a secondary *dialogue summarisation* dataset. This paper presents the dialogue and summarisation corpora, as well as the software used for data collection. The dialogue dataset consists of 43 conversations (13.7k words) and 61 summaries (2.0k words).[2]

## 1   Introduction

The preservation of minority languages and the development of Natural Language Processing (NLP) systems in low-resource settings have gained increasing attention in recent years (Howcroft and Gkatzia, 2022; Castro Ferreira et al., 2020; Zhao et al., 2022), fueled by efforts to safeguard linguistic diversity and cultural heritage (Bella et al., 2020) as well as efforts to create inclusive and fairer systems (Nee et al., 2021; Joshi et al., 2020). One such minority language is Scottish Gaelic, which despite being a recognised national language under the European Charter for Regional or Minority Languages and a rise in Gaelic-medium education, faces challenges in terms of linguistic resources for the development of natural language generation (NLG) systems. To bridge this gap, we present a

novel dialogue and dialogue summarization corpus for Scottish Gaelic, laying the foundation for further developments in NLG in this language.

To construct the dataset, we adopted a task setup centered around dialogues between proficient speakers, engaging them in conversations about museum exhibits. By grounding the conversations in images and texts, we aimed to create a contextually rich conversational dataset. Subsequently, each participant summarised the dialogue, resulting in a secondary dialogue summarisation dataset.

The primary contributions of this work are twofold: the dual-corpus, comprising the dialogue and dialogue summarisation datasets; and the software employed for data collection. These contributions are pivotal in advancing research in NLG for Scottish Gaelic and hold significant potential for future developments in the field. Moreover, this paper sheds light on the challenges and complexities encountered when gathering high-quality dialogue datasets involving native speakers of low-resource languages. By addressing these challenges and presenting a robust corpus and data collection methodology, this work enables further progress in low-resource NLG, within and beyond Scottish Gaelic.

## 2   Corpus Collection

We ground our data collection in multi-modal sources about exhibits found at the National Museum of Scotland. Participants discuss each exhibit through a chat interface based on slurk (Götze et al., 2022).

### 2.1   Task Description

We adopt a conversational question-answering task grounded in short texts combined with images. Pairs of participants engage in conversation about museum exhibits, with each acting as a museum visitor (the QUESTIONER) or a museum guide (the RESPONDER). Both participants see the same heading and image for the exhibit, but in addition to

---

[1] https://www.scotlandscensus.gov.uk/census-results/at-a-glance/languages/

[2] The datasets, along with code for the interface, are available at https://github.com/NapierNLP/sgge.

this, the QUESTIONER sees also a list of keywords and definitions relating to the exhibit while the RESPONDER has access to a text (mean length: 405 words, std. dev. 55) providing more information about the exhibit. These textual grounding materials are provided in Gàidhlig to avoid additional influence from English and to avoid breaking immersion in the task. Participants always have access to a short version of the experiment instructions for their current role. Each conversation focuses on a single exhibit, and after each conversation, each participant summarises the discussion. For a more detailed discussion of the task and the motivation behind it, refer to (Chandu et al., 2023).

## 2.2 Interface

We extend the `slurk` server (Götze et al., 2022) as the basis for our experimental interface, depicted in Figure 1. To begin, participants are assigned a role, a login token and receive a 'username' for the duration of the experiment session. This username includes an integer assigning them to a particular list of exhibits to be discussed, ensuring that each participant in a pair (with matching integers in their usernames) sees the same exhibit at the same time as their partner.

In addition to the interface itself, we made several modifications to the underlying server. To facilitate participants reconnecting when poor internet connectivity interrupted their session, we added support for returning users. We also modified the code for handling user commands (e.g. `/tòiseachadh` to begin the experiment) to account for whitespace errors.

## 2.3 Participant Recruitment

We set participant compensation at ∼£15/hour, with each experimental session booked as a 2-hour timeslot. Recruitment posts mentioned this rate (£30 for participation in a 2-hour study).

Our initial recruitment took place through social media (Twitter, Facebook, & Discord) and a blog post on the NLP Research Group page for Edinburgh Napier University. Potential participants completed a short (3-question) comprehension quiz based on a passage about a museum exhibit as part of the pre-screening process. Those with 100% accuracy were immediately invited to complete the consent form and scheduling poll, and those with a single error were contacted over email to assess whether the error was inadvertent or actual. The initial wave of recruitment over the course of 3

months resulted in 43 people completing the screening quiz, of whom 40 were invited to join the study. Scheduling via pre-designated timeslots[3] proved to be too challenging given the relative scarcity of participants, so we switched to a general availability model. This manual process did require more correspondence compared to the pre-designated timeslots; however, this extra workload for the research team greatly increased the ease of scheduling.

With this change, we began the second three-month phase of recruitment, adding Mastodon to our set of recruitment channels. In the final month of this phase, we also published a notice on Facebook via the Edinburgh Napier University page and paid 100 GBP to promote the post throughout the Scottish Highlands and Islands. Our second wave of recruitment yielded another 26 potential participants, of whom 19 were invited to join the study.

Once scheduled, participants received an email with the full instructions for the experiment along with a copy of the articles for the exhibits for which they would play the RESPONDER role. This way each participant could be familiar with the exhibits about which they would answer questions, making it easier for them to know where to find the answers they needed. The instruction email also provided two links, one for each half of the experiment session, so that the participants could switch roles after completing a number of dialogues. Each pair of participants was assigned 12 exhibits in total, 6 in each role, although depending on the session participants might not get to every exhibit.

Our study received standard institutional ethical and data management oversight.

**Challenges with Recruitment** We include a detailed breakdown of the attrition rate during recruitment to highlight the biggest challenge we encountered in our work: reaching the relatively small number of speakers of Scottish Gaelic. From the 59 participants invited to join the study, 42 completed the consent form and 19 were successfully paired and scheduled. This problem is pronounced despite the language being spoken in a highly developed country with reasonably good internet connectivity across the region. Researchers working with limited groups of speakers will need to spend considerably more time, effort, and funds on recruitment compared to what they may be used to

---

[3]fixed times and dates in a Doodle poll

444

Figure 1: Our task interface. The left side of the screen contains a chat interface. Participants enter their messages to send in the text field at the bottom left. On the top right, a summary of the instructions for the user's role is displayed. Below this, the exhibit label and image are followed by the textual grounding material. This screenshot shows the RESPONDER view with enlarged text for legibility.

| | |
|---|---|
| **Q**: Dè th' ann an tàileasg? | **Q**: What is chess? |
| **R**: 'S e geama a th' ann. Bidh dà chluicheadair a' gluasad phìosan timcheall air bòrd a' feuchainn ri 'rìgh' an neach eile a ghlacadh. | **R**: It's a game. Two players move pieces around a board trying to capture the other's 'king'. |
| **Q**: Cò na daoine a chruthaich na pìosan Thàileisg Leòdhais? | **Q**: Who were the creators of the Lewis Chess pieces? |
| **R**: Thathas den bheachd gun deach an cruathachadh ann an Nirribhidh. Tha iad air an dèanamh le ìbhri each-mara a thàinig à Graonlainn. | **R**: It is believed that they were hardened in Norway. They are made from walrus ivory that came from Greenland. |
| **Q**: Tha seo inntinneach. Carson a lorgar rudeigin air a dhèanamh ann an Nirribhidh ann an Leòdhas? | **Q**: This is interesting. Why find something made in Norway in Lewis? |
| **R**: Aig an àm bha Leòdhas, agus Innse Gall air fad, na phàirt den t-saoghal Lochlannach. Bhiodh daoine à Nirribhidh a' siubhal air ais is air adhart eadar Leòdhas is na dùthchannan Lochlannach agus a' toirt rudan mar seo leotha. | **R**: At the time Lewis, and the whole of Hebrides, was part of the Scandinavian world. People from Norway would travel back and forth between Lewis and the Scandinavian countries and bring things like this with them. |
| **Q**: Agus ciamar a chaidh na pìosan seo a lorg ann an Leòdhas? An robh iad am falach badeigin? | **Q**: And how were these pieces found in Lewis? Were they hiding somewhere? |
| **R**: Is dòcha gun deach an tìodhlachadh, gun deach an tasgadh. Chaidh an lorg faisg air Camas Ùige. | **R**: Maybe they were buried, they were deposited. They were found near Camas Uige. |
| **Q**: Gabhaibh mo leisgeul nach do dh'fhaighnich mi roimhe seo, ach dè gu dearbh a th' ann an ìbhri? | **Q**: Please excuse me for not asking before, but what exactly is ivory? |
| **R**: 'S e na fiaclan fada aig eich-mhara a th' ann. Canaidh sinn ivory ris sa Bheurla. Bidh daoine nas cleachdte ri ìbhri a gheibhear bho ailbheanan, 's dòcha! | **R**: It is the long teeth of walruses. We call it ivory in English. People will be more used to ivory from elephants, perhaps! |
| **Q**: Agus ceist mu dheireadh: cùin a chruthaich daoine na pìosan seo? Cò an linn? | **Q**: And a final question: when did people create these pieces? Whose era? |
| **R**: Chaidh an dèanamh uaireigin ron 12mh linn. Uaireigin sna meadhan aoisean. | **R**: They were made sometime before the 12th century. Sometime in the Middle Ages. |
| **Q**: Glè mhath. Taing mhòr. | **Q**: Very Good. Thanks a lot. |

Table 1: Example dialogue between a QUESTIONER (Q) and RESPONDER (R) about the Lewis Chess Pieces in Scottish Gaelic, along with a translation to English.

**Q** (summary): 'S e geama-bùird a th' ann an Tàileasg, far am bhios dà chluicheadair a' gluasad phìosan timcheall air bòrd a' feuchainn ri 'rìgh' an neach eile a ghlacadh. Rinn cudeigin ann an Nirribhidh na pìosan Tàileisg Leòdhais, à ìbhri each-mara a thàinig a Graonlainn. Thàinig iad gu Leòdhas oir bha Leòdhas agus na h-Innse Gall air fad nam pàirt dhen 'shaoghal Lochlannach' aig an àm. Chaidh na pìosan an cruthachadh uaireigin ron 12mh linn, anns na meadhan aoisean. Chaidh na pìosan Thàileisg a lorg faisg air Camas Ùige. 'S e fiaclan fada aig eich-mhara a th' ann an ìbhri, ged a thigeadh ìbhri bhon ailbheann cuideachd.

**R** (summary): Bhruidhinn sinn mu fhir Thàileasg Leòdhais. Dè th' ann an Tàileasg, cò chruthaich na pìosan, carson a lorgadh rudan a chaidh a dhèanamh ann an Nirribhidh ann an Leòdhas agus ciamar a chaidh an lorg. Bhruidhinn sinn cuideachd air dè th' ann an ìbhir agus air cuin a chaidh na pìosan a chruthachadh.

---

**Q** (summary): Chess is a board game, where two players move pieces around a board trying to capture the other person's 'king'. Someone in Norway made the Lewis Chess pieces, from walrus ivory that came to Greenland. They came to Lewis because Lewis and the whole of the Hebrides were part of the 'Scandinavian world' at the time. The pieces were created sometime before the 12th century, in the Middle Ages. The pieces were found near Camas Uíge. Ivory is the long teeth of walruses, although ivory could also come from the elephant.

**R** (summary): We talked about the Lewis chess men. What Chess is, who created the pieces, why items made in Norway were found in Lewis and how they were found. We also discussed what ivory is and when the pieces were created.

Table 2: Summaries for the conversation in Table 1 written by the same users, along with translation to English.

with crowdsourcing-based studies on platforms like Prolific and Amazon's Mechanical Turk.

## 3 Corpus Analysis[4]

**Demographics** We recruited 19 participants, 13 of whom live in the Highlands and Islands and 4 of whom live in the Central Belt in Scotland. Nearly all report speaking (16) and hearing (17) Scottish Gaelic daily, with 15 also reading it daily and 11 writing daily. All participants used the language across all modalities at least monthly. Fourteen started learning the language before the age of 25. We had only one participant under 25 years old, with participants otherwise fairly evenly split among 25-34, 35-44, 45-54, and 55 or older.

**Backing Data Statistics** The corpus contains information about 12 different museum exhibits drawn from a pool of exhibits which were popular and/or relevant to Gaelic culture as recommended by a museum curator. We created summary texts in English based on these materials and translated them into Scottish Gaelic, using 1–3 images for each exhibit. QUESTIONERs received 7 keywords on average, with the text of the definitions about 170 words (std. dev. 69). The grounding texts for RESPONDERs typically contained 6 paragraphs and 385 words (std. dev. 53).

**Conversation Statistics** The dataset includes 43 conversations consisting of 870 turns. Table 1 shows one of the dialogues from our corpus. Conversations were 16.8 turns long on average (std. dev. 9.4) lasting about 11 min 53 seconds (std. dev.

5:31), with each turn taking an average of 51.6 seconds (std. dev. 18.4). The dataset has about 13.7k tokens (2.3k types), for an average of 16 words per turn (std. dev. 24). The QUESTIONER averaged 12 words/turn (std. dev. 24) versus 19 (std. dev. 24) for the RESPONDER, in line with roles each participant played.

**Summary Statistics** Across the 43 conversations, QUESTIONERs wrote 26 summaries and RESPONDERs wrote 35 summaries. The summaries contain 2.0k tokens (573 types), with an average of 33.7 words per message (std. dev. 16.8).

## 4 Discussion & Conclusion

Recruitment was a major obstacle to our data collection (cf. Sec. 2.3), with this kind of study being more akin to linguistic field work or a psycholinguistic lab experiment than typical crowdsourcing tasks. Our research greatly benefited from the expertise of our Gaelic specialist, who contributed invaluable cultural and linguistic knowledge.

This first dataset for Scottish Gaelic NLG is of the same order of magnitude as NLG datasets for English just a decade ago (e.g. Wen et al.'s (2015) restaurant corpus of 5k utterances), providing a solid starting point for developing prototype systems. Indeed, we are currently exploring response generation for grounded question answering, dialogue summarisation, and text summarisation in Scottish Gaelic using this dataset.

Beyond the dataset, our experimental interface enables other researchers to build multipurpose datasets combining summarisation and conversation and/or grounding in text and images, especially for low-resource languages. We look forward to ex-

---

[4]In addition to these summaries provided here, detailed statistics can be found in Appendix A.

panding this corpus in the future and enabling others to develop more datasets for more low-resource languages.

## 5 Ethical Implications & Limitations

Working on low-resource languages which have smaller communities of use raises a variety of important ethical considerations in common with other work in linguistics (cf. Rice, 2006; Eckert, 2014; D'Arcy and Bender, 2023). As part of the standard ethical considerations for research with human participants (e.g. consent, compensation, etc), we recognise both the increased importance of demographic information for identifying speakers' position within the language community and the increased risk of being de-anonymised based on exactly that information. To address these problems, we focused on broad geographic regions and age bands when collecting demographic information and did not collect additional information which was less important to our analyses (e.g. gender). This allows us to differentiate between speakers in regions where the language is spoken by larger proportions of the local populace and speakers elsewhere, as well as to see differences between younger learners and older, established speakers of the language while avoiding de-anonymising our participants. In addition to standard ethical considerations, we also gave participants the opportunity to be associated with the dataset by name, without having their name directly linked to their contributions, as a way of recognising their contributions to the development of language technologies for their language community. We also found that by working with a small number of participants who care about the language, we were able to collect high quality data.

Our dataset is representative of the way a relatively small number of speakers would discuss a dozen museum exhibits; while this is a useful starting point for developing NLG systems, we cannot claim that it is representative of conversational Scottish Gaelic more broadly.

## Acknowledgements

## References

Gábor Bella, Fiona McNeill, Rody Gorman, Caoimhin O Donnaile, Kirsty MacDonald, Yamini Chandrashekar, Abed Alhakim Freihat, and Fausto Giunchiglia. 2020. A major Wordnet for a minority language: Scottish Gaelic. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 2812–2818, Marseille, France. European Language Resources Association.

Thiago Castro Ferreira, Claire Gardent, Nikolai Ilinykh, Chris van der Lee, Simon Mille, Diego Moussallem, and Anastasia Shimorina. 2020. The 2020 bilingual, bi-directional WebNLG+ shared task: Overview and evaluation results (WebNLG+ 2020). In *Proceedings of the 3rd International Workshop on Natural Language Generation from the Semantic Web (WebNLG+)*, pages 55–76, Dublin, Ireland (Virtual). Association for Computational Linguistics.

Khyathi Raghavi Chandu, David M. Howcroft, Dimitra Gkatzia, Yi-Ling Chung, Yufang Hou, Chris Emezue, Pawan Rajpoot, and Tosin Adewumi. 2023. Lowrecorp: the low-resource nlg corpus building challenge. In *Proceedings of the 16th International Conference on Natural Language Generation*, Prague, Czech Republic and virtual meeting. Association for Computational Linguistics.

Alexandra D'Arcy and Emily Bender. 2023. Ethics in linguistics. *Annual Review of Linguistics*, 9(1):49–69.

Penelope Eckert. 2014. *Ethics in linguistic research*, page 11–26. Cambridge University Press.

Jana Götze, Maike Paetzel-Prüsmann, Wencke Liermann, Tim Diekmann, and David Schlangen. 2022. The slurk interaction server framework: Better data for better dialog models. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 4069–4078, Marseille, France. European Language Resources Association.

David M. Howcroft and Dimitra Gkatzia. 2022. Most NLG is low-resource: here's what we can do about it. In *Proceedings of the 2nd Workshop on Natural Language Generation, Evaluation, and Metrics (GEM)*, pages 336–350, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. The state and fate of linguistic diversity and inclusion in the NLP world. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6282–6293, Online. Association for Computational Linguistics.

Julia Nee, Genevieve Macfarlane Smith, Alicia Sheares, and Ishita Rustagi. 2021. Advancing social justice through linguistic justice: Strategies for building equity fluent nlp technology. In *Equity and Access in Algorithms, Mechanisms, and Optimization*, EAAMO

'21, New York, NY, USA. Association for Computing Machinery.

Keren Rice. 2006. Ethical issues in linguistic fieldwork: An overview. *Journal of Academic Ethics*, 4(1-4):123–155.

Tsung-Hsien Wen, Milica Gašić, Nikola Mrkšić, Pei-Hao Su, David Vandyke, and Steve Young. 2015. Semantically conditioned LSTM-based natural language generation for spoken dialogue systems. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1711–1721, Lisbon, Portugal. Association for Computational Linguistics.

Yingxiu Zhao, Zhiliang Tian, Huaxiu Yao, Yinhe Zheng, Dongkyu Lee, Yiping Song, Jian Sun, and Nevin Zhang. 2022. Improving meta-learning for low-resource text classification and generation via memory imitation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 583–595, Dublin, Ireland. Association for Computational Linguistics.

# Generating Multiple Questions from Presentation Transcripts: A Pilot Study on Earnings Conference Calls

**Yining Juan,**[1] **Chung-Chi Chen,**[2] **Hen-Hsen Huang,**[3] **Hsin-Hsi Chen** [1]

[1] Department of Computer Science and Information Engineering,
National Taiwan University, Taiwan
[2] AIST, Japan
[3] Institute of Information Science, Academia Sinica, Taiwan
ynjuan@nlg.csie.ntu.edu.tw, c.c.chen@acm.org,
hhhuang@iis.sinica.edu.tw, hhchen@ntu.edu.tw

## Abstract

In various scenarios, such as conference oral presentations, company managers' talks, and politicians' speeches, individuals often contemplate the potential questions that may arise from their presentations. This common practice prompts the research question addressed in this study: to what extent can models generate multiple questions based on a given presentation transcript? To investigate this, we conduct pilot explorations using earnings conference call transcripts, which serve as regular meetings between professional investors and company managers. We experiment with different task settings and methods and evaluate the results from various perspectives. Our findings highlight that incorporating key points retrieval techniques enhances the accuracy and diversity of the generated questions.

## 1 Introduction

Preparing for an oral presentation is a common task in various domains, particularly in professional settings. For instance, researchers who have had their papers accepted at conferences need to deliver either an oral or poster presentation to share their findings with fellow researchers. Politicians must prepare for debates during election periods, while company managers are required to deliver speeches to update investors on company operations. When crafting their presentation drafts, a fundamental concern arises: what kinds of questions might the audience ask? This paper introduces a novel task, Multi-Question Generation (MQG), to assist presenters in preparing for Q&A sessions.

Diverging from previous studies that predominantly focused on one-to-one question generation tasks (Du et al., 2017; Song et al., 2018), the proposed MQG task is a one-to-many question generation task. In other words, after the presentation, audiences typically pose multiple questions. Table 1 showcases examples of one-to-one question

| **Presentation** |
| --- |
| Good day, and welcome to the Apple Q4 fiscal year 2022 earnings conference call... |
| **One-to-One Question Generation** |
| *Condition: gross margin* |
| Can you talk a bit about gross margin puts and takes? |
| **Proposed MQG** |
| 1. Can you talk a bit about gross margin puts and takes? |
| 2. How you think about balancing the consumer price versus your own costs and kind of the associated follow-through? |
| 3. Any preliminary thoughts around capital intensity into fiscal 2023? |

Table 1: Examples of one-to-one question generation and the proposed MQG.

generation and the proposed MQG task. The objective of the one-to-one question generation task is to generate a question based on a given condition (e.g., gross margin). In contrast, the proposed MQG task aims to generate multiple possible questions concurrently. We argue that exploring question generation tasks in a one-to-many setting closely aligns with real-world scenarios. However, the one-to-many setting also presents numerous challenges and research questions, including:

1. **Retrieving Keypoints from Long Documents**: Language model limitations prevent the entire speech draft from being inputted into the models. This raises the research question of identifying which parts of the speech are important and likely to prompt questions. Consequently, keypoint retrieval becomes a crucial aspect for question generation. Can these keypoints improve the performance of the MQG task?

2. **Task Setting**: Differing from the one-to-one task setting, which involves generating one question given a passage and a condition, the proposed MQG task requires generating multiple questions. This leads to the following research questions: Can models generate all questions at once? Does generating questions

sequentially yield better results?

3. **Evaluation**: As previously generated terms can influence the output of models, evaluating the accuracy and diversity of the generated questions becomes challenging. Specifically, can models generate several questions on the same topic, or can they generate questions from different perspectives?

To investigate these research questions, we collect earnings conference call transcriptions, regular meetings between company managers and professional analysts. Our aim with the proposed MQG task is to generate questions similar to those posed by analysts after listening to managers' presentations. We provide the collected dataset for future research endeavors.

To address these research questions, we propose the MQG-KR approach, combining MQG with Keypoint Retriever (KR). Keypoints are retrieved using BERT, enhancing question generation performance. We explore two task settings: generating all questions simultaneously and generating questions sequentially. Preliminary evaluations show that the MQG-KR approach improves the diversity of generated questions.

## 2    Related Work

Generating good questions is a challenging task for both humans and machines. Previous studies have primarily focused on one-to-one question generation, often centered around generating questions for reading comprehension tests. Heilman and Smith (Heilman and Smith, 2010) introduced syntactic transformations to convert given statement sentences into questions. Jia et al. (Jia et al., 2020) improved performance by incorporating a paraphrase module into their model. Wang et al. (Wang et al., 2020) generated questions based on the knowledge graph path of the input sentence. Song et al. (Song et al., 2018) matched given answers and paragraphs to augment context information for question generation.

In contrast to previous studies focusing on generating questions for machine reading comprehension tests, our work proposes generating questions in live presentations. There are distinct differences between these tasks. Questions for machine reading comprehension tests inquire about content, with most answers explicitly provided within the text. Complex questions may require some common
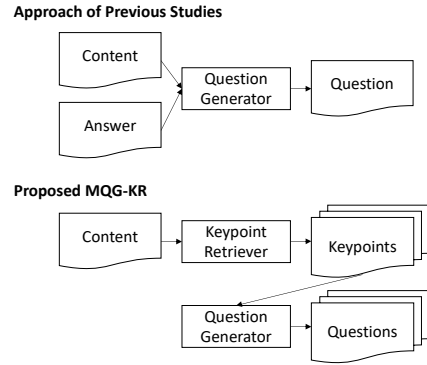


Figure 1: Overview of the proposed MQG-KR.

sense inference. However, professional analysts posing questions during presentations will not ask for information already provided; instead, they seek clarification or further explanation (Palmieri et al., 2015). Consequently, models generating such questions must first identify unclear or insufficiently detailed portions of the presentation. Although earnings conference calls have been widely studied for various tasks such as stock movement prediction (Medya et al., 2022), volatility forecasting (Qin and Yang, 2019; Sawhney et al., 2020), and summarization (Mukherjee et al., 2022), the question generation task has received little attention. Thus, our paper aims to fill this research gap by providing an initial exploration of question generation from earnings conference calls.

## 3    Method

Approach comparison is illustrated in Figure 1, highlighting two key distinctions between the proposed MQG-KR and previous approaches. Firstly, instead of providing the answer (condition) as input to the models, the proposed approach identifies keypoints that are likely to prompt questions from the audience and generates questions based on these keypoints. Secondly, while previous studies primarily focused on generating a single question, the proposed MQG-KR aims to generate multiple questions. In this section, we present two modules within the proposed MQG-KR framework: the Keypoint Retriever and the Question Generator.

### 3.1    Keypoint Retriever

As highlighted in Section 1, the length of the entire presentation exceeds the input capacity of most models. Additionally, not all sentences within the presentation hold investment-related significance, and professional analysts may not base their ques-

tions on these unimportant sentences. For instance, the greeting sentences in Table 1 do not provide valuable information to investors and may introduce noise to the models. To address this concern, we propose a keypoint retriever to refine the input.

Ideally, manual annotations for keypoints would be beneficial. However, the process of manual annotation is both expensive and time-consuming. Moreover, annotators without financial backgrounds may find it challenging to identify relevant keypoints. Therefore, we leverage the managers' answers to pinpoint the related section in their presentation. Specifically, we employ BM25 (Robertson et al., 2009) to calculate the similarity between the answer and each sentence in the presentation. Subsequently, we select the top-5 most similar sentences as the keypoints. This approach allows us to obtain fuzzy annotations for keypoints. Importantly, this process is applied solely to the training data, ensuring there are no issues related to data leakage. Once the keypoint labels are established, we train a classifier to identify the keypoints within the presentation, enabling the generation of questions based on these keypoints. For the keypoint retriever, we employ BERT (Devlin et al., 2019).

## 3.2 Question Generator

Our objective is to assess the effectiveness of the proposed approach, MQG-KR, in the novel task of Multi-Question Generation (MQG). We employ the well-performing generative model, FROST (Narayan et al., 2021), as our question generator. During the training process of FROST question generator, the entity chain of the presentation and questions is provided. In the inference (test) process, the models are required to generate both the entity chains and the questions. This approach has shown promise in abstractive summarization tasks (Narayan et al., 2021). In this paper, we present an initial exploration of FROST in the context of the one-to-many question generation task.

## 4 Experiment

### 4.1 Dataset

We compile a dataset of 995 transcriptions of earnings conference calls obtained from Seeking Alpha[1]. This dataset encompasses presentations from 18 different companies. During these 995 earnings conference calls, a total of 32,115 questions were

---

[1] https://seekingalpha.com/earnings/earnings-call-transcripts

asked. On average, each presentation received approximately 32 questions from the audience. In our task setting, models are required to generate multiple questions based on the provided presentation. We split the dataset into an 80% training set and a 20% test set for evaluation purposes, respectively.

### 4.2 Baselines

In addition to the vanilla FROST model, we consider two other baselines: Longformer (Beltagy et al., 2020) and LongT5 (Guo et al., 2022). These models are specifically designed to handle longer documents. Longformer employs sparsity in the attention matrix and utilizes a global and sliding window approach for encoding longer sequences. On the other hand, LongT5 is an extension of the T5 model (Raffel et al., 2020) and adopts a similar approach to Longformer. Notably, LongT5 has demonstrated superior performance compared to Longformer across six summarization datasets (Guo et al., 2022).

### 4.3 Evaluation

We employ the ROUGE-L score for evaluating the generated results (Lin, 2004). Additionally, as the proposed MQG task involves generating multiple questions, we propose two additional evaluation metrics, namely ROUGE-AMG and ROUGE-AMR, to assess the results from different perspectives. Each generated question ($GQ_i$) is assigned a list of ROUGE-L scores ($GL_i$) with each reference question (ground truth).

ROUGE-AMG is calculated using the following equation:

$$ROUGE - AMG = \frac{\sum_{i=1}^{N} \max(GL_i)}{N}, \quad (1)$$

where $N$ is the number of generated questions. ROUGE-AMG measures the extent to which the generated question is similar to the reference question.

On the other hand, we also evaluate the results from the perspective of reference questions using the ROUGE-AMR metric. Each reference question ($RQ_j$) receives a list of ROUGE-L scores ($RL_j$) with each generated question. ROUGE-AMR is calculated as follows:

$$ROUGE - AMR = \frac{\sum_{j=1}^{M} \max(RL_j)}{M}, \quad (2)$$

where $M$ is the number of reference questions.

As mentioned in Section 1, models may generate questions on the same topic by merely rephrasing

451

| | Question Generator | Max Input Length | ROUGE-L (↑) | ROUGE-AMG (↑) | ROUGE-AMR (↑) | Diversity (↓) |
|---|---|---|---|---|---|---|
| | Longformer | 4,096 | 19.37 | 18.21 | 15.54 | 100.00% |
| Baseline | LongT5 | 4,096 | 20.48 | 19.23 | 15.37 | 100.00% |
| | FROST | 1,024 | 23.08 | 22.20 | 17.95 | 100.00% |
| | LongFormer | 4,096 | 24.26 | 21.82 | 18.29 | 96.48% |
| MQG-KR | LongT5 | 4,096 | 24.43 | 22.65 | 18.66 | 96.48% |
| | FROST | 1,024 | **26.93** | **25.79** | **21.33** | **95.47%** |

Table 2: Experimental results. ↑ and ↓ denote the higher the better and the lower the better, respectively.

| | ROUGE-L |
|---|---|
| DialogueVED | 22.08 |
| PLATO | 22.13 |
| MQG-KR (FROST) | **26.93** |

Table 3: Results of sequential generation.

the question. However, since our goal is to generate diverse questions that could potentially be asked, we further evaluate the diversity by examining whether the most similar reference question for each generated question is the same or not. We calculate the ratio based on the maximum number of questions similar to the same reference question. Therefore, from a diversity perspective, a lower ratio is considered a better evaluation metric.

### 4.4 Experimental Results

The experimental results are presented in Table 2. Firstly, we observe that the proposed MQG-KR consistently yields improved performance across different question generators. These findings emphasize the significance of the keypoint retriever in the context of the proposed MQG task. Secondly, despite FROST having a shorter maximum input length compared to the other two models, it achieves the best performance among the baselines. This outcome may be attributed to the design of the entity chain prediction task within the decoder component, indicating the importance of entity chains in the proposed MQG task. Lastly, with respect to diversity, we find that all questions generated by the vanilla question generators pertain to the same topic. Conversely, the proposed MQG-KR models exhibit the ability to generate a wider range of diverse questions. Notably, MQG-KR combined with FROST demonstrates the highest diversity performance based on the evaluation conducted.

### 4.5 Sequential Generation

In earnings conference calls, analysts typically ask questions one by one, with later questions often following up on the previous ones. To simulate this process, we also experiment with a sequential ques-

tion generation setting. After generating a question, it is then used as input to the question generator to generate the subsequent question. For the sequential setting, we employ two well-performing dialogue generation models: PLATO (Bao et al., 2020) and DialogueVED (Chen et al., 2022). PLATO utilizes discrete latent variables to address one-to-many problems, while DialogueVED incorporates a multi-task pre-finetune process to enhance generated results and employs continuous latent variables for one-to-many generation tasks.

The experimental results are presented in Table 3. First, PLATO outperforms DialogueVED in the MQG task. Second, the performance of the models in the sequential generation setting falls short of FROST and the proposed MQG-KR, which operate under the simultaneous generation setting. However, since the models are designed for different purposes, we refrain from determining the best-performing model. Our aim in this paper is to explore the proposed MQG task from various perspectives, sharing our findings and establishing baselines for future studies.

## 5 Conclusion

This paper introduces the novel task of MQG and explores its potential applications in generating multiple questions based on a given presentation transcript. Our experiments demonstrate that incorporating a keypoint retriever improves the accuracy and diversity of the generated questions. This research contributes to the fields of NLP and Information Retrieval (IR) by offering insights into the MQG task and its relevance in various scenarios.

Our findings provide a valuable starting point for future research in MQG. By better anticipating and preparing for question and answer sessions, presenters can benefit from the generated questions. We believe that this work stimulates further discussions, advancements, and collaborations in the exciting field of Multi-Question Generation, driving the development of more effective and efficient question generation models.

## Acknowledgments

## Limitations

One limitation of this work is the focus on a specific application scenario, namely earnings conference calls, while only mentioning other real-world presentation and QA scenarios. Acquiring transcriptions for conference oral presentations or numerous meetings can be challenging and resource-intensive. Although video recordings may be available, the transcription process is time-consuming and costly. Moreover, data from politicians' debates are relatively scarce compared to earnings conference call data, as elections are not held every year. In contrast, quarterly earnings conference calls yield four transcriptions per company annually. Future studies can either develop new methods utilizing the proposed dataset or explore MQG tasks using alternative resources. This paper aims to serve as a starting point for the IR and NLP communities to delve into MQG tasks by improving keypoint retrieval performance and enhancing the ability to generate insightful questions.

Furthermore, we propose a more granular future direction for MQG. Subsequent work can involve annotating questions based on the question taxonomy proposed by Palmieri et al. (Palmieri et al., 2015). This taxonomy classifies analysts' questions into three levels, each comprising two to six labels. We believe that incorporating these labels can aid in automatically understanding analysts' questions and generating high-quality questions. Additionally, future research can explore assisting presenters in preparing answers for the generated questions, thereby progressing towards the development of a Q&A session tutor or assistant. This work highlights the significance of not only NLP techniques but also the relevance of IR considerations in this application domain.

## References

Siqi Bao, Huang He, Fan Wang, Hua Wu, and Haifeng Wang. 2020. PLATO: Pre-trained dialogue generation model with discrete latent variable. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 85–96, Online. Association for Computational Linguistics.

Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.

Wei Chen, Yeyun Gong, Song Wang, Bolun Yao, Weizhen Qi, Zhongyu Wei, Xiaowu Hu, Bartuer Zhou, Yi Mao, Weizhu Chen, Biao Cheng, and Nan Duan. 2022. DialogVED: A pre-trained latent variable encoder-decoder model for dialog response generation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4852–4864, Dublin, Ireland. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Xinya Du, Junru Shao, and Claire Cardie. 2017. Learning to ask: Neural question generation for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1342–1352, Vancouver, Canada. Association for Computational Linguistics.

Mandy Guo, Joshua Ainslie, David Uthus, Santiago Ontanon, Jianmo Ni, Yun-Hsuan Sung, and Yinfei Yang. 2022. LongT5: Efficient text-to-text transformer for long sequences. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 724–736, Seattle, United States. Association for Computational Linguistics.

Michael Heilman and Noah A. Smith. 2010. Good question! statistical ranking for question generation. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 609–617, Los Angeles, California. Association for Computational Linguistics.

Xin Jia, Wenjie Zhou, Xu Sun, and Yunfang Wu. 2020. How to ask good questions? try to leverage paraphrases. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6130–6140, Online. Association for Computational Linguistics.

Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.

Sourav Medya, Mohammad Rasoolinejad, Yang Yang, and Brian Uzzi. 2022. An exploratory study of stock price movements from earnings calls. In *Companion Proceedings of the Web Conference 2022*, WWW '22, page 20–31, New York, NY, USA. Association for Computing Machinery.

Rajdeep Mukherjee, Abhinav Bohra, Akash Banerjee, Soumya Sharma, Manjunath Hegde, Afreen Shaikh, Shivani Shrivastava, Koustuv Dasgupta, Niloy Ganguly, Saptarshi Ghosh, and Pawan Goyal. 2022. ECT-Sum: A new benchmark dataset for bullet point summarization of long earnings call transcripts. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10893–10906, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Shashi Narayan, Yao Zhao, Joshua Maynez, Gonçalo Simões, Vitaly Nikolaev, and Ryan McDonald. 2021. Planning with learned entity prompts for abstractive summarization. *Transactions of the Association for Computational Linguistics*, 9:1475–1492.

Rudi Palmieri, Andrea Rocci, and Nadzeya Kudrautsava. 2015. Argumentation in earnings conference calls. corporate standpoints and analysts' challenges. *Studies in communication sciences*, 15(1):120–132.

Yu Qin and Yi Yang. 2019. What you say and how you say it matters: Predicting stock volatility using verbal and vocal cues. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 390–401, Florence, Italy. Association for Computational Linguistics.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551.

Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389.

Ramit Sawhney, Piyush Khanna, Arshiya Aggarwal, Taru Jain, Puneet Mathur, and Rajiv Ratn Shah. 2020. VolTAGE: Volatility forecasting via text audio fusion with graph convolution networks for earnings calls. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8001–8013, Online. Association for Computational Linguistics.

Linfeng Song, Zhiguo Wang, Wael Hamza, Yue Zhang, and Daniel Gildea. 2018. Leveraging context information for natural question generation. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 569–574, New Orleans, Louisiana. Association for Computational Linguistics.

Siyuan Wang, Zhongyu Wei, Zhihao Fan, Zengfeng Huang, Weijian Sun, Qi Zhang, and Xuanjing Huang. 2020. PathQG: Neural question generation from facts. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9066–9075, Online. Association for Computational Linguistics.

# Mod-D2T: A Multi-layer Dataset for Modular Data-to-Text Generation

**Simon Mille**
ADAPT, Dublin City University
simon.mille@adaptcentre.ie

**François Lareau**
OLST, Université de Montréal
francois.lareau@umontreal.ca

**Stamatia Dasiopoulou**
Independent Researcher
stamatia.dasiopoulou@gmail.com

**Anya Belz**
ADAPT, Dublin City University
anya.belz@adaptcentre.ie

## Abstract

Rule-based text generators lack the coverage and fluency of their neural counterparts, but have two big advantages over them: (i) they are entirely controllable and do not hallucinate; and (ii) they can fully explain how an output was generated from an input. In this paper we leverage these two advantages to create large and reliable synthetic datasets with multiple human-intelligible intermediate representations. We present the Modular Data-to-Text (Mod-D2T) Dataset which incorporates ten intermediate-level representations between input triple sets and output text; the mappings from one level to the next can broadly be interpreted as the traditional modular tasks of an NLG pipeline. We describe the Mod-D2T dataset, evaluate its quality via manual validation and discuss its applications and limitations. Data, code and documentation are available at https://github.com/mille-s/Mod-D2T.

## 1 Introduction

Multi-level linguistic representations are inherent to several linguistic formalisms, including Lexical Functional Grammar (Dalrymple, 2001) and Meaning-Text Theory (Mel'čuk, 1973). They have been widely used in Natural Language Understanding (NLU), e.g. in Enhanced Universal Dependencies (Schuster and Manning, 2016) and Prague DT (Bejček et al., 2013), as well as in Natural Language Generation (NLG), e.g. in Enhanced WebNLG (Castro Ferreira et al., 2018), SRST (Mille et al., 2018), and Wikifluent (Kasner and Dusek, 2022). Reference architectures have been proposed that define modules and/or levels of representation, with the first theoretical architecture probably dating back to the work of Žolkovskij and Mel'čuk (1965), while a widely accepted applied NLG architecture was described by Reiter and Dale (1997). While end-to-end generators are very efficient (Dušek et al., 2018; Castro Ferreira et al.,

2020), there is evidence that splitting the generation process into sub-steps can lead to improvements (Castro Ferreira et al., 2019; Moryossef et al., 2019; Puduppully and Lapata, 2021; Kasner and Dusek, 2022). However, corresponding datasets with multiple intermediate representational levels are scarce.

In this paper, we present the Modular Data-to-Text (Mod-D2T) dataset, which comprises the inputs from the WebNLG 2020 shared task data (Castro Ferreira et al., 2020), paired with new output texts and ten intermediate-level representations that incrementally specify the output. The mappings from one level to the next can broadly be interpreted as the modular tasks of linguistic structuring, text planning/sentence aggregation, lexicalisation, communicative structure determination, deep sentence structuring, surface sentence structuring, surface aggregation, referring expression generation (REG), linearisation/morphology resolution and surface form generation. Unlike existing multi-level datasets, which were created by adding annotated layers on top of existing text, we leverage the FORGe rule-based pipeline generator (Mille et al., 2019b) to produce multiple human-intelligible intermediate (semantic, syntactic, morphological) and final (text) representations starting from abstract structures.

## 2 The Mod-D2T Dataset

In this section, we describe the dataset and how it was built. Table 1 lists the 10 intermediate levels of representation with associated tasks and approximate[1] correspondence to Reiter and Dale (1997). All examples in this section are for the same output text: *103 Colmore Row, designed by John Madin, is in Birmingham. It has 23 floors and was completed in 1976.*

---

[1] Our Surface sentence structuring spans Reiter and Dale (1997)'s Lexicalisation, REG and Linguistic realisation.

| Reiter&Dale Tasks | Mod-D2T Tasks | Mod-D2T Input | Mod-D2T Output |
|---|---|---|---|
| Content determination | — | — | — |
| Discourse planning | Linguistic structuring | WebNLG | PredArg |
| Sentence aggregation | Text planning* | PredArg | PredArg-Agg |
| Lexicalisation | Lexicalisation<br>Comm. structuring<br>Deep sent. structuring<br>Surf. sent. structuring<br>Synt. aggregation* | PredArg(-Agg)<br>PredArg-Lex<br>PredArg-Th<br>DSynt<br>SSynt | PredArg-Lex<br>PredArg-Th<br>DSynt<br>SSynt<br>SSynt-Agg |
| REG | REG* | SSynt(-Agg) | SSynt-Pro |
| Linguistic realisation | Word ord. and agree. resolution<br>Surface form retrieval | SSynt(-Agg/-Pro)<br>DMorph | DMorph<br>SMorph |

Table 1: The Mod-D2T layers (Mod-D2T Output) and tasks, and their correspondence with Reiter and Dale (1997)'s tasks; * Denotes optional modules, i.e., it is possible to generate grammatical texts without activating them.

```
<entry category="Building" eid="Id10" shape="(x (x) (x) (x) (x))"
    shape_type="sibling" size="4">
  <modifiedtripleset>
    <mtriple> 103_Colmore_Row | location | Birmingham </mtriple>
    <mtriple> 103_Colmore_Row | architect | Jonh_Madin </mtriple>
    <mtriple> 103_Colmore_Row | floorCount | 23 </mtriple>
    <mtriple> 103_Colmore_Row | completionDate | 1976 </mtriple>
  </modifiedtripleset>
</entry>
```

Figure 1: WebNLG'20 input triples.

## 2.1 Format

Intermediate representations in Mod-D2T are represented as CoNLL-U tables,[2] but not all CoNLL-U columns are used exactly as intended or at all. Because CoNLL-U is a linear format that we use to represent unordered graphs and trees, we delimit sentences by a <SENT> tag at the end of a group of nodes. All lines before <SENT> tag belong to the same sentence, but their relative order in the ConNLL-U file is not relevant. However, the order in which the sentences appear does correspond to their order in the text (see Table 2 for an example). For levels that are chains (in the sense explained below), the order of the lines is the order of the elements in the sentence.

## 2.2 Levels of representation

All ten intermediate representations in Mod-D2T are multi-sentence graphs that can be grouped into three main types: (i) **directed acyclic graphs (DAGs)** for semantic information; (ii) **unordered dependency trees** for syntactic information; and (iii) **chains** for morphological information. Nodes are connected across layers through individual IDs,

and coreference is explicitly marked (see the Misc column of, e.g., Table 2). Below, we describe each level of representation in turn, showing the last DAG and the last dependency tree for our running example in full here (Tables 2–3), while the other levels are shown in Appendix C. Appendix A presents some dataset statistics, and Appendix B the tag sets used.

### 2.2.1 WebNLG'20 inputs

The dataset is fully aligned with the WebNLG 3.0 release,[3] in which the inputs are sets of DBpedia triples (Subject|Property|Object), as described by Gardent et al. (2017); an example is shown in Figure 1. The labels from the WebNLG properties are stored in our annotations for a one-to-one mapping between properties and linguistic sub-structures.

### 2.2.2 Semantic levels: DAGs

**Predicate-argument graphs (PredArg)** are basic predicates linked to their arguments, mainly via numbered relations in the style of PropBank (Kingsbury and Palmer, 2002), with a few exceptions for common modifiers such as time and loca-

| ID | Lexeme | POS | Features | Head | Rel | Misc |
|----|--------|-----|----------|------|-----|------|
| 1 | design | VB | past | 0 | root | src=4 |
| 2 | John_Madin | NP | person\|ne | 1 | A0 | src=6 |
| 3 | be | VB | **rheme** | 0 | root | src=1 |
| 4 | 103_Colmore_Row | NP | ne | 3,1 | A1,A1 | src=2\|coref=0 |
| 5 | Birmingham | NP | location\|ne | 3 | A2 | src=3 |
| 6 | <SENT> | _ | _ | _ | _ | _ |
| 7 | 23 | CD | _ | 0 | root | src=9 |
| 8 | have | VB | **rheme** | 0 | root | src=7 |
| 9 | floor | NN | _ | 7,8 | A1,A2 | src=10 |
| 10 | 103_Colmore_Row | NP | ne | 8 | A1 | src=8\|coref=0 |
| 11 | <SENT> | _ | _ | _ | _ | _ |
| 12 | point_time_year | _ | _ | 0 | root | src=14 |
| 13 | 1976 | NP | year\|ne | 12 | A2 | src=14 |
| 14 | complete | VB | past\|**rheme** | 12 | A1 | src=12 |
| 15 | 103_Colmore_Row | NP | ne | 14 | A1 | src=13\|coref=0 |
| 16 | <SENT> | _ | _ | _ | _ | _ |

Table 2: Predicate-argument structure with thematicity (PredArg-Th).

tion. Due to the nature of WebNLG triples, nearly all predicates at this level of representation are binary. The main difference between PredArg graphs and the RDF input is that in the former, the content is structured linguistically, in terms of language-oriented representations based on meanings and predicate/argument relations between them. Table 9 gives an example of a text with four "sentences", i.e., elementary blocks of information, corresponding to the four input triples (Figure 1).

**Aggregated PredArg graphs (PredArg-Agg)** represent content packaging, where predicates that have common arguments can be merged into a sentence. Table 10 shows an example: the first two sentences from the previous level have been merged into one that will express both who designed the building and where it is located.

**Lexicalised PredArg graphs (PredArg-Lex)** replace meanings from the previous level with specific lexical units with an associated part of speech. In multilingual generation, this is where we pivot to the target language. Table 11 shows an example.

**PredArg graphs with thematicity (PredArg-Th)** give communicative structure to the text, i.e., establishes what each sentence asserts (the rheme) and what it asserts it about (the theme), as proposed by Mel'čuk (2001). In practice, we usually only add a rheme feature to the main node of the rheme, which is essential because it identifies the syntactic root. Table 2 shows an example, where each sentence has its main node identified. Multiple predicates that share an argument (e.g. $have(103\_Colmore\_Row, floor) \land 23(floor)$), are represented with heads and relations separated by a comma, as on line 9.

### 2.2.3 Syntactic levels: Dependency trees

**Deep-syntactic trees (DSynt).** Based on thematicity, we establish hierarchy and introduce deep syntactic relations between the meaningful lexical units of the sentence only, mostly distinguishing between complements (numbered) and modifiers (ATTR) (Mel'čuk, 1988; Kahane, 2009). This type of tree is roughly equivalent to a non-ordered UD representation (de Marneffe et al., 2021) stripped of cases, determiners and auxiliaries. This is also where we introduce semantically motivated features such as tense and number (see Table 13).

**Surface-syntactic trees (SSynt)** introduce function words and surface (usually language-specific) syntactic functions *à la* Mel'čuk (1988). This type of tree is similar to Surface-Syntactic UD (Gerdes et al., 2018), minus the linearity, or to Prague dependencies (Bejček et al., 2013), minus the morphematic nodes; Table 14 shows an example.

**Aggregated SSynt trees (SSynt-Agg)** introduce a more surface-oriented kind of aggregation that can only be performed once the syntactic structure has been computed. For example, in Table 15, the last two sentences of the previous structure share the same Subject, so they are merged into one sentence with coordinated main clauses. This operation is intended to increase fluency.

**Pronominalised SSynt trees (SSynt-Pro)** introduce pronouns where linguistically needed; see Table 3. This operation is intended to increase fluency and is sometimes needed for grammaticality.

### 2.2.4 Word-based levels: Chains

**Deep morphological chains (DMorph)** form the first linear layer. It introduces agreement resolution,

| ID | Lexeme | POS | Features | Head | Rel | Misc |
|----|--------|-----|----------|------|-----|------|
| 1  | be | VB | decl\|fin\|ind\|pres\|rheme | 0 | root | src=1 |
| 2  | 103_Colmore_Row | NP | sg\|ne | 1 | SBJ | src=2\|coref=0 |
| 3  | design | JJ | part | 2 | NMOD | src=4 |
| 4  | by | IN | _ | 3 | LGS | src=6 |
| 5  | John_Madin | NP | masc\|sg\|person\|ne | 4 | PMOD | src=6 |
| 6  | in | IN | _ | 1 | PRD | src=3 |
| 7  | Birmingham | NP | sg\|location\|ne | 6 | PMOD | src=3 |
| 8  | <SENT> | _ | _ | – | – | – |
| 9  | have | VB | decl\|fin\|ind\|pres\|rheme | 0 | root | src=7 |
| 10 | and | CC | _ | 9 | COORD | src=- |
| 11 | be | VB | decl\|fin\|ind\|past\|rheme | 10 | CONJ | src=12 |
| 12 | in | IN | _ | 11 | ADV | src=14 |
| 13 | 1976 | NP | year\|ne | 12 | PMOD | src=14 |
| 14 | floor | NN | pl | 9 | OBJ | src=10 |
| 15 | 23 | CD | _ | 14 | NMOD | src=9 |
| 16 | **_PRO_** | **PP** | sg\|ne | 9 | SBJ | src=8\|coref=0 |
| 17 | complete | VB | decl\|part\|rheme | 11 | VC | src=12 |
| 18 | **_PRO_** | **PP** | sg\|ne | 11 | SBJ | src=13\|coref=0 |
| 19 | <SENT> | _ | _ | – | – | – |

Table 3: Pronominalised surface syntactic representation (SSynt-Pro).

sentence-final punctuation (typically, a period), and ellipsis earmarking, as shown in Table 17.

**Surface morphological chains (SMorph)** list all tokens, including non-final punctuation (parentheses, commas, etc.), together with the POS tag, source ID and coreference ID, as in Table 18.

## 2.3 Dataset construction

We process automatically each WebNLG'20 input triple set by running 5 modules that consecutively (1) enrich the triple set, (2) populate PredArg templates and pre-order them based on the elements they have in common, (3) generate the text while saving intermediate layers in the process, (4) clean the generated text, and (5) clean and map the intermediate representations to CoNLL-U format.

Code for steps (4-5) was created for this paper, and for steps (1-3) we use as a starting point the FORGe pipeline (Mille et al., 2019b), which we tailored to our needs as follows. For (2), to maintain alignment, we copy from the WebNLG'20 inputs the information relative to the `category` and the `eid` to each input structure, and the respective property names to each <SENT> of each input structure (not shown in Appendix C). For (3), we modified FORGe by (i) separating the REG submodule from the linearisation submodule it was part of, and (ii) adding functionalities to maintain node and coreference alignments across levels; we also implemented a component that enables us to call each FORGe (group of) submodule(s) separately and store our 10 intermediate representations. Our 5 modules and the produced data

can be found in the following GitHub repository `https://github.com/mille-s/Mod-D2T`.

## 3 Qualitative Evaluation

We evaluated the quality of the intermediate representations by counting and classifying errors in PredArg-Th (Figure 2) and SSynt-Pro (Figure 3) for 30 randomly selected inputs of the WebNLG'20 dev split. These inputs contain 1 to 7 triples each, and are rendered as texts of up to 5 sentences.[4]

**PredArg-Th:** Out of the 30 semantic graphs, 20 (66.7%) were considered correct by a semantics expert (an author). The most common error was an overactive aggregation pattern that coordinated incompatible elements, typically a territory and one of its constituents, e.g., *Aarhus University is in Aarhus and Denmark*. This problem was found in 8 structures (26.7%) and it resulted in sentences that were still grammatical and faithful, but lacked fluency. We also identified one case of unnecessary predicate 'be' ('be' was used for *Abraham A. Ribicoff is American*; 'American' itself being predicative, we could do without the copula in the PredArg structure). The rest of the structures did not exhibit problems, but some representations could raise debate, such as (i) the use of "phantom

---

[4]For the random sample, we gave an equal probability to inputs of every size, although inputs of different sizes are not equally represented in the dataset: there are 10 to 20 times less 6- and 7-triple inputs than smaller sized ones. Large inputs are more challenging, and over-representing them as we did allowed us to detect more potential issues, but it should be noted that the numbers reported here are probably worse than they would be with a representative sample.

agents", as in *AIT is affiliated with VTU*, where the arguments are numbered A2 and A3, supposing an A1 that would affiliate one with the other (6 cases in 4 structures), and (ii) the choice of some edge labelling, e.g. for *Hypermarcas is in the pharmaceuticals sector*, which is represented with a Location relation between 'Hypermarcas' and 'sector'; a more solid representation would treat 'sector' as a predicate with 'Hypermarcas' as its argument.

**SSynt-Pro:** Out of the 30 syntactic structures, 28 (93%) were considered correct given the PredArg-Th representation by a syntax expert (another author). The two problems found were one case of superfluous determiner *the* next to a genitive complement (which would produce *the Baku's memorial* in the final text),[5] and one case of an underspecified dependency DEP where ADV would have been more appropriate.

**Text:** A previous human evaluation of the text quality of the FORGe generator that we use was provided by Castro Ferreira et al. (2020). In the "seen" scenario, they reported the following raw scores (out of 100): Data coverage: 95.3 (human-written text: 95.5), Relevance = 94.6 (94.1) , Correctness: 93.6 (93.4), Text structure: 87.0 (91.2), Fluency: 82.7 (88.1).

## 4 Related Work and Limitations

Moryossef et al. (2019) and Castro Ferreira et al. (2019) aligned WebNLG triples with the corresponding reference texts, making their datasets particularly appropriate for learning Reiter and Dale's Sentence aggregation. Castro Ferreira et al. (2019) also replaced the mentions of Subject and Object values with placeholders and lemmatised verbs in texts, allowing for learning both Lexicalisation as a whole and REG. Kasner and Dusek (2022) split Wikipedia paragraphs and rephrase splits into autonomous minimal sentences, replacing pronouns by their referent. They thus end up with two layers used to train a Sentence aggregation module and a paragraph compression module, which includes REG. They then apply their approach to the WebNLG dataset by crafting minimal sentence (as opposed to PredArg in our case) templates that they instantiate with the input triples. Mille et al. (2018) propose one syntactic and one predicate-argument levels using Universal Dependency annotations (de Marneffe et al., 2021) as a source.

Our dataset differs from the previous work in that we do not use human-written texts, and that we provide richer linguistic structures, with multiple semantic, syntactic and morphological levels that are (to the extent of our knowledge) not currently available for triple-to-text generation. The main limitation of our approach is that since the texts are synthetic and produced by a deterministic generator, their variety and quality is limited by the knowledge encoded in the generator (in particular, they generally lack the naturalness of human-written texts), and they represent only a fraction of what is possible for a language to express. Another (current) limitation is that Mod-D2T only contains English; generating texts in other languages requires crafting lexical resources and retrieving the Subject and Object values in the target language, adding rules to cover language-specific phenomena, and adapting PredArg templates; see (Mille et al., 2019a). We are currently completing Irish and French versions of Mod-D2T and will report on the multilingual aspects in future publications.

## 5 Conclusions and Future work

With the Mod-D2T dataset, we are making available a large amount of rich and reliable linguistic structures at several levels of representation for a sizeable set of D2T input/output pairs. These can be used, e.g., for experimenting with plug-and-play NLU and NLG, facilitating (human) language learning, or teaching linguistics.[6] The dataset construction process is flexible enough to allow the controlled production of a myriad of variants of the dataset in terms of verbalisation, sentence grouping/structuring, output simplicity/complexity, etc., simply by (de)activating optional modules (see Table 1) or introducing variation during the linguistic structuring task –thus providing multiple ways of verbalising each input triple. In contrast to neural generation, our approach ensures that the final text is faithful to the input, and will not contain inaccuracies, biases or offensive language.

In future work, in addition to the multilingual extension mentioned above, we will make the generation pipeline available for researchers to generate their own datasets, and provide mappings from our representations to standard representation schemes (e.g., Surface-syntactic UD).

---

[5] It is not the case since some rules filter out the superfluous *the* at a later stage in the pipeline.

[6] Thanks to their accuracy, the syntactic representations in particular can be used as teaching material, since unlike automatically parsed sentences, they will contain a negligible amount of errors.

## Acknowledgements

## Ethics statement

Given that we do not resort to using language models nor to human evaluation with people who are not authors of this paper, our work has no ethics implication that we are aware of.

## References

Eduard Bejček, Eva Hajičová, Jan Hajič, Pavlína Jínová, Václava Kettnerová, Veronika Kolářová, Marie Mikulová, Jiří Mírovský, Anna Nedoluzhko, Jarmila Panevová, Lucie Poláková, Magda Ševčíková, Jan Štěpánek, and Šárka Zikánová. 2013. Prague dependency treebank 3.0.

Thiago Castro Ferreira, Claire Gardent, Nikolai Ilinykh, Chris van der Lee, Simon Mille, Diego Moussallem, and Anastasia Shimorina. 2020. The 2020 bilingual, bi-directional WebNLG+ shared task: Overview and evaluation results (WebNLG+ 2020). In *Proceedings of the 3rd International Workshop on Natural Language Generation from the Semantic Web (WebNLG+)*, pages 55–76, Dublin, Ireland (Virtual). Association for Computational Linguistics.

Thiago Castro Ferreira, Diego Moussallem, Emiel Krahmer, and Sander Wubben. 2018. Enriching the WebNLG corpus. In *Proceedings of the 11th International Conference on Natural Language Generation*, pages 171–176, Tilburg University, The Netherlands. Association for Computational Linguistics.

Thiago Castro Ferreira, Chris van der Lee, Emiel van Miltenburg, and Emiel Krahmer. 2019. Neural data-to-text generation: A comparison between pipeline and end-to-end architectures. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 552–562, Hong Kong, China. Association for Computational Linguistics.

Mary Dalrymple. 2001. *Lexical functional grammar*, volume 34. Brill.

Marie-Catherine de Marneffe, Christopher D. Manning, Joakim Nivre, and Daniel Zeman. 2021. Universal Dependencies. *Computational Linguistics*, 47(2):255–308.

Ondřej Dušek, Jekaterina Novikova, and Verena Rieser. 2018. Findings of the E2E NLG challenge. In *Proceedings of the 11th International Conference on Natural Language Generation*, pages 322–328, Tilburg University, The Netherlands. Association for Computational Linguistics.

Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez-Beltrachini. 2017. Creating training corpora for NLG micro-planners. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 179–188, Vancouver, Canada. Association for Computational Linguistics.

Kim Gerdes, Bruno Guillaume, Sylvain Kahane, and Guy Perrier. 2018. SUD or surface-syntactic Universal Dependencies: An annotation scheme near-isomorphic to UD. In *Proceedings of the Second Workshop on Universal Dependencies (UDW 2018)*, pages 66–74, Brussels, Belgium. Association for Computational Linguistics.

Richard Johansson and Pierre Nugues. 2007. Extended constituent-to-dependency conversion for english. In *Proceedings of the 16th Nordic Conference of Computational Linguistics (NODALIDA 2007)*, pages 105–112.

Sylvain Kahane. 2009. Defining the deep syntactic structure: How the signifying units combine. In *Proceedings of the Meaning-Text Conference (MTT)*, Montréal, Canada.

Zdeněk Kasner and Ondrej Dusek. 2022. Neural pipeline for zero-shot data-to-text generation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3914–3932, Dublin, Ireland. Association for Computational Linguistics.

Paul Kingsbury and Martha Palmer. 2002. From Tree-Bank to PropBank. In *Proceedings of the Third International Conference on Language Resources and Evaluation (LREC'02)*, Las Palmas, Canary Islands - Spain. European Language Resources Association (ELRA).

Igor A. Mel'čuk. 1973. Towards a linguistic 'Meaning ↔ Text' model. *Trends in Soviet theoretical linguistics*, pages 33–57.

Igor A. Mel'čuk. 1988. *Dependency Syntax: Theory and Practice*. State University of New York Press, Albany, NY.

Igor A. Mel'čuk. 2001. *Communicative organization in natural language: the semantic-communicative structure of sentences*. John Benjamins, Amsterdam/ Philadelphia.

Simon Mille, Anja Belz, Bernd Bohnet, and Leo Wanner. 2018. Underspecified Universal Dependency structures as inputs for multilingual surface realisation. In *Proceedings of the 11th International Conference on Natural Language Generation*, pages 199–209, Tilburg University, The Netherlands. Association for Computational Linguistics.

Simon Mille, Stamatia Dasiopoulou, Beatriz Fisas, and Leo Wanner. 2019a. Teaching FORGe to verbalize DBpedia properties in Spanish. In *Proceedings of the 12th International Conference on Natural Language Generation*, pages 473–483, Tokyo, Japan. Association for Computational Linguistics.

Simon Mille, Stamatia Dasiopoulou, and Leo Wanner. 2019b. A portable grammar-based nlg system for verbalization of structured data. In *Proceedings of the 34th ACM/SIGAPP Symposium on Applied Computing*, pages 1054–1056.

Amit Moryossef, Yoav Goldberg, and Ido Dagan. 2019. Step-by-step: Separating planning from realization in neural data-to-text generation. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2267–2277, Minneapolis, Minnesota. Association for Computational Linguistics.

Ratish Puduppully and Mirella Lapata. 2021. Data-to-text generation with macro planning. *Transactions of the Association for Computational Linguistics*, 9:510–527.

Ehud Reiter and Robert Dale. 1997. Building applied natural language generation systems. *Natural Language Engineering*, 3(1):57–87.

Sebastian Schuster and Christopher D. Manning. 2016. Enhanced English Universal Dependencies: An improved representation for natural language understanding tasks. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 2371–2378, Portorož, Slovenia. European Language Resources Association (ELRA).

Aleksandr K. Žolkovskij and Igor A. Mel'čuk. 1965. O vozmožnom metode i instrumentax semantičeskogo sinteza [Method and instruments of semantic synthesis]. *Naučnotexničeskaja informacija*, 5:23–28.

## A   Statistics

There are 13,211, 1,667 and 1,779 texts in the training, development and test splits respectively. Tables 4-5 provide an overview of the number of nodes and sentences per text for all splits. Our 10 intermediate layers contain over 1.9 million nodes.

## B   Tagsets used

The edge labels for semantic graphs come mainly from PropBank (Kingsbury and Palmer, 2002), plus some generic labels such as Location and Time; see Table 6. The ones for deep syntactic trees come from Meaning-Text Theory (Mel'čuk, 1988); see Table 7. As for surface syntactic edge labels, they are a subset of the dependency Penn Treebank labels (Johansson and Nugues, 2007); see Table 8.

| Layer | N | S |
|---|---|---|
| PredArg | 152,664 | 48,776 |
| PredArg-Agg | 134,188 | 31,204 |
| PredArg-Lex | 134,188 | 31,204 |
| PredArg-Comm | 143,448 | 31,204 |
| DSynt | 169,325 | 31,204 |
| SSynt | 219,962 | 31,204 |
| SSynt-Agg | 222,970 | 27,557 |
| REG | 220,218 | 27,557 |
| DMorph | 247,795 | 27,557 |
| Text | 268,267 | 27,557 |

Table 4: Total number of nodes (N) and sentences (S) per layer.

| Layer | N | S | N/S |
|---|---|---|---|
| PredArg | 9.2 | 2.9 | 3.1 |
| PredArg-Agg | 8.1 | 1.9 | 4.4 |
| PredArg-Lex | 8.1 | 1.9 | 4.4 |
| PredArg-Th | 8.6 | 1.9 | 4.7 |
| DSynt | 10.2 | 1.9 | 5.5 |
| SSynt | 13.2 | 1.9 | 7.1 |
| SSynt-Agg | 13.4 | 1.7 | 8.2 |
| SSynt-Pro | 13.2 | 1.7 | 8.1 |
| DMorph | 14.9 | 1.7 | 9.1 |
| SMorph | 16.1 | 1.7 | 9.9 |

Table 5: Average number of nodes (N), sentences (S) and nodes per sentence (N/S) for each text, per layer.

## C   Sample structures

The annotations are released in CoNLL-U format, but because of space constraints, we have truncated the data in Tables 9–18 below:

- we dropped unused columns and renamed the remaining ones for readability;
- we removed feature names to retain only their values;
- we omit the metadata, which specifies the text ID (irrelevant here), the level of representation (see the captions) and the corresponding text string (see below).

The elements in bold below highlight the main changes between levels. The showcased structures all correspond to the following text:

*103 Colmore Row, designed by John Madin, is in Birmingham. It has 23 floors and was completed in 1976.*

| Label | Description | Example |
|---|---|---|
| A0–A6 | $n$-th argument of a predicate or quasi-predicate | speak→ English |
| Location | location | born→ Paris |
| Time | time | build→ 1932 |
| NonCore | inverted first argument of a predicate | runway→ second |
| Set | list of elements | and→ speak |
| Elaboration | (i) none of governor or dependent are argument of the other (ii) unknown argument slot | above me→ 610m |

Table 6: Edge labels: semantic graphs

| Label | Description | Example |
|---|---|---|
| I–VI | $n$-th complement of a syntactic predicate | speak→ English |
| ATTR | modifier | runway→ second |
| COORD | coordination | staff members→ and |
| APPEND | parenthetical modifier | Hypermarcas Brazil→ (s.a.) |

Table 7: Edge labels: deep syntactic trees

| Label | Description | Example |
|---|---|---|
| ADV | adverbial (broadly) | built→ in 1932 |
| AMOD | argument or modifier of an adjective | similar→ to |
| AMOD_COMP | argument of a comparative adjective | higher→ than |
| COORD | between conjunct and conjunction | and→ speak |
| DEP | underspecified | — |
| EXT | prepositional object (not *to*) | ask→for |
| IM | infinitive marker | to→ ask |
| IOBJ | dative object (after OBJ) | give→ her |
| LGS | logical subject | owned→ by |
| NMOD | argument or modifier of a noun | runway→ fifth |
| OBJ | non-prepositional object | give→ medal |
| OPRD | prepositional object (*to*) | give→ to |
| PMOD | complement of a preposition | to→ her |
| PRD | predicative complement | be→ president |
| SBJ | syntactic subject | play→ Beatles |
| SUB | complement of a conjunction | while→ be |

Table 8: Edge labels: surface syntactic trees

```
<entry category="Building" eid="Id10" shape="(x (x) (x) (x) (x))"
    shape_type="sibling" size="4">
  <modifiedtripleset>
    <mtriple> 103_Colmore_Row | location | Birmingham </mtriple>
    <mtriple> 103_Colmore_Row | architect | Jonh_Madin </mtriple>
    <mtriple> 103_Colmore_Row | floorCount | 23 </mtriple>
    <mtriple> 103_Colmore_Row | completionDate | 1976 </mtriple>
  </modifiedtripleset>
</entry>
```

Figure 2: WebNLG'20 input triples (same as Figure 1).

| ID | Semanteme | Features | Head | Rel | Misc |
|----|-----------|----------|------|-----|------|
| 1 | be | _ | 0 | root | src=1 |
| 2 | Birmingham | location\|ne | 1 | A2 | src=3 |
| 3 | 103_Colmore_Row | ne | 1 | A1 | src=2\|coref=0 |
| 4 | &lt;SENT&gt; | _ | _ | _ | _ |
| 5 | design | past | 0 | root | src=4 |
| 6 | John_Madin | person\|ne | 5 | A1 | src=6 |
| 7 | 103_Colmore_Row | ne | 5 | A2 | src=5\|coref=0 |
| 8 | &lt;SENT&gt; | _ | _ | _ | _ |
| 9 | have | _ | 0 | root | src=7 |
| 10 | 23 | _ | 0 | root | src=9 |
| 11 | floor | _ | 9,10 | A2,A1 | src=10 |
| 12 | 103_Colmore_Row | ne | 9 | A1 | src=8\|coref=0 |
| 13 | &lt;SENT&gt; | _ | _ | _ | _ |
| 14 | complete | past | 0 | root | src=12 |
| 15 | 1976 | year\|ne | 14 | Time | src=14 |
| 16 | 103_Colmore_Row | ne | 14 | A2 | src=13\|coref=0 |
| 17 | &lt;SENT&gt; | _ | _ | _ | _ |

Table 9: Predicate-argument structure (PredArg).

| ID | Semanteme | Features | Head | Rel | Misc |
|----|-----------|----------|------|-----|------|
| 1 | be | _ | 0 | root | src=1 |
| 2 | design | past | 0 | root | src=4 |
| 3 | 103_Colmore_Row | ne | **1,2** | **A1,A2** | src=2\|coref=0 |
| 4 | Birmingham | location\|ne | 1 | A2 | src=3 |
| 5 | John_Madin | person\|ne | 2 | A1 | src=6 |
| 6 | &lt;SENT&gt; | _ | _ | _ | _ |
| 7 | 23 | _ | 0 | root | src=9 |
| 8 | have | _ | 0 | root | src=7 |
| 9 | floor | _ | 7,8 | A1,A2 | src=10 |
| 10 | 103_Colmore_Row | ne | 8 | A1 | src=8\|coref=0 |
| 11 | &lt;SENT&gt; | _ | _ | _ | _ |
| 12 | complete | past | 0 | root | src=12 |
| 13 | 1976 | year\|ne | 12 | Time | src=14 |
| 14 | 103_Colmore_Row | ne | 12 | A2 | src=13\|coref=0 |
| 15 | &lt;SENT&gt; | _ | _ | _ | _ |

Table 10: Aggregated predicate-argument structure (PredArg-Agg).

| ID | Lexeme | **POS** | Features | Head | Rel | Misc |
|----|--------|---------|----------|------|-----|------|
| 1 | be | **VB** | _ | 0 | root | src=1 |
| 2 | Birmingham | **NP** | location\|ne | 1 | A2 | src=3 |
| 3 | design | **VB** | past | 0 | root | src=4 |
| 4 | John_Madin | **NP** | person\|ne | 3 | A1 | src=6 |
| 5 | 103_Colmore_Row | **NP** | ne | 1,3 | A1,A2 | src=2\|coref=0 |
| 6 | &lt;SENT&gt; | _ | _ | _ | _ | _ |
| 7 | 23 | **CD** | _ | 0 | root | src=9 |
| 8 | have | **VB** | _ | 0 | root | src=7 |
| 9 | floor | **NN** | _ | 7,8 | A1,A2 | src=10 |
| 10 | 103_Colmore_Row | **NP** | ne | 8 | A1 | src=8\|coref=0 |
| 11 | &lt;SENT&gt; | _ | _ | _ | _ | _ |
| 12 | complete | **VB** | past | 0 | root | src=12 |
| 13 | 1976 | **NP** | year\|ne | 12 | Time | src=14 |
| 14 | 103_Colmore_Row | **NP** | ne | 12 | A2 | src=13\|coref=0 |
| 15 | &lt;SENT&gt; | _ | _ | _ | _ | _ |

Table 11: Lexicalised predicate-argument structure (PredArg-Lex).

| ID | Lexeme | POS | Features | Head | Rel | Misc |
|---|---|---|---|---|---|---|
| 1 | design | VB | past | 0 | root | src=4 |
| 2 | John_Madin | NP | person\|ne | 1 | A0 | src=6 |
| 3 | be | VB | **rheme** | 0 | root | src=1 |
| 4 | 103_Colmore_Row | NP | ne | 3,1 | A1,A1 | src=2\|coref=0 |
| 5 | Birmingham | NP | location\|ne | 3 | A2 | src=3 |
| 6 | <SENT> | _ | _ | _ | _ | _ |
| 7 | 23 | CD | _ | 0 | root | src=9 |
| 8 | have | VB | **rheme** | 0 | root | src=7 |
| 9 | floor | NN | _ | 7,8 | A1,A2 | src=10 |
| 10 | 103_Colmore_Row | NP | ne | 8 | A1 | src=8\|coref=0 |
| 11 | <SENT> | _ | _ | _ | _ | _ |
| 12 | point_time_year | _ | _ | 0 | root | src=14 |
| 13 | 1976 | NP | year\|ne | 12 | A2 | src=14 |
| 14 | complete | VB | past\|**rheme** | 12 | A1 | src=12 |
| 15 | 103_Colmore_Row | NP | ne | 14 | A1 | src=13\|coref=0 |
| 16 | <SENT> | _ | _ | _ | _ | _ |

Table 12: Predicate-argument structure with thematicity (PredArg-Th, same as Table 2).

| ID | Lexeme | POS | Features | Head | Rel | Misc |
|---|---|---|---|---|---|---|
| 1 | be | VB | **decl\|act\|fin\|pres**\|rheme | 0 | root | src=1 |
| 2 | Birmingham | NP | location | 1 | **II** | src=3 |
| 3 | 103_Colmore_Row | NP | _ | **1** | **I** | src=2\|coref=0 |
| 4 | design | VB | **part**\|past | **3** | **ATTR** | src=4 |
| 5 | John_Madin | NP | person | 4 | **I** | src=6 |
| 6 | <SENT> | _ | _ | _ | _ | _ |
| 7 | have | VB | **decl\|act\|fin\|pres**\|rheme | 0 | root | src=7 |
| 8 | floor | NN | **pl** | **7** | **II** | src=10 |
| 9 | 103_Colmore_Row | NP | _ | 7 | **I** | src=8\|coref=0 |
| 10 | 23 | CD | _ | **8** | **ATTR** | src=9 |
| 11 | <SENT> | _ | _ | _ | _ | _ |
| 12 | complete | VB | **decl\|pass\|fin**\|past\|rheme | 0 | root | src=12 |
| 13 | in | IN | _ | 12 | **ATTR** | src=14 |
| 14 | 1976 | NP | year | 13 | **II** | src=14 |
| 15 | 103_Colmore_Row | NP | _ | 12 | **II** | src=13\|coref=0 |
| 16 | <SENT> | _ | _ | _ | _ | _ |

Table 13: Deep syntactic representation (DSynt).

| ID | Lexeme | POS | Features | Head | Rel | Misc |
|---|---|---|---|---|---|---|
| 1 | be | VB | decl\|fin\|**ind**\|pres\|rheme | 0 | root | src=1 |
| 2 | 103_Colmore_Row | NP | **sg**\|ne | 1 | **SBJ** | src=2\|coref=0 |
| 3 | design | JJ | part | 2 | **NMOD** | src=4 |
| 4 | in | IN | _ | 1 | **PRD** | src=3 |
| 5 | Birmingham | NP | **sg**\|location\|ne | 4 | **PMOD** | src=3 |
| 6 | by | IN | _ | 3 | **LGS** | src=6 |
| 7 | John_Madin | NP | **masc\|sg**\|person\|ne | 6 | **PMOD** | src=6 |
| 8 | <SENT> | _ | _ | _ | _ | _ |
| 9 | have | VB | decl\|fin\|**ind**\|pres\|rheme | 0 | root | src=7 |
| 10 | floor | NN | **pl** | 9 | **OBJ** | src=10 |
| 11 | 23 | CD | | 10 | **NMOD** | src=9 |
| 12 | 103_Colmore_Row | NP | **sg**\|ne | 9 | **SBJ** | src=8\|coref=0 |
| 13 | <SENT> | _ | _ | _ | _ | _ |
| 14 | be | VB | decl\|fin\|**ind**\|past\|rheme | 0 | root | src=12 |
| 15 | in | IN | _ | 14 | **ADV** | src=14 |
| 16 | 1976 | NP | year\|ne | 15 | **PMOD** | src=14 |
| 17 | 103_Colmore_Row | NP | **sg**\|ne | 14 | **SBJ** | src=13\|coref=0 |
| 18 | complete | VB | decl\|**part**\|rheme | 14 | **VC** | src=12 |
| 19 | <SENT> | _ | _ | _ | _ | _ |

Table 14: Surface syntactic representation (SSynt).

| ID | Lexeme | POS | Features | Head | Rel | Misc |
|----|--------|-----|----------|------|-----|------|
| 1 | be | VB | decl\|fin\|ind\|pres\|rheme | 0 | root | src=1 |
| 2 | 103_Colmore_Row | NP | sg\|ne | 1 | SBJ | src=2\|coref=0 |
| 3 | design | JJ | part | 2 | NMOD | src=4 |
| 4 | by | IN | _ | 3 | LGS | src=6 |
| 5 | in | IN | _ | 1 | PRD | src=3 |
| 6 | Birmingham | NP | sg\|location\|ne | 5 | PMOD | src=3 |
| 7 | John_Madin | NP | masc\|sg\|person\|ne | 4 | PMOD | src=6 |
| 8 | <SENT> | _ | _ | _ | _ | _ |
| 9 | have | VB | decl\|fin\|ind\|pres\|rheme | 0 | root | src=7 |
| **10** | **and** | **CC** | _ | **9** | **COORD** | src=- |
| 11 | be | VB | decl\|fin\|ind\|past\|rheme | 10 | CONJ | src=12 |
| 12 | in | IN | _ | 11 | ADV | src=14 |
| 13 | 1976 | NP | year\|ne | 12 | PMOD | src=14 |
| 14 | 103_Colmore_Row | NP | sg\|ne | 9 | SBJ | src=8\|coref=0 |
| 15 | floor | NN | pl | 9 | OBJ | src=10 |
| 16 | 23 | CD | _ | 15 | NMOD | src=9 |
| 17 | complete | VB | decl\|part\|rheme | 11 | VC | src=12 |
| 18 | 103_Colmore_Row | NP | sg\|ne | 11 | SBJ | src=13\|coref=0 |
| 19 | <SENT> | _ | _ | _ | _ | _ |

Table 15: Aggregated surface syntactic representation (SSynt-Agg).

| ID | Lexeme | POS | Features | Head | Rel | Misc |
|----|--------|-----|----------|------|-----|------|
| 1 | be | VB | decl\|fin\|ind\|pres\|rheme | 0 | root | src=1 |
| 2 | 103_Colmore_Row | NP | sg\|ne | 1 | SBJ | src=2\|coref=0 |
| 3 | design | JJ | part | 2 | NMOD | src=4 |
| 4 | by | IN | _ | 3 | LGS | src=6 |
| 5 | John_Madin | NP | masc\|sg\|person\|ne | 4 | PMOD | src=6 |
| 6 | in | IN | _ | 1 | PRD | src=3 |
| 7 | Birmingham | NP | sg\|location\|ne | 6 | PMOD | src=3 |
| 8 | <SENT> | _ | _ | _ | _ | _ |
| 9 | have | VB | decl\|fin\|ind\|pres\|rheme | 0 | root | src=7 |
| 10 | and | CC | _ | 9 | COORD | src=- |
| 11 | be | VB | decl\|fin\|ind\|past\|rheme | 10 | CONJ | src=12 |
| 12 | in | IN | _ | 11 | ADV | src=14 |
| 13 | 1976 | NP | year\|ne | 12 | PMOD | src=14 |
| 14 | floor | NN | pl | 9 | OBJ | src=10 |
| 15 | 23 | CD | _ | 14 | NMOD | src=9 |
| 16 | **_PRO_** | **PP** | sg\|ne | 9 | SBJ | src=8\|coref=0 |
| 17 | complete | VB | decl\|part\|rheme | 11 | VC | src=12 |
| 18 | **_PRO_** | **PP** | sg\|ne | 11 | SBJ | src=13\|coref=0 |
| 19 | <SENT> | _ | _ | _ | _ | _ |

Table 16: Pronominalised surface syntactic representation (SSynt-Pro, same as Table 3).

| ID | Word | POS | Features | Misc |
|---|---|---|---|---|
| **1** | 103_Colmore_Row | NP | _ | src=2\|coref=0 |
| **2** | design | JJ | part | src=4 |
| **3** | by | IN | _ | src=6 |
| **4** | John_Madin | NP | _ | src=6 |
| **5** | be | VB | decl\|fin\|ind\|pres\|**sg** | src=1 |
| **6** | in | IN | _ | src=3 |
| **7** | Birmingham | NP | _ | src=3 |
| **8** | . | _ | _ | src=- |
| **9** | _PRO_ | PP | sg | src=8\|coref=0 |
| **10** | have | VB | decl\|fin\|ind\|pres\|**sg** | src=7 |
| **11** | 23 | CD | _ | src=9 |
| **12** | floor | NN | pl | src=10 |
| **13** | and | CC | _ | src=- |
| **14** | _PRO_ | PP | sg\|**delete** | src=13\|coref=0 |
| **15** | be | VB | decl\|fin\|ind\|past\|**sg** | src=12 |
| **16** | complete | VB | decl\|part | src=12 |
| **17** | in | IN | _ | src=14 |
| **18** | 1976 | NP | _ | src=14 |
| **19** | . | _ | _ | src=- |

Table 17: Deep morphological representation (DMorph).

| ID | Word | POS | Misc |
|---|---|---|---|
| 1 | 103_Colmore_Row | NP | src=2\|coref=0 |
| **2** | , | _ | src=- |
| 3 | **designed** | JJ | src=4 |
| 4 | by | IN | src=6 |
| 5 | John_Madin | NP | src=6 |
| **6** | , | _ | src=- |
| 7 | **is** | VB | src=1 |
| 8 | in | IN | src=3 |
| 9 | Birmingham | NP | src=3 |
| 10 | . | _ | src=- |
| 11 | **it** | PP | src=8\|coref=0 |
| 12 | **has** | VB | src=7 |
| 13 | 23 | CD | src=9 |
| 14 | **floors** | NN | src=10 |
| **15** | and | CC | src=- |
| **16** | **was** | VB | src=12 |
| 17 | **completed** | VB | src=12 |
| 18 | in | IN | src=14 |
| 19 | 1976 | NP | src=14 |
| 20 | . | _ | src=- |

Table 18: Surface morphological representation (SMorph).

# Author Index