

# Beyond Traditional Benchmarks: Analyzing Behaviors of Open LLMs on Data-to-Text Generation

Zdeněk Kasner and Ondřej Dušek

Charles University, Faculty of Mathematics and Physics

Institute of Formal and Applied Linguistics

Prague, Czech Republic

{kasner, odusek}@ufal.mff.cuni.cz

## Abstract

We analyze the behaviors of open large language models (LLMs) on the task of data-to-text (D2T) generation, i.e., generating coherent and relevant text from structured data. To avoid the issue of LLM training data contamination with standard benchmarks, we design QUINTD – a tool for collecting novel structured data records from public APIs. We find that open LLMs (Llama 2, Mistral, and Zephyr) can generate fluent and coherent texts in zero-shot settings from data in common formats collected with QUINTD. However, we show that the semantic accuracy of the outputs is a major issue: both according to human annotators and our reference-free metric based on GPT-4, more than 80% of the outputs of open LLMs contain at least one semantic error. We publicly release the code, data, and model outputs.<sup>1</sup>

## 1 Introduction

Large language models (LLMs; Ouyang et al., 2022; Touvron et al., 2023a,b; Jiang et al., 2023; Tunstall et al., 2023) have already left a mark in many areas of natural language processing (NLP). Surprisingly, their applicability to the task of data-to-text (D2T) generation (Reiter and Dale, 1997; Gatt and Krahmer, 2018) remains underexplored, with limited evaluation on a handful of well-established benchmarks only (Axelsson and Skantze, 2023; Yuan and Färber, 2023). Generating text from structured data is arguably challenging for LLMs, given the specifics of D2T generation, such as long inputs, complex non-linear structure, and strict requirements on semantic accuracy. However, a more significant issue is the lack of testing grounds. The current D2T generation benchmarks are not only getting saturated (Van Miltenburg et al., 2023), but also promote optimization towards traditional reference-based evaluation metrics, which

<sup>1</sup><https://d2t-llm.github.io/>

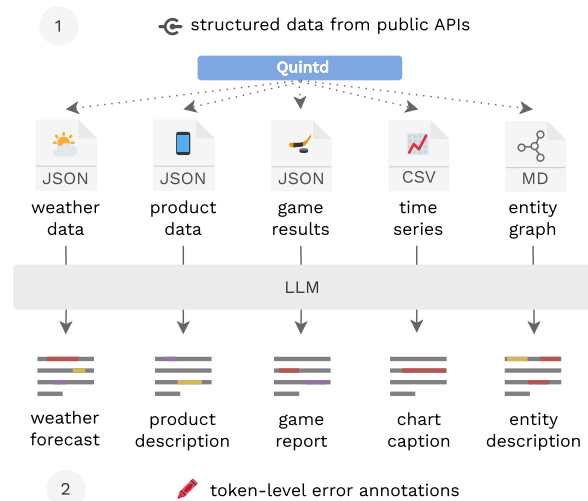


Figure 1: To benchmark LLMs, we download unlabeled structured data from public APIs and prompt LLMs to generate texts based on the data. We annotate semantic errors in the outputs using reference-free metrics.

were shown to correlate poorly with human judgment (Gehrmann et al., 2023; van der Lee et al., 2021; Novikova et al., 2017). When it comes to the models, using closed LLMs (OpenAI, 2023a,b) is increasingly considered a bad research practice due to its non-reproducibility (Rogers, 2023; Chen et al., 2023). On top of that, contamination of LLM training data with standard benchmarks further restricts the space for experiments (Golchin and Surdeanu, 2023; Aiyappa et al., 2023; Balloccu et al., 2024).

In this paper, we propose an approach that allows us to analyze model behavior in D2T generation on novel, real-world structured data records with reference-free evaluation metrics. We begin by realizing that *unlabeled data are plentiful*. To leverage the data for our experiments, we introduce QUINTD<sup>2</sup> – a tool for collecting structured data from five domains in standard formats: JSON,

<sup>2</sup>Quintet of Unlabeled Inputs for Natural Tasks in Data-to-text, pronounced as “quintet”

Task Id	Domain	Task Description	Source	Format
openweather	Weather	Generating a weather forecast from weather data.	OpenWeather	JSON
gsmarena	Technology	Describing a product based on its attributes.	GSMarena	JSON
ice_hockey	Sport	Describing an outcome of an ice-hockey game.	RapidAPI	JSON
owid	Health	Generating a caption for a time series.	OurWorldInData	CSV
wikidata	World facts	Describing entities and relations in a knowledge graph.	Wikidata	Markdown

Table 1: The domains and tasks included in the QUINTD data collection tool we use for testing D2T generation with LLMs. In our experiments, we download 100 development and 100 test examples of input data for each domain.

CSV, and Markdown. We choose the domains so that the data can be directly used as input for five distinct D2T generation tasks. Our tasks include generating weather forecasts, sports reports, product descriptions, chart captions, and entity descriptions (see Table 1). Next, we collect a set of 1,000 inputs with QUINTD and use the inputs as an ad-hoc benchmark (called QUINTD-1) for testing the abilities of LLMs for D2T generation. We assume that the data formats in QUINTD-1 are common in the LLMs’ pretraining corpora, so we specify the task using instructions instead of standard finetuning with human-written outputs, capitalizing on the zero-shot abilities of instruction-tuned LLMs (§2).

We push towards better reproducibility by *focusing on open LLMs*, which – apart from being more accessible – also achieve increasingly better results across tasks (Zheng et al., 2023; Beeching et al., 2023). For our experiments, we use three open LLMs with 7B parameters: Llama 2 (Touvron et al., 2023b; TogetherAI, 2023), Mistral (Jiang et al., 2023), and Zephyr (Tunstall et al., 2023). We also use GPT-3.5 (OpenAI, 2023b) as a closed model baseline for the final experiments. Given the behavioral nature of the experiments with LLMs (Holtzman et al., 2023), we put emphasis on reporting model behavior throughout the process (§3).

Another piece of the puzzle is *reference-free evaluation*: using the input data as a ground for comparison instead of reference outputs (§4). We focus on identifying semantic errors in the model outputs, i.e., the information that is not supported by the input data. We use two separate evaluation methods: manual annotations from human crowdworkers (van der Lee et al., 2021) and a custom automatic metric based on GPT-4 (Liu et al., 2023; Chiang and Lee, 2023; Kocmi and Federmann, 2023a). We annotate the errors on the level of individual words, getting fine-grained annotations of error spans in several categories (Thomson and Reiter, 2020; Thomson et al., 2023).

Based on our results, we provide general recommendations for D2T generation with open LLMs across tasks and formats (§5). Our main findings are as follows:

- **Open LLMs can generate fluent outputs from structured data** in common formats under zero-shot settings.
- **Semantic accuracy is a major obstacle**: both human annotators and GPT-4-based metric report that over 80% of outputs of open LLMs on our data contain a semantic error.
- **Long data inputs cause practical issues**, including the need for long-context models, increased GPU memory requirements, and unavailability of few-shot approaches.
- **Outputs can be empirically improved by following several rules-of-thumb** for preprocessing the model input, such as including units, removing unnecessary fields, or prefixing the model answer.

## 2 Reference-Free D2T Generation

### 2.1 Data Collection Tool

We introduce QUINTD,<sup>3</sup> a tool for collecting ad-hoc test sets using public APIs in five different domains. Our main reasons for departing from the traditional scheme of benchmarking on well-established datasets are:

1. Any published test sets may be potentially included in the training data of LLMs.
2. Public sources of structured data offer enough resources for creating ad-hoc test sets.
3. Without human references, our data collection scheme is lightweight and replicable.

Given the available public sources of data, we settled on the five tasks which are described in Table 1 (see Appendix A for more details). The tasks

<sup>3</sup><https://github.com/kasnerz/quintd>

### Prompt

```
Based on the given data:
...
{data}
...

Your task is to write a brief, fluent,
and coherent single-paragraph {output_type}
in natural language. The text should be
balanced and neutral. Make sure that all the
facts mentioned in the text can be derived
from the input data, do *not* add any extra
information.
```

### Output prefix

```
Sure! Here is the {output_type}:
"
```

Figure 2: The prompt  $\mathcal{P}$  and the model output prefix we used for the experiments in this paper. `data` is filled with the data record  $x$  and `output_type` is filled accordingly for each domain  $\mathcal{D}$  (see Table 1 and Table 6 in the Appendix).

are based on structured data in common formats: JSON, CSV, and Markdown.

## 2.2 QUINTD-1 Dataset

Using QUINTD, we collected the dataset for our experiments in this paper (QUINTD-1). The dataset contains 500 examples in the development set and 500 examples in the test set (100 examples per domain for each split). We downloaded the data between November 2023 and January 2024. Note that the dataset contains only **unlabeled** data without any reference outputs (e.g., weather data, but not a textual weather forecast), so the outputs need to be evaluated using reference-free metrics. New versions of the benchmark can be easily generated with our QUINTD tool.

## 2.3 Task Definition

Each example in QUINTD-1 consists of a structured data record  $x$  from a domain  $\mathcal{D} \in \{\text{openweather, gsmarena, ice\_hockey, owid, wikidata}\}$ . Given  $x$  and a prompt  $\mathcal{P}$ , the goal is to generate natural language output  $y$  faithful to the data  $x$ , according to the instructions in the prompt  $\mathcal{P}$  (see Figure 2).

# 3 Experiments

## 3.1 Experimental Process

Our goal is to avoid extensive data preprocessing and prompt engineering since these steps could harm the reproducibility and generalizability of our

experiments. With this goal in mind, we decided to use the same prompt template  $\mathcal{P}$  for all the domains and models.

For a set of preliminary experiments, we first wrote down the initial version of the prompt and used the data without further preprocessing. We then iteratively improved our experimental setup by observing outputs on the development set. In §3.3, we describe all the observations and modifications we made before generating the final outputs on the test set.

## 3.2 Models

For our experiments, we selected the following LLMs available under an open license:

- **Llama 2** (Touvron et al., 2023a; TogetherAI, 2023),  
togethercomputer/Llama-2-7B-32K-Instruct
- **Mistral** (Jiang et al., 2023),  
mistralai/Mistral-7B-Instruct-v0.1
- **Zephyr** (Tunstall et al., 2023).  
HuggingFaceH4/zephyr-7b-beta

The models are instruction-tuned, operate with 32k context, and perform well on recent benchmarks (Beeching et al., 2023). All the models have 7B parameters and thus fit on a single NVIDIA A40 (48G VRAM) in 16-bit precision. The models are available through HuggingFace (Wolf et al., 2020).

We accessed the models via an API provided by the text-generation-webui framework<sup>4</sup> running locally. For the final experiments, we also included GPT-3.5 (gpt-3.5-turbo-1106) accessed through the OpenAI API (OpenAI, 2023b).<sup>5</sup>

## 3.3 Observations from Preliminary Experiments

During development, we made several observations which we took into account for our final experimental setup:

**Any input field may appear in the output.** The models do not always select the most relevant fields for the given output. For example, we observed that

<sup>4</sup><https://github.com/oobabooga/text-generation-webui>

<sup>5</sup>We only included GPT-3.5 in our final experiments as proprietary models were not our focus. We did not include GPT-4 since we use the same model for evaluation (see §4.2) and LLMs tend to be biased towards their own outputs (Koo et al., 2023; Stureborg et al., 2024).

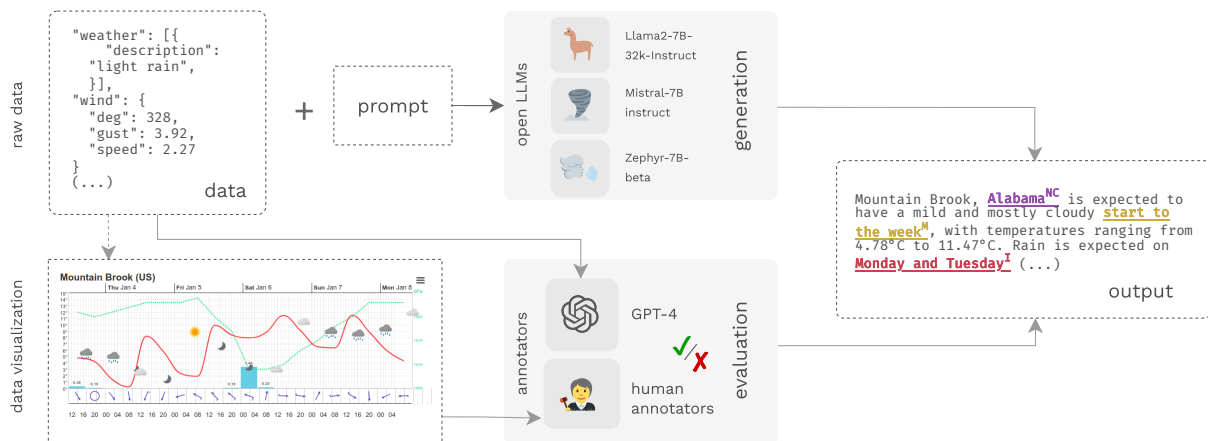


Figure 3: Our experimental setup. We first generate the outputs using LLMs that are given raw data and a task-specific prompt. We annotate the word-level semantic errors in the LLM outputs with (a) an automatic metric based on GPT-4 that matches the output to the raw data, and (b) human annotators, who annotate the errors in the output given the data visualization.

the models commonly mention identifiers, timestamps, files, and other metadata, leading to unnatural outputs. To mitigate these issues, we manually picked irrelevant fields and filtered them out from the input.

**Units need to be specified explicitly.** If the units are not specified in the data record, the models tend to resort to their best guess. This may go unnoticed if the unit is evident from the context (e.g., the model will usually not report the temperature in Fahrenheit instead of Celsius), but it may get problematic if the value is ambiguous (e.g., wind speed in km/h versus m/s). Therefore, we explicitly add units to all data records where appropriate.

**Understandable field names are enough.** On the flip side, we decided not to add extra descriptions to field names in the data if the field was understandable from its name (e.g., homeTeam or dimensions). As discussed by Kasner et al. (2023), pretrained models interpret field names correctly as long as they are human-readable. We only include chart metadata for the CSV files in the owid domain.

**Long inputs can be troublesome.** The inputs in some domains can easily get longer than 10-20k tokens. This issue is amplified by the fact that the evaluated LLMs tokenize numbers into individual digits. To accommodate for the long inputs, we picked models that accept up to 32k tokens.<sup>6</sup> How-

ever, with long inputs, the GPU memory consumption also gets considerably higher, so we needed to downsample the data in owid and openweather to keep their length under ~8k tokens.

**Few-shot experiments are infeasible.** Due to the above context-length limitations, we were not able to run few-shot experiments since we could not robustly fit an additional  $(x_{example}, y_{example})$  pair in the prompt. We attempted to include only  $y_{example}$  (making the setup “half-shot”), but we observed that the models then used entities from the example (unrelated to the actual input) in their outputs. Therefore, we decided to leave this line of experiments for future work (see §5.3 for discussion).

**Deterministic decoding and sampling are on par.** In our preliminary experiments, we observed a roughly similar output quality for both deterministic decoding and sampling.<sup>7</sup> For the final experiments, we decided to use deterministic decoding, which is non-parametric and conceptually more suitable for D2T generation.

**Prefixing the output makes parsing easier.** Even with variations of a “generate only the output” instruction appended to the prompt, the models (especially Llama 2) tended to first confirm the request. For that reason, we decided to prefix the input for all the models with “Sure! Here is the {output\_type}: ”. The opening quote at the end of the prefix allowed us to robustly parse the text sim-

<sup>6</sup>For this reason, we use Llama-2-7B-32k with 32k token context (TogetherAI, 2023) instead of the official Llama-2-7B-Instruct, which only supports 4k context (Touvron et al.,

2023b).

<sup>7</sup>We used the text-generation-webui default decoding parameters: temperature=0.7, top\_p=0.9, and top\_k=20.

ply by stripping the closing quote from the model output.

**The outputs are fluent but inaccurate.** We observed that the vast majority of model outputs were grammatically and stylistically correct, capturing the output type specified in the prompt. However, we also noticed that the outputs contained many semantic errors (even after emphasizing the focus on semantic accuracy in the prompt, see Figure 2). This observation led us to evaluate the model outputs using word-level annotations focused on semantic accuracy errors (see §4).

### 3.4 Final Experiments

Taking the observations in §3.3 into account, we proceeded to generate the outputs on the test set of QUINTD-1 for word-level error analysis. We first preprocessed the data as mentioned: we stripped out unnecessary fields, added units, and downsampled the data to fit the context. For all the models mentioned in §3.2, we used the prompt in Figure 2 and deterministic decoding with a maximum length of 512 tokens.

For comparison, we also generated outputs for the same inputs and identical prompts with GPT-3.5. Note that even though we fixed the temperature and seed to 0, the rest of the decoding parameters are inaccessible to us and may differ from the parameters we used for the open models.

## 4 Evaluation

For evaluation, we focus on identifying *semantic errors* in model outputs. We compare the generated texts to the input data, looking for parts of texts that are not faithful to the input data. We annotate the errors on the word level, considering all the words in the output text as potential sources of errors.

We use two complementary referenceless evaluation methods:

- $\mathcal{E}_{\text{hum}}$ : **human evaluation** based on crowdsourcing (§4.1),
- $\mathcal{E}_{\text{gpt}}$ : **an automatic metric** based on GPT-4 (§4.2).

The methods use similar instructions and produce outputs with equivalent format. The main idea is to compensate for the shortcomings of each approach: while human evaluation is costly (about ten times more expensive than automatic evaluation), using only an automatic metric based on a closed LLM would make the evaluation potentially

non-reproducible and biased (Kocmi and Federmann, 2023a; Wang et al., 2023b). Reporting the results of both methods should hopefully increase the robustness and replicability of our results.

Our error taxonomy and its notation are inspired by Thomson and Reiter (2020) and Thomson et al. (2023). After preliminary examination of the outputs, we settled on four error categories: **INCORRECT<sup>I</sup>**, **NOT\_CHECKABLE<sup>NC</sup>**, **MISLEADING<sup>M</sup>**, and **OTHER<sup>O</sup>**. To set clear boundaries between the categories, simplify the annotation interface and reach better inter-annotator agreement, we decided to keep our taxonomy more high-level than Thomson and Reiter and not to distinguish between fine-grained categories (e.g., *incorrect name* vs. *incorrect number*). The descriptions of our error categories, as presented in the instructions for annotation, are included in Table 2.

### 4.1 Human-based Evaluation

For the human annotation metric, we prepared a custom web interface, where an annotator is instructed to annotate text spans with respective error categories. We created custom visualizations for each data format (see Figure 3 and Appendix E for examples).<sup>8</sup>

We hired annotators on the Prolific<sup>9</sup> crowdsourcing platform. In total, we hired 100 annotators, each annotating 20 examples (four model outputs for each of the five domains). We selected annotators with at least 10 completed tasks, a 100% approval rate and English as their primary language. We paid the annotators £9 per hour, according to the platform’s recommendations. The median time for completing the annotations was 47 minutes. See Appendix B for the instructions for the annotators and the annotation interface.

### 4.2 GPT-4-based Evaluation

For automatic evaluation, we leverage the fact that LLM-based metrics can be customized for a particular task without the need for training data. In our experiments, we use a metric based on GPT-4 (gpt-4-1106-preview, OpenAI, 2023a), which was shown to be superior to other LLMs in following fine-grained instructions, reaching high correlations with human judgment on evaluating generated

<sup>8</sup>We open-sourced our annotation framework as a stand-alone software package, see <https://github.com/kasnerz/factgenie>.

<sup>9</sup><https://prolific.com>

Error	Description
<b>INCORRECT<sup>I</sup></b>	The fact in the text contradicts the data.
<b>NOT_CHECKABLE<sup>NC</sup></b>	The fact in the text cannot be checked given the data.
<b>MISLEADING<sup>M</sup></b>	The fact in the text is misleading in the given context.
<b>OTHER<sup>O</sup></b>	The text is problematic for another reason, e.g., grammatically or stylistically incorrect, irrelevant, or repetitive.
<b>Example</b>	
<i>data</i>	<b>Nokia 3310</b>   <i>color</i> : black, blue, grey   <i>display</i> : 320x240px
<i>text</i>	Nokia 3310 is <b>produced in Finland<sup>NC</sup></b> and features a <b>320x320<sup>I</sup></b> display. It is <b>available in black color<sup>M</sup></b> . <u>The data seem to provide only partial information about the phone.<sup>O</sup></u>

Table 2: Categories of errors annotated in our evaluation and an example demonstrating the error types. See Appendix B for an explanation of individual errors in the example.

texts (Zhao et al., 2023; Sottana et al., 2023; Kocmi and Federmann, 2023a,b).<sup>10</sup>

We instantiate  $\mathcal{E}_{\text{gpt}}$  with a prompt and a system message describing the task. We instruct the model to produce a JSON output with sequentially ordered errors using the following format:

```
{
  "errors": [{
    "reason": [REASON],
    "text": [TEXT_SPAN],
    "type": [ERROR_CATEGORY]
  },
  ...]
}
```

Note that we require that the model first generates the free-form text *reason* for the error. Generating the reason comes at almost no extra cost and our cursory observations suggest that requiring it leads to more precise outputs.<sup>11</sup>

We align the model outputs with the original text by string matching on TEXT\_SPAN, moving the current position forward after each match. We ensure that the model response is a valid JSON using OpenAI’s response\_format parameter. See Appendix C for more details about the metric, including the prompt and the system message.

## 5 Results and Discussion

A summary of the word-level annotations is in Table 3 and 4, with detailed results per domain provided in Appendix F.

<sup>10</sup>We confirmed that GPT-3.5 and Llama 3 have lower correlations with human judgments also in our scenario, see Appendix D.

<sup>11</sup>We did not ask the crowdworkers for free-form reasoning about the errors since that would make the annotation notably more complex.

### 5.1 How Accurate Are the Model Outputs?

Depending on the model, between 76-86% of examples contain an error according to  $\mathcal{E}_{\text{hum}}$ , suggesting that open LLMs make semantic errors very often. According to  $\mathcal{E}_{\text{gpt}}$ , the number is as high as 89-94%.

The most common error type is **INCORRECT<sup>I</sup>**. As shown in Table 3, all the open LLMs make more than **two statements contradicting the data per output on average**. The **NOT\_CHECKABLE<sup>NC</sup>** errors are also relatively common: more than one per output on average according to  $\mathcal{E}_{\text{hum}}$ , and at least one being present in more than 25% of examples according to both metrics.

The results vary widely according to the domain (see Appendix F). For example, the outputs in wikidata contain much more **NOT\_CHECKABLE<sup>NC</sup>** errors on average (1.01 per output according to  $\mathcal{E}_{\text{hum}}$ ) than **INCORRECT<sup>I</sup>** errors (0.11 per output according to  $\mathcal{E}_{\text{hum}}$ ), suggesting that with simpler inputs, the models tend to introduce extra information. The openweather domain seems to be the most complex with the longest outputs (~164 tokens), more than eight errors in the output on average, and >90% of outputs containing an error.

The differences between the open LLMs are not major. Out of the open LLMs, Zephyr has the best results across categories and metrics, followed by Llama 2. However, the outputs of Mistral are longer on average, leaving more space for errors. GPT-3.5 (which we consider separately) does generally better according to both  $\mathcal{E}_{\text{gpt}}$  and  $\mathcal{E}_{\text{hum}}$ , although it still makes an error in 60-75% of examples (2 errors per example on average). In general, the results show that LLMs make too many semantic errors to be usable in practice for D2T generation in a zero-shot setting.

	Incorrect		Not Checkable		Misleading		Other		All categories		# Tok.
	$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}}$	$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}}$	$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}}$	$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}}$	$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}}$	
<b>Llama 2</b>	1.57	<b>2.79</b>	1.25	0.91	0.25	<b>0.12</b>	<b>0.10</b>	0.09	3.18	3.90	83.8
<b>Mistral</b>	2.03	3.23	1.12	0.54	0.44	0.26	0.25	0.10	3.85	4.12	114.9
<b>Zephyr</b>	<b>1.44</b>	2.84	<b>0.77</b>	<b>0.40</b>	<b>0.20</b>	0.29	0.16	<b>0.05</b>	<b>2.58</b>	<b>3.58</b>	98.0
<b>GPT-3.5</b>	0.65	1.76	0.49	0.38	0.18	0.26	0.07	0.02	1.39	2.42	84.9

Table 3: The average *numbers of errors per output* (lower is better) based on human annotators ( $\mathcal{E}_{\text{hum}}$ ) and GPT-4 ( $\mathcal{E}_{\text{gpt}}$ ). We also include the average number of tokens per output in the rightmost column. The results of the best open LLM are emphasized.

	Incorrect		Not Checkable		Misleading		Other		All categories	
	$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}}$	$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}}$	$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}}$	$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}}$	$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}}$
<b>Llama 2</b>	53.2%	80.0%	57.4%	44.8%	17.4%	<b>8.8%</b>	<b>7.6%</b>	7.6%	85.6%	94.0%
<b>Mistral</b>	53.6%	80.2%	49.6%	31.8%	20.6%	17.0%	13.6%	8.4%	81.2%	93.0%
<b>Zephyr</b>	<b>46.8%</b>	<b>78.0%</b>	<b>42.2%</b>	<b>25.0%</b>	<b>16.2%</b>	20.6%	11.6%	<b>4.2%</b>	<b>75.6%</b>	<b>89.4%</b>
<b>GPT-3.5</b>	38.0%	65.0%	28.8%	19.6%	12.6%	16.2%	6.2%	2.2%	60.6%	75.8%

Table 4: The percentage of *outputs containing at least one error* (lower is better) based on human annotators ( $\mathcal{E}_{\text{hum}}$ ) and GPT-4 ( $\mathcal{E}_{\text{gpt}}$ ). The results of the best open LLM are emphasized.

## 5.2 Do Evaluation Methods Agree?

To quantify the agreement of  $\mathcal{E}_{\text{hum}}$  and  $\mathcal{E}_{\text{gpt}}$ , we computed the Pearson correlation coefficient between the error counts on the level of words, examples, and domains as follows (note that each error category was considered separately):

- For  $r_{\text{domain}}$ , we used the average error counts per domain (see Table 13).
- For  $r_{\text{example}}$ , we used the count of errors per example.
- For  $r_{\text{word}}$ , we used the binary indicators marking an error per word.

The correlation on the level of words is weak ( $r_{\text{word}} = 0.26$ ) but gets better on the example-level ( $r_{\text{example}} = 0.52$ ) and even better on the domain-level ( $r_{\text{domain}} = 0.93$ ). In Table 5, we show the percentage of words marked as errors by individual metrics. The metrics agree on the specific words in less than 6%, although they both mark around 21% of words as erroneous.

We also measure inter-annotator agreement between human annotators. For that, we obtained annotations from two annotators for 100 model outputs. The results are similar: the annotators agree weakly on the word level ( $r_{\text{word}} = 0.36$ ), stronger on the example level ( $r_{\text{example}} = 0.53$ ), and even stronger on the domain level ( $r_{\text{domain}} = 0.85$ ). We conclude that while the details regarding error spans and categories may vary, the annotators as well as GPT-4 generally agree on the accuracy

	$\mathcal{E}_{\text{gpt}}$	$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}} + \mathcal{E}_{\text{hum}}$
<b>Incorrect</b>	10.1%	14.2%	4.1%
<b>Not checkable</b>	7.8%	4.3%	2.0%
<b>Misleading</b>	2.2%	1.5%	0.1%
<b>Other</b>	1.8%	0.7%	0.1%
<b>Total</b>	21.9%	20.7%	6.3%

Table 5: The percentage of *words marked as erroneous* by human annotators ( $\mathcal{E}_{\text{hum}}$ ), GPT-4 ( $\mathcal{E}_{\text{gpt}}$ ), and both approaches at the same time ( $\mathcal{E}_{\text{hum}} + \mathcal{E}_{\text{gpt}}$ ).

of model outputs for a given set of examples. In the future, the agreement could be improved by measuring errors on the phrase level (Vamvas and Sennrich, 2022).

## 5.3 Recommendations for Future Work

**Focus on semantic accuracy.** The output of LLMs is satisfactory regarding the style, format, and purpose of the text. However, the amount of semantic errors remains very high. Improving the semantic accuracy of the models (Li et al., 2022), along with new model-based evaluation metrics (Liu et al., 2023; Xu et al., 2023), could thus help to bring improve LLM-based D2T generation systems where it is most needed.

**Use efficient and long-context models.** The memory issues with long context, making few-shot experiments infeasible, can potentially be solved by using more efficient models equipped with Flash Attention (Dao et al., 2022) and fast inference li-

libraries such as `llama.cpp`<sup>12</sup>. The recent emergence of capable long-context models (Bai et al., 2023; Munkhdalai et al., 2024) also seems to play in favor of LLM-based D2T generation with long inputs.

**Be careful about subtle bugs.** During our preliminary experiments, we fixed several subtle bugs in our API calls such as incorrect instruction templates<sup>13</sup> or involuntary input truncation. Therefore, we recommend careful checks of API calls, as with the apparent ease of API access and robustness of LLMs, such bugs could go unnoticed and artificially skew the model performance.

**Test the models in the wild.** Except for using an ad-hoc dataset of real-world data as we did in our work, the ecological validity of D2T evaluation can also be ensured by continuous evaluation with human users (Zheng et al., 2023) and evaluating the real-world impact of the systems (Reiter, 2023).

**Multilinguality is an opportunity.** With the recent efforts in extending D2T generation to low-resource languages (Cripwell et al., 2023), multilingual D2T generation with open LLMs seems a promising direction. Although we did not go beyond English, initial steps were already done by works such as Lorandi and Belz (2023).

## 6 Related Work

### 6.1 Evaluation of Generated Texts

Evaluation of generated texts is a complex task lacking a generally accepted solution (Celikyilmaz et al., 2020). While researchers are acknowledging the importance of combining multiple evaluation metrics (Gehrmann et al., 2021, 2023), most evaluation is still based on comparing the model outputs to human-written references, which tend to be noisy and expensive to collect (Dušek et al., 2019, 2020).

Many works recently investigated the potential of using LLMs for automatic reference-free evaluation of generated texts, generally achieving high correlations with human judgment (Zhao et al., 2023; Sottana et al., 2023; Kocmi and Federmann, 2023a,b; Chiang and Lee, 2023; Wang et al., 2023a; Fu et al., 2023). However, they also voice concerns about its non-reproducibility (Kocmi and Federmann, 2023a) and potential bias of these models (Wang et al., 2023b).

<sup>12</sup><https://github.com/ggerganov/llama.cpp>

<sup>13</sup>[https://huggingface.co/docs/transformers/chat\\_templating](https://huggingface.co/docs/transformers/chat_templating)

Human evaluation is an essential component of natural language generation experiments (van der Lee et al., 2019, 2021). The closest human evaluation protocol to our scenario is the reference-free word-level annotation of complex D2T generation output proposed in Thomson and Reiter (2020) and Thomson et al. (2023).

### 6.2 D2T Generation Tasks

**Weather Forecasts** First attempts for generating weather forecasts include template-based and statistical approaches (Belz, 2005, 2008; Angeli et al., 2010) for the Sumtime-meteo and WeatherGov datasets (Sripada et al., 2002; Liang et al., 2009). More recently, Balakrishnan et al. (2019) introduced a weather forecast dataset with tree-structured meaning representations. Our weather forecasts are less structured and based on a 5-day weather outlook.

**Product Descriptions** Our phone specifications are closest to Wen et al. (2015, 2016), who introduced a dataset for generating descriptions of laptops and TVs. Their solution was based on recurrent neural networks, although templates remained a go-to approach for the task (Wang et al., 2017). Recently, Shao et al. (2021) and Koto et al. (2022) also proposed specialized architectures based on pretrained language models for the data from big e-commerce platforms.

**Sport Reports** All the D2T generation datasets from the Rotowire family (Wiseman et al., 2017; Wang, 2019), including SportSett:Basketball (Thomson et al., 2021), and ESPN-NBA (Nie et al., 2018) focus on generating basketball reports. Along with MLB (Puduppully et al., 2019b), these datasets belong among the most challenging D2T datasets, attracting various neural-based solutions (Puduppully et al., 2019a, 2022; Puduppully and Lapata, 2021; Rebuffel et al., 2020). We use instead simpler data covering ice hockey game summaries.

**Chart Captions** Following the early rule-based approaches (Demir et al., 2008, 2012), the approaches for chart captioning recently tackle large-scale datasets from data analytic institutions (Obeid and Hoque, 2020; Kantharaj et al., 2022). We focus on one of the tasks from Sharma et al. (2021), which is generating descriptions of time series in the health domain.

**Entity Descriptions** The task of generating descriptions for a knowledge graph has been covered



extensively in D2T generation (Gardent et al., 2017; Castro Ferreira et al., 2020; Agarwal et al., 2021; Chen et al., 2020; Ribeiro et al., 2020, *inter alia*). Our task is to describe an entity provided a list of its properties, which is closely related to generating entity descriptions from Wikipedia infotables (Lebret et al., 2016).

### 6.3 D2T Generation with LLMs

Recent works have focused on exploring the capabilities of closed LLMs on existing D2T generation datasets. Axelsson and Skantze (2023) evaluated GPT-3.5 (OpenAI, 2023b) on WebNLG, along with Yuan and Färber (2023), who also tested the model on the AGENDA dataset (Koncel-Kedziorski et al., 2019). Both works found that regardless of potential data contamination, the LLMs rank behind state-of-the-art finetuned models on automatic metrics. Zhao et al. (2023) tested closed models on modified table-to-text generation datasets and found out that in terms of faithfulness, GPT-4 can outperform state-of-the-art models.

## 7 Conclusion

We provided an exploratory study into D2T generation with open LLMs. We proposed new directions for D2T generation, including using ad-hoc test sets, data in common formats, and reference-free evaluation. By a combination of GPT-4-based metric and human evaluation, we evaluated the performance of LLMs on five domains, providing word-level annotations of model outputs across five domains and recommendations for future directions in D2T generation.

### Acknowledgements

This work was funded by the European Union (ERC, NG-NLG, 101039303) and Charles University project SVV 260 698. It used resources of the LINDAT/CLARIAH-CZ Research Infrastructure (Czech Ministry of Education, Youth, and Sports project No. LM2018101).

### Limitations

In our work, we do not include a comparison to other D2T generation approaches. The main reason is that our benchmark is reference-free, while a large majority of prior approaches are based on models finetuned on reference outputs. However, we believe that our work still satisfies our main

goal of providing insights into behaviors of open LLM models on D2T generation.

We acknowledge that reference-free metrics currently have various shortcomings, including reliance on closed models and specific human annotation protocols, leading to limited replicability and a high price of execution. The approaches occasionally produce incorrect outcomes themselves and they have only moderate correlations with each other. We believe that these shortcomings will be addressed in the future with open model-based metrics.

Our choice of models is limited to 7B-parameter open LLMs due to our limited computational resources. Also, unlike some other LLMs such as GPT-Neo (Black et al., 2022) or BLOOM (Big-Science Workshop et al., 2022), the models we used do not disclose the data they were trained on. For this reason, we find it ever more important to test the models on benchmarks whose labels could have *not* been included in their training data.

The approaches based on LLMs may produce factually and semantically incorrect information. Any text produced by the LLMs therefore needs to be carefully examined, and no decisions should be based on the generated text alone. Discovering the *causes* of LLMs’ hallucinations is out of scope of this paper, but is currently a major topic under investigation (Ji et al., 2023).

### Ethical Considerations

The human evaluation study was approved by the internal ethics committee of our institution. The annotators were hired over Prolific and paid the platform-recommended wage of 9 GBP/hour. The annotators were preselected based on their primary language (English) and their country of residence (US, UK, Ireland, Australia, New Zealand). All annotators were shown detailed instructions and explanation of the data types, data sources, and the purpose of the research (see Appendix B for details). The domains in QUINTD were selected so that they do not contain any sensitive or potentially offensive content. We do not collect any demographic data about the participants.

### 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, NAACL-HLT 2021, pages 3554–3565, Online.
- Rachith Aiyappa, Jisun An, Haewoon Kwak, and Yongyeol Ahn. 2023. [Can we trust the evaluation on ChatGPT?](#) In *Proceedings of the 3rd Workshop on Trustworthy Natural Language Processing (TrustNLP 2023)*, pages 47–54, Toronto, Canada.
- Gabor Angeli, Percy Liang, and Dan Klein. 2010. [A Simple Domain-Independent Probabilistic Approach to Generation](#). In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing, EMNLP 2010, 9-11 October 2010, MIT Stata Center, Massachusetts, USA, A meeting of SIG-DAT, a Special Interest Group of the ACL*, pages 502–512.
- Agnes Axelsson and Gabriel Skantze. 2023. [Using Large Language Models for Zero-Shot Natural Language Generation from Knowledge Graphs](#). *CoRR*, abs/2307.07312.
- Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2023. [LongBench: A bilingual, multitask benchmark for long context understanding](#). *CoRR*, abs/2308.14508.
- 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 Conference of the Association for Computational Linguistics, ACL 2019, Volume 1: Long Papers*, pages 831–844, Florence, Italy.
- Simone Balloccu, Patrícia Schmidtová, Mateusz Lango, and Ondrej Dušek. 2024. [Leak, cheat, repeat: Data contamination and evaluation malpractices in closed-source LLMs](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2024 - Volume 1: Long Papers*, pages 67–93, St. Julian’s, Malta.
- Edward Beeching, Clémentine Fourier, Nathan Habib, Sheon Han, Nathan Lambert, Nazneen Rajani, Omar Sanseviero, Lewis Tunstall, and Thomas Wolf. 2023. [Open LLM Leaderboard](https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard). [https://huggingface.co/spaces/HuggingFaceH4/open\\_llm\\_leaderboard](https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard).
- Anja Belz. 2005. [Corpus-driven generation of weather forecasts](#). In *Proc. 3rd Corpus Linguistics Conference*.
- Anja Belz. 2008. [Automatic generation of weather forecast texts using comprehensive probabilistic generation-space models](#). *Nat. Lang. Eng.*, 14(4):431–455.
- BigScience Workshop, Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, et al. 2022. [Bloom: A 176b-parameter open-access multilingual language model](#). *arXiv preprint arXiv:2211.05100*.
- Sidney Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, et al. 2022. [GPT-NeoX-20B: An open-source autoregressive language model](#). In *Proceedings of BigScience Episode# 5–Workshop on Challenges & Perspectives in Creating Large Language Models*, pages 95–136.
- Armand Bosch. 2019. [WikiDataSets : Standardized sub-graphs from WikiData](#). *CoRR*, abs/1906.04536.
- 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.
- Asli Celikyilmaz, Elizabeth Clark, and Jianfeng Gao. 2020. [Evaluation of text generation: A survey](#). *CoRR*, abs/2006.14799.
- Lingjiao Chen, Matei Zaharia, and James Zou. 2023. [How is ChatGPT’s behavior changing over time?](#) *CoRR*, abs/2307.09009.
- 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 2020*, pages 8635–8648, Online.
- 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*, pages 15607–15631, Toronto, Canada.
- Liam Cripwell, Anya Belz, Claire Gardent, Albert Gatt, Claudia Borg, Marthese Borg, John Judge, Michela Lorandi, Anna Nikiforovskaya, and William Soto Martinez. 2023. [The 2023 WebNLG Shared Task on Low Resource Languages. Overview and Evaluation Results \(WebNLG 2023\)](#). In *Proceedings of the Workshop on Multimodal, Multilingual Natural Language Generation and Multilingual WebNLG Challenge (MM-NLG 2023)*, pages 55–66, Prague, Czech Republic.
- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. [FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness](#). In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural*

- Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA.*
- Seniz Demir, Sandra Carberry, and Kathleen F. McCoy. 2008. [Generating Textual Summaries of Bar Charts](#). In *INLG 2008 - Proceedings of the Fifth International Natural Language Generation Conference, June 12-14, 2008, Salt Fork, Ohio, USA*.
- Seniz Demir, Sandra Carberry, and Kathleen F. McCoy. 2012. [Summarizing Information Graphics Textually](#). *Comput. Linguistics*, 38(3):527–574.
- Ondrej 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, INLG 2019*, pages 421–426, Tokyo, Japan.
- Ondrej Dušek, Jekaterina Novikova, and Verena Rieser. 2020. [Evaluating the state-of-the-art of end-to-end natural language generation: The E2E NLG challenge](#). *Comput. Speech Lang.*, 59:123–156.
- Jinlan Fu, See-Kiong Ng, Zhengbao Jiang, and Pengfei Liu. 2023. [GPTScore: Evaluate as You Desire](#). *CoRR*, abs/2302.04166.
- Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez-Beltrachini. 2017. [The WebNLG Challenge: Generating Text from RDF Data](#). In *Proceedings of the 10th International Conference on Natural Language Generation, INLG 2017, Santiago de Compostela*, pages 124–133, Spain.
- Albert Gatt and Emiel Krahmer. 2018. [Survey of the State of the Art in Natural Language Generation: Core tasks, applications and evaluation](#). *J. Artif. Intell. Res.*, 61:65–170.
- Sebastian Gehrmann, Tosin P. Adewumi, Karmanya Aggarwal, Pawan Sasanka Ammanamanchi, Aremu Anuoluwapo, Antoine Bosselut, Khyathi Raghavi Chandu, Miruna-Adriana Clinciu, Dipanjan Das, Kaustubh D. Dhole, Wanyu Du, Esin Durmus, Ondrej Dušek, Chris Emezue, Varun Gangal, Cristina Garbacea, Tatsunori Hashimoto, Yufang Hou, Yacine Jernite, Harsh Jhamtani, Yangfeng Ji, Shailza Jolly, Dhruv Kumar, Faisal Ladhak, Aman Madaan, Mounica Maddela, Khyati Mahajan, Saad Mahamood, Bodhisattwa Prasad Majumder, Pedro Henrique Martins, Angelina McMillan-Major, Simon Mille, Emiel van Miltenburg, Moin Nadeem, Shashi Narayan, Vitaly Nikolaev, Rubungo Andre Niyongabo, Salomey Osei, Ankur P. Parikh, Laura Perez-Beltrachini, Niranjan Ramesh Rao, Vikas Raunak, Juan Diego Rodriguez, Sashank Santhanam, João Sedoc, Thibault Sellam, Samira Shaikh, Anastasia Shimorina, Marco Antonio Sobrevilla Cabezudo, Hendrik Strobelt, Nishant Subramani, Wei Xu, Diyi Yang, Akhila Yerukola, and Jiawei Zhou. 2021. [The GEM benchmark: Natural language generation, its evaluation and metrics](#). *CoRR*, abs/2102.01672.
- Sebastian Gehrmann, Elizabeth Clark, and Thibault Sellam. 2023. [Repairing the Cracked Foundation: A Survey of Obstacles in Evaluation Practices for Generated Text](#). *J. Artif. Intell. Res.*, 77:103–166.
- Shahriar Golchin and Mihai Surdeanu. 2023. [Time Travel in LLMs: Tracing Data Contamination in Large Language Models](#). *CoRR*, abs/2308.08493.
- Ari Holtzman, Peter West, and Luke Zettlemoyer. 2023. [Generative Models as a Complex Systems Science: How can we make sense of large language model behavior?](#) *CoRR*, abs/2308.00189.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea Madotto, and Pascale Fung. 2023. [Survey of hallucination in natural language generation](#). *ACM Comput. Surv.*, 55(12):248:1–248:38.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. [Mistral 7B](#). *CoRR*, abs/2310.06825.
- Shankar Kantharaj, Rixie Tiffany Ko Leong, Xiang Lin, Ahmed Masry, Megh Thakkar, Enamul Hoque, and Shafiq R. 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), ACL 2022*, pages 4005–4023, Dublin, Ireland.
- Zdeněk Kasner, Ioannis Konstas, and Ondřej Dušek. 2023. [Mind the Labels: Describing Relations in Knowledge Graphs With Pretrained Models](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2023, Dubrovnik*, pages 2390–2407, Croatia.
- Tom Kocmi and Christian Federmann. 2023a. [GEMBA-MQM: Detecting Translation Quality Error Spans with GPT-4](#). In *Proceedings of the Eighth Conference on Machine Translation, WMT 2023*, pages 768–775, Singapore.
- Tom Kocmi and Christian Federmann. 2023b. [Large Language Models Are State-of-the-Art Evaluators of Translation Quality](#). In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation, EAMT 2023*, pages 193–203, Tampere, Finland.
- 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, NAACL-HLT 2019, Minneapolis, MN, Volume 1 (Long and Short Papers)*, pages 2284–2293, USA.

- Ryan Koo, Minhwa Lee, Vipul Raheja, Jong Inn Park, Zae Myung Kim, and Dongyeop Kang. 2023. **Benchmarking cognitive biases in large language models as evaluators.** *CoRR*, abs/2309.17012.
- Fajri Koto, Jey Han Lau, and Timothy Baldwin. 2022. **Can Pretrained Language Models Generate Persuasive, Faithful, and Informative Ad Text for Product Descriptions?** In *Proceedings of the Fifth Workshop on E-Commerce and NLP (ECNLP 5)*, pages 234–243, Dublin, Ireland.
- 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, EMNLP 2016*, pages 1203–1213, Austin, Texas, USA.
- 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.** *CoRR*, abs/2203.05227.
- Percy Liang, Michael I. Jordan, and Dan Klein. 2009. **Learning Semantic Correspondences with Less Supervision.** In *ACL 2009, Proceedings of the 47th Annual Meeting of the Association for Computational Linguistics and the 4th International Joint Conference on Natural Language Processing of the AFNLP, 2-7 August 2009*, pages 91–99, Singapore.
- Yang Liu, Dan Iter, Yichong Xu, Shuhang Wang, Ruochen Xu, and Chenguang Zhu. 2023. **G-Eval: NLG Evaluation using Gpt-4 with Better Human Alignment.** In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023*, pages 2511–2522, Singapore.
- Michela Lorandi and Anja Belz. 2023. **Data-to-text Generation for Severely Under-Resourced Languages with GPT-3.5: A Bit of Help Needed from Google Translate (WebNLG 2023).** In *Proceedings of the Workshop on Multimodal, Multilingual Natural Language Generation and Multilingual WebNLG Challenge (MM-NLG 2023)*, pages 80–86.
- Tsendsuren Munkhdalai, Manaal Faruqui, and Siddharth Gopal. 2024. **Leave no context behind: Efficient infinite context transformers with infinite attention.** *CoRR*, abs/2404.07143.
- 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.
- Jekaterina Novikova, Ondřej Dušek, Amanda Cercas Curry, and Verena Rieser. 2017. **Why We Need New Evaluation Metrics for NLG.** In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017*, pages 2241–2252, Copenhagen, Denmark.
- 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, INLG 2020*, pages 138–147, Dublin, Ireland.
- OpenAI. 2023a. **GPT-4 Technical Report.** *CoRR*, abs/2303.08774.
- OpenAI. 2023b. **Introducing ChatGPT.** <https://openai.com/blog/chatgpt>. Accessed on January 9, 2024.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022. **Training language models to follow instructions with human feedback.** In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022*, New Orleans, LA, USA.
- Ratish Puduppully, Li Dong, and Mirella Lapata. 2019a. **Data-to-Text Generation with Content Selection and Planning.** 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*, pages 6908–6915, Honolulu, Hawaii, USA.
- Ratish Puduppully, Li Dong, and Mirella Lapata. 2019b. **Data-to-text Generation with Entity Modeling.** In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Volume 1: Long Papers*, pages 2023–2035, Florence, Italy.
- Ratish Puduppully, Yao Fu, and Mirella Lapata. 2022. **Data-to-text Generation with Variational Sequential Planning.** *Trans. Assoc. Comput. Linguistics*, 10:697–715.
- Ratish Puduppully and Mirella Lapata. 2021. **Data-to-text Generation with Macro Planning.** *Trans. Assoc. Comput. Linguistics*, 9:510–527.
- Clément Rebuffel, Laure Soulier, Geoffrey Scouteeten, 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, Lisbon, Portugal, April 14-17, 2020, Proceedings, Part I*, volume 12035 of *Lecture Notes in Computer Science*, pages 65–80.
- Ehud Reiter. 2023. **We should evaluate real-world impact!** <https://ehudreiter.com/2023/11/13/evaluate-real-world-impact/>. Accessed on January 11, 2024.
- Ehud Reiter and Robert Dale. 1997. **Building applied natural language generation systems.** *Nat. Lang. Eng.*, 3(1):57–87.

- Leonardo F. R. Ribeiro, Martin Schmitt, Hinrich Schütze, and Iryna Gurevych. 2020. [Investigating Pretrained Language Models for Graph-to-Text Generation](#). *CoRR*, abs/2007.08426.
- Anna Rogers. 2023. Closed AI Models Make Bad Baselines. <https://hackingsemantics.xyz/2023/closed-baselines/>. Accessed on January 11, 2024.
- Huajie Shao, Jun Wang, Haohong Lin, Xuezhou Zhang, Aston Zhang, Heng Ji, and Tarek F. Abdelzaher. 2021. [Controllable and Diverse Text Generation in E-commerce](#). In *WWW '21: The Web Conference 2021*, pages 2392–2401, Virtual Event / Ljubljana, Slovenia.
- Mandar Sharma, John S. Brownstein, and Naren Ramakrishnan. 2021. [TCube: Domain-Agnostic Neural Time-series Narration](#). *CoRR*, abs/2110.05633.
- Andrea Sottana, Bin Liang, Kai Zou, and Zheng Yuan. 2023. [Evaluation Metrics in the Era of GPT-4: Reliably Evaluating Large Language Models on Sequence to Sequence Tasks](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023*, pages 8776–8788, Singapore.
- Somayajulu Sripada, Ehud Reiter, Jim Hunter, and Jin Yu. 2002. Sumtime-meteo: Parallel corpus of naturally occurring forecast texts and weather data. *Computing Science Department, University of Aberdeen, Aberdeen, Scotland, Tech. Rep. AUCS/TR0201*.
- Rickard Stureborg, Dimitris Alikaniotis, and Yoshi Suhara. 2024. [Large Language Models are Inconsistent and Biased Evaluators](#).
- 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, INLG 2020*, pages 158–168, Dublin, Ireland.
- Craig Thomson, Ehud Reiter, and Somayajulu Sripada. 2021. SportSett:Basketball - A Robust and Maintainable Dataset for Natural Language Generation. page 9.
- Craig Thomson, Ehud Reiter, and Barkavi Sundararajan. 2023. [Evaluating factual accuracy in complex data-to-text](#). *Comput. Speech Lang.*, 80:101482.
- TogetherAI. 2023. Preparing for the era of 32K context: Early learnings and explorations. <https://www.together.ai/blog/llama-2-7b-32k>. Accessed on January 2, 2024.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. [LLaMA: Open and Efficient Foundation Language Models](#). *CoRR*, abs/2302.13971.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. [Llama 2: Open Foundation and Fine-Tuned Chat Models](#). *CoRR*, abs/2307.09288.
- Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M. Rush, and Thomas Wolf. 2023. [Zephyr: Direct Distillation of LM Alignment](#). *CoRR*, abs/2310.16944.
- Jannis Vamvas and Rico Sennrich. 2022. [As little as possible, as much as necessary: Detecting over- and undertranslations with contrastive conditioning](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), ACL 2022*, pages 490–500, Dublin, Ireland.
- Chris van der Lee, Albert Gatt, Emiel van Miltenburg, and Emiel Kraemer. 2021. [Human evaluation of automatically generated text: Current trends and best practice guidelines](#). *Comput. Speech Lang.*, 67:101151.
- Chris van der Lee, Albert Gatt, Emiel van Miltenburg, Sander Wubben, and Emiel Kraemer. 2019. [Best practices for the human evaluation of automatically generated text](#). In *Proceedings of the 12th International Conference on Natural Language Generation, INLG 2019*, pages 355–368, Tokyo, Japan.
- 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.
- Hongmin Wang. 2019. [Revisiting Challenges in Data-to-Text Generation with Fact Grounding](#). In *Proceedings of the 12th International Conference on Natural Language Generation, INLG 2019*, pages 311–322, Tokyo, Japan.

- Jiaan Wang, Yunlong Liang, Fandong Meng, Haoxiang Shi, Zhixu Li, Jinan Xu, Jianfeng Qu, and Jie Zhou. 2023a. *Is ChatGPT a Good NLG Evaluator? A Preliminary Study*. *CoRR*, abs/2303.04048.
- Jinpeng Wang, Yutai Hou, Jing Liu, Yunbo Cao, and Chin-Yew Lin. 2017. *A Statistical Framework for Product Description Generation*. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing, IJCNLP 2017, Volume 2: Short Papers*, pages 187–192, Taipei, Taiwan.
- Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. 2023b. *Large Language Models are not Fair Evaluators*. *CoRR*, abs/2305.17926.
- Tsung-Hsien Wen, Milica Gašić, Nikola Mrkšić, Lina M Rojas-Barahona, Pei-Hao Su, David Vandyke, and Steve Young. 2015. *Toward multi-domain language generation using recurrent neural networks*. In *NIPS Workshop on Machine Learning for Spoken Language Understanding and Interaction*.
- Tsung-Hsien Wen, Milica Gasic, Nikola Mrksic, Lina Maria Rojas-Barahona, Pei-Hao Su, David Vandyke, and Steve J. Young. 2016. *Multi-domain Neural Network Language Generation for Spoken Dialogue Systems*. In *NAACL HLT 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 120–129, San Diego California, USA.
- Sam Wiseman, Stuart M. Shieber, and Alexander M. Rush. 2017. *Challenges in Data-to-Document Generation*. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017*, pages 2253–2263, Copenhagen, Denmark.
- 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, EMNLP 2020 - Demos*, pages 38–45, Online.
- Wenda Xu, Danqing Wang, Liangming Pan, Zhenqiao Song, Markus Freitag, William Yang Wang, and Lei Li. 2023. *INSTRUCTSCORE: Explainable Text Generation Evaluation with Finegrained Feedback*.
- Shuzhou Yuan and Michael Färber. 2023. *Evaluating Generative Models for Graph-to-Text Generation*. In *Proceedings of the 14th International Conference on Recent Advances in Natural Language Processing, RANLP 2023*, pages 1256–1264, Varna, Bulgaria.
- Yilun Zhao, Haowei Zhang, Shengyun Si, Linyong Nan, Xiangru Tang, and Arman Cohan. 2023. *Investigating Table-to-Text Generation Capabilities of LLMs in Real-World Information Seeking Scenarios*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. *Judging LLM-as-a-judge with MT-Bench and Chatbot Arena*. *CoRR*, abs/2306.05685.

## A QUINTD Data

Here, we describe the data sources we include in the QUINTD collection tool and the procedure of collecting the QUINTD-1 benchmark. To replicate the data collection, please refer to the scripts we provide.<sup>14</sup>

### A.1 Selection of Data Sources

When selecting the data sources, we had the following desiderata:

- Data needs to be publicly available.
- Data needs to represent a common data-to-text task.
- Data needs to be in a common format (or straightforwardly transformable to one).

We settled on the data sources described in Appendix A.2. All the sources can be accessed using an API. Note that some of the APIs have access limits, either for the requests made from a single account per day or for a number of requests from an IP address within a time window. However, these limits do not severely limit the data collection process on the scale we use here.

### A.2 Data Collection

Table 6 summarizes the output types for each domain.

Domain Id	Output type
openweather	five-day weather forecast
gsmarena	product description
ice_hockey	ice hockey game summary
owid	chart caption
wikidata	entity description

Table 6: The output types for individual domains in QUINTD.

<sup>14</sup><https://github.com/kasnerz/quintd>

### A.2.1 OpenWeather

OpenWeather ([OpenWeatherMap.org](https://openweathermap.org)) is an online service that provides global weather data via web interface and API. The API responses are in the JSON format [documented](#) at the official website. For our experiments, we used the [forecast5](#) API, which allows to download a 5-day forecast with 3-hour resolution for any location specified by its GPS coordinates.

The free tier is limited to 1,000 API calls per day, which is enough to download our whole test set in one bulk. However, at the time of experiments, the free API only allowed to download the data for the *time when the request was made*. At the time of writing, OpenWeather is pushing a new [One Call API 3.0](#) which allows to download weather data for any timestamp, but only *4 days ahead* (instead of 5). These restrictions somehow limit the replicability of our QUINTD-1 dataset (at least with the free API) but do not limit downloading a new batch of data with a similar format.

For the QUINTD-1 dataset, we randomly sampled 100 cities for each split from the [list of cities with a population over 1000](#) and used their coordinates in the queries to OpenWeather API. All the data forecasts were downloaded on Jan 3, 2024.

### A.2.2 GSMArena

[GSMArena](#) is a website providing specifications and reviews for mobile devices. For downloading the data, we used the unofficial [gsmarena-api](#) tool, which returns the data in a JSON format. Note that GSMArena imposes limitations on the number of requests per IP address, which may induce delays when downloading a larger amount of data.

To create a balanced sample, we downloaded detailed specifications of 10 products from each available brand and randomly selected 100 products for each split from the downloaded set.

### A.2.3 RapidAPI Ice Hockey

[RapidAPI](#) is a service that provides API access to data from multiple domains, including sport, finance, entertainment, and others. Most APIs are provided in a freemium mode, i.e., with a limited number of daily API calls.

For QUINTD, we selected the [IceHockeyAPI](#) (popularity 9.1 / 10), which provides access to ice hockey games from world top leagues. Our choice was influenced by our own personal preferences, combined with the desire to cover a sport that has

not been covered previously in sports report generation.

We used the [matches](#) endpoint which returns high-level details about a game. Note that the API allows only 50 requests per day, but that does not limit the data collection since the endpoint returns *all the games* played on a particular day in a single request. We downloaded the games played on 27 November 2023 for the development set (184 games) and 29 November 2023 for the test set (216 games), taking a random sample of 100 for each split.

### A.2.4 OurWorldInData

[OurWorldInData](#) is a public database and web interface for data about world developments in various domains and sources. We used the official API (currently experimental), which is accessible through the Python package [owid-catalog](#). The package allows accessing individual CSV tables as Pandas dataframes.

For our data collection, we decided to limit ourselves to time series, i.e., a single column with values changing over time. Besides the simplicity of visualizing such a chart (which is used by human annotators for checking the correctness of the output), there is also a clear goal for the target chart description: describing the developments of a value over time. We also limited ourselves to the health domain. In particular, we selected the tables [COVID data](#) (columns `new_cases_smoothed_per_million`, `new_tests_smoothed_per_thousand`, `people_vaccinated_per_hundred`, `reproduction_rate`, and `positive_rate`) and [Life expectancy data](#) (column `life_expectancy_0`).

We downloaded the data for all countries with non-empty entries in the table, taking a random sample of 100 examples for each split. On model input, we formatted the data for each time series as a two-column CSV, including the title, the description, and the unit for each example as a comment (#) at the beginning of the input.

### A.2.5 Wikidata

[Wikidata](#) is a large open-source knowledge graph containing factual information about entities and their properties. Wikidata provides access through an [official API](#), but we instead decided to extract our data using the [wikidatasets](#) ([Boschin, 2019](#)) Python library, which provides access to preprocessed properties of entities from particular do-

mains. It allowed us to avoid crawling and filtering the knowledge graph, and its offline processing made the data collection faster.<sup>15</sup>

For our dataset, we selected the entities from the companies, countries, films, and humans domains. For each entity, we randomly extracted between 2 to 10 properties in the knowledge graph. We extracted up to 100 subgraphs for each domain and took a random sample of 100 subgraphs for each split. On model input, we formatted each subgraph as a simple Markdown-formatted text snippet, using the entity as a title and including a bullet point for each key-value pair.

## B Human Evaluation

As described in §4.1, we set up the human evaluation campaign on Prolific. To make the data more accessible to the annotators, we created custom data visualizations for each domain. For the data in openweather and owid, we used interactive graphs from [Highcharts.com](https://highcharts.com), and we manually created the tables for other domains. You can find the full instructions for human annotators in [Figure 5](#) and the examples of data visualizations in [Appendix E](#).

## C GPT-4 Evaluation

We used the prompt in [Figure 4](#) for instantiating the GPT-4-based metric.<sup>16</sup> We set the temperature to 0 to improve the replicability of the process. We ensured that the output is a valid JSON using the parameter `response_format` in the [OpenAI API](#). At the price of \$0.01 per 1k input tokens and \$0.03 per 1k generated tokens, the evaluation process costs approximately \$45 in total.

### C.1 Aligning the Errors

For aligning the errors with the original text, we perform string matching on the text span decoded by GPT-4 in the `TEXT_SPAN` field. In our preliminary experiments, this method proved to be more robust than either asking for start and end indices of the error span (which would rely on the model’s ability to count characters) or performing sequence tagging on the copy of the input (which would rely on the model’s ability to perfectly copy the input).

<sup>15</sup>All the entities and properties are linked with an identifier to the Wikidata database, making the process also replicable through the official API.

<sup>16</sup>Note that the example in the prompt differs from the example used for human annotators (see [Figure 5](#)). We revised the example to be more instructive, but we were not able to re-run the GPT-4 evaluation due to our limited budget.

We tried to respect the monotonic ordering of text spans but fell back to full-text search if the span is not found following the previous one. We consider this approach successful since matching completely failed only in a minority of cases (137 out of 6927). Based on our manual examination, these mostly include cases where GPT-4 tried to suggest a *missing* piece of text as an error or did not manage to copy the input text verbatim.

## D Experiments with Open LLMs as Evaluators

To select the most suitable LLM for automatic evaluation, we compared correlations with human judgment of the following models:

- **GPT-4** ([OpenAI, 2023a](#)) used via OpenAI API (`gpt-4-1106-preview`),
- **GPT-3.5** ([OpenAI, 2023b](#)) used via OpenAI API (`gpt-3.5-turbo-1106`),
- **Llama-3-70B**<sup>17</sup> running locally via [Ollama](#) in 4-bit quantization (`meta-llama/Meta-Llama-3-70B`).

We used all the models with the same prompts, temperature 0, and force-decoded JSON outputs. In [Table 7](#), we show Pearson correlation coefficients computed as described in §5.2. We can see that the strongest model is GPT-4, followed by Llama-3-70B and GPT-3.5. As the gap between the models is substantial, we opted for using GPT-4 which is the strongest model.

model	$r_{\text{word}}$	$r_{\text{example}}$	$r_{\text{domain}}$
GPT-4	<b>0.26</b>	<b>0.52</b>	<b>0.93</b>
GPT-3.5	0.07	0.33	0.82
Llama3-70B	0.09	0.44	0.92

Table 7: Pearson correlation coefficients of the model annotations as compared with human annotations (cf. §5.2).

## E Examples

Here, we present an example of inputs and model outputs (along with annotations) for each domain:

- openweather: [Figure 7](#) (in) and [Table 8](#) (out),
- gsmarena: [Figure 8](#) (in) and [Table 8](#) (out),

<sup>17</sup><https://llama.meta.com/llama3/>



- ice\_hockey: [Figure 9](#) (in) and [Table 10](#) (out),
- owid: [Figure 10](#) (in) and [Table 11](#) (out),
- wikidata: [Figure 11](#) (in) and [Table 12](#) (out).

Note that the graphs for openweather and owid are interactive when accessed through the web interface.

## **F Full Results**

Here, we include the tables with results for individual domains:

- [Table 13](#) presents the average *numbers of errors per output* separately for each domain (the aggregated results are in [Table 3](#)),
- [Table 14](#) presents the ratio of *outputs containing at least one error* separately for each domain (the aggregated results are in [Table 4](#)).

### System Message

You are an expert data-to-text error annotation system. You understand structured data and you can correctly operate with units and numerical values. You are designed to output token-level annotations in JSON.

### Prompt

```
Given the data:
```
data
```
Annotate all the errors in the following text:
```
text
```
Output the errors as a JSON list "errors" in which each object contains fields "reason", "text",
and "type". The value of "text" is the text of the error. The value of "reason" is the reason
for the error. The value of "type" is one of 0, 1, 2, 3 based on the following list:
- 0: Incorrect fact: The fact in the text contradicts the data.
- 1: Not checkable: The fact in the text cannot be checked in the data.
- 2: Misleading: The fact in the text is misleading in the given context.
- 3: Other: The text is problematic for another reason, e.g. grammatically or stylistically
incorrect, irrelevant, or repetitive.
The list should be sorted by the position of the error in the text.
*Example:*
data:
```
[[ "Aditi Bhagwat", "occupation", "television actor" ], [ "Aditi Bhagwat", "date of birth", "18
January 1981" ] ]
```
text:
```
Aditi Bhagwat, born on January 18, 1991, used to be a popular Indian television actor. The data
comes from a knowledge graph.
```
output:
```
"errors": [{"reason": "The data mentions that the actor was born on 1981", "text": "1991",
"type": 0, "reason": "Misleadingly suggests that the actor is not alive", "text": "used to
be", "type": 2, "reason": "Popularity is not mentioned in the data", "text": "popular", "type": 1,
"reason": "Nationality is not mentioned in the data", "text": "Indian", "type": 1, "reason": "The
note is superfluous", "text": "The data comes from a knowledge graph.", "type": 3}
```
Note that some details may not be mentioned in the text: do not count omissions as errors. Also
do not be too strict: some facts can be less specific than in the data (rounded values, shortened
or abbreviated text, etc.), do not count these as errors. If there are no errors in the text,
"errors" will be an empty list.
```

Figure 4: The prompt we used for the GPT-4 evaluation metric.

In this task, you will annotate **20** examples in total. For each example, you will see **data** on the left side and the corresponding generated **text** on the right side. Your task is to **annotate errors** in the text with respect to the data. There are five types of errors that you can mark in the generated text:

1. **Incorrect fact<sup>I</sup>**: The fact in the text contradicts the data.
2. **Not checkable<sup>NC</sup>**: The fact in the text cannot be checked given the data.
3. **Misleading<sup>M</sup>**: The fact in the text is misleading in the given context.
4. **Other<sup>O</sup>**: The text is problematic for another reason, e.g. grammatically or stylistically incorrect, irrelevant, or repetitive.

### How to mark and submit the annotations?

Use your mouse to **highlight specific parts of the text** containing the errors. To switch between error categories, repeatedly click on the highlighted text (the last click removes the highlight). Note that highlighting from the right to left can work better for longer spans.

Once you think you have marked all the errors present in the text, click the **Mark example as complete** button (you can still update the annotation later). You will be able to submit the annotations once they are all are marked as complete.

### How should I decide on the errors?

- Each error span should include all the words related to the error (but nothing else).
- If you think the fact is probably true, but cannot be derived from the data, mark it as not checkable.
- If you are not really sure if the fact should be marked as an error, leave it unmarked.

### An example of the annotated output

An example of the data input and the corresponding text annotated with errors:

#### data

##### Nokia 3310

- **color**: black, blue, grey
- **display**: 320x240px

#### text (product description)

Nokia 3310 is **produced in Finland<sup>NC</sup>** and features a **320x320 display<sup>I</sup>**. It is **available in black color<sup>M</sup>**. The data seem to provide only partial information about the phone.<sup>O</sup>

#### explanation

- **produced in Finland**: The country where the phone is produced is not mentioned in the data.
- **320x320**: The data mentions that the display has resolution 320x240px.
- **available in black color**: Misleadingly suggests that the phone is not available in other colors.
- **The data seem to provide only partial information about the phone.**: The note is irrelevant for the phone description.

*Note that for the sake of brevity, this particular example is based on a small data input and contains many errors, which may not always be the case.*

### What kinds of data and text can I encounter?

You can encounter the following kinds of texts:

- a 5-day weather forecast generated from weather data,
- a description of a product generated from product specifications
- an ice hockey game report generated from information about the game,
- a caption of a health-related chart,
- a description of an entity (human, country, film, or company) based on its properties.

### What is the source of the data and the texts?

The data is downloaded from public sources ([openweathermap.org](https://openweathermap.org), [gsmarena.com](https://gsmarena.com), [rapidapi.com](https://rapidapi.com), [ourworldindata.com](https://ourworldindata.com), [wikidata.com](https://wikidata.com)). The texts are generated by large language models.

### What is the purpose of the research?

The research is conducted by the NLP research group at Charles University in Prague. The research is a part of the ERC project [Next-Generation Natural Language Generation](#) and was approved by the ethical committee of the project.

By continuing, you agree that the anonymous annotations will be made available to the research community to help evaluating and improving data-to-text generation systems based on large language models.

Figure 5: The instructions given to the human annotators.

### New Jersey Devils – New York Islanders

**metadata**

id	11385979
startDatetime	2023-11-29 01:00:00
startTimestamp	1701216000

**season**

name	NHL
year	23/24
id	52528

**homeTeam**

name	New Jersey Devils
shortName	Devils
nameCode	NJ
id	3704
country	USA

**homeScore**

display	5
period1	1
period2	1
period3	3
normaltime	5

**awayScore**

display	4
period1	1
period2	3
period3	0
normaltime	4

**awayTeam**

name	New York Islanders
shortName	Islanders
nameCode	NYI
id	3703
country	USA

**Instructions**

The following text is a(n) **ice hockey game summary**. Please annotate all the errors in the text:

In the NHL 23/24 season, the New Jersey Devils hosted the New York Islanders at their home arena. The game was closely contested, with both teams **scoring frequently in the first and second periods**. The Devils took a narrow lead into the final frame, but the Islanders fought back hard, scoring three unanswered goals in the third period to secure a 5-4 victory. The final scoreline was **5-4 in favor of the Islanders**, with the game ending in normal time. The matchup between the Devils and Islanders **was intense and exciting, with both teams showcasing their skills and determination throughout the game**.

After annotating the errors in the text, please also **check if you agree with any of the following statements**, then mark the example as complete:

- The text is severely off-topic (seems completely unrelated to the data).
- The text is missing or incomplete.
- The text is 100% accurate and does not contain any errors.

✔ Mark example as complete

Figure 6: The annotation interface for human annotators.

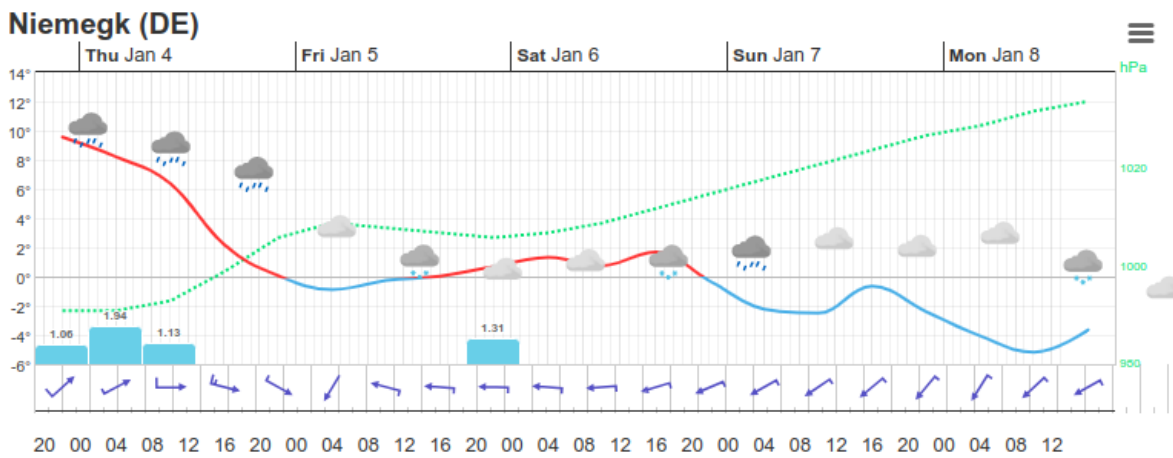


Figure 7: An example of an input from openweather (see the outputs in Table 8).

Model	Human annotations ( $\mathcal{E}_{\text{hum}}$ )	GPT-4 annotations ( $\mathcal{E}_{\text{gpt}}$ )
Llama2	The weather in Niemegek, Germany, will be cloudy and rainy for the next few days. On January 3rd, the temperature will be around 9.55°C, with a feels-like temperature of 6.59°C. The weather will be light rain, with a wind speed of 6.24 m/s and a gust of 12.1 m/s. <u>The rain will last for 3 hours<sup>NC</sup></u> , and the pressure will be 991 hPa.	The weather in Niemegek, Germany, will be cloudy and rainy for the next few days. On January 3rd, the temperature will be around 9.55°C, with a feels-like temperature of 6.59°C. The weather will be light rain, with a wind speed of 6.24 m/s and a gust of 12.1 m/s. <u>The rain will last for 3 hours<sup>M</sup></u> , and the pressure will be 991 hPa.
Mistral	From January 3rd to January 7th, the city of Niemegek in Germany will experience a range of weather conditions. On January 3rd, the temperature will be 9.55°C with a feels like temperature of 6.59°C. There will be light rain with a <u>3-hour precipitation<sup>NC</sup></u> of 1.06mm. The wind speed will be 6.24m/s with a gust of 12.1m/s. On January 4th, the temperature will be 8.2°C with a feels like temperature of 5.64°C. There will be light rain with a <u>3-hour precipitation<sup>I</sup></u> of 1.94mm. The wind speed will be 4.31m/s with a gust of 11.6m/s. On January 5th, the temperature will be <u>6.37°C with a feels like temperature of 2.25°C<sup>I</sup></u> . There will be light rain with a <u>3-hour precipitation of 1.13mm<sup>I</sup></u> . <u>The wind speed will be 7.07m/s with a gust of 15.51m/s<sup>I</sup></u> . On January 6th, the temperature will be <u>2.17°C<sup>I</sup></u> with a <u>feels like temperature of -3.55°C<sup>I</sup></u> . There will be overcast clouds with a <u>3-hour precipitation of 0.31mm<sup>I</sup></u> . The wind speed will be <u>3.54m/s with a gust of 7.11m/s<sup>I</sup></u> . On January 7th, the temperature will be <u>-0.89°C with a feels like temperature of -0.89°C<sup>I</sup></u> . There will be overcast clouds with a <u>3-hour precipitation of 0.31mm. The wind speed will be 0.63m/s with a gust of 0.7m/s<sup>I</sup></u> .	From January 3rd to January 7th, the city of Niemegek in Germany will experience a range of weather conditions. On January 3rd, the temperature will be 9.55°C with a feels like temperature of 6.59°C. There will be light rain with a 3-hour precipitation of 1.06mm. The wind speed will be 6.24m/s with a gust of 12.1m/s. On January 4th, the temperature will be 8.2°C with a feels like temperature of 5.64°C. There will be light rain with a 3-hour precipitation of 1.94mm. The wind speed will be 4.31m/s with a gust of 11.6m/s. On January 5th, the temperature will be 6.37°C with a feels like temperature of 2.25°C. There will be light rain with a 3-hour precipitation of <u>1.13mm<sup>I</sup></u> . The wind speed will be 7.07m/s with a gust of 15.51m/s. On January 6th, the temperature will be 2.17°C with a feels like temperature of -3.55°C. There will be overcast clouds with a 3-hour precipitation of <u>0.31mm<sup>I</sup></u> . The wind speed will be <u>3.54m/s<sup>I</sup></u> with a gust of <u>7.11m/s<sup>I</sup></u> . On January 7th, the temperature will be <u>-0.89°C<sup>I</sup></u> with a feels like temperature of <u>-0.89°C<sup>I</sup></u> . There will be overcast clouds with a 3-hour precipitation of <u>0.31mm<sup>I</sup></u> . The wind speed will be <u>0.63m/s<sup>I</sup></u> with a gust of <u>0.7m/s<sup>I</sup></u> .
Zephyr	Over the next five days, the weather in Niemegek, Germany will be characterized by fluctuating temperatures and varying precipitation. On Thursday, January <u>3<sup>I</sup></u> , light rain and a <u>low<sup>I</sup></u> temperature of 9.55°C can be expected, with a chance of rain continuing into Friday. Snowfall is possible on Friday, with a low temperature of -0.89°C. Saturday will be overcast with a low temperature of <u>0.03°C<sup>I</sup></u> , and <u>light rain is forecasted<sup>I</sup></u> for Saturday evening. Sunday will be overcast as well, with a low temperature of <u>-2.23°C<sup>I</sup></u> . Monday will see a drop in temperatures, with a low of <u>-4.07°C<sup>I</sup></u> and scattered clouds. <u>Temperatures will continue to fall throughout the week, with lows of -5.19°C and -9.51°C on Tuesday and Wednesday, respectively<sup>NC</sup></u> . Wind speeds will vary from 3.1 to <u>15.51<sup>I</sup></u> meters per second, with gusts reaching up to <u>15.51<sup>I</sup></u> meters per second on Monday.	Over the next five days, the weather in Niemegek, Germany will be characterized by fluctuating temperatures and varying precipitation. On <u>Thursday, January 3<sup>I</sup></u> , light rain and a <u>low temperature of 9.55°C<sup>I</sup></u> can be expected, with a chance of rain continuing into Friday. Snowfall is possible on Friday, with a <u>low temperature of -0.89°C<sup>I</sup></u> . Saturday will be overcast with a low temperature of 0.03°C, and <u>light rain is forecasted for Saturday evening<sup>I</sup></u> . Sunday will be overcast as well, with a <u>low temperature of -2.23°C<sup>I</sup></u> . Monday will see a drop in temperatures, with a <u>low of -4.07°C<sup>I</sup></u> and scattered clouds. Temperatures will continue to fall throughout the week, with <u>lows of -5.19°C and -9.51°C on Tuesday and Wednesday<sup>I</sup></u> , respectively. Wind speeds will vary from 3.1 to 15.51 meters per second, with <u>gusts reaching up to 15.51 meters per second on Monday<sup>I</sup></u> .
GPT-3.5	Over the next five days, the weather in Niemegek, Germany is expected to be quite varied. On January 4th, there will be light rain with temperatures ranging from <u>2.17°C to 9.55°C<sup>I</sup></u> . The following day, January 5th, there will be <u>light<sup>I</sup></u> snow with temperatures dropping to as low as -0.89°C. On January 6th, the weather will be overcast with temperatures ranging from -0.32°C to 1.69°C. The trend continues on January 7th with overcast clouds and temperatures dropping to as low as <u>-2.51°C<sup>I</sup></u> . Finally, on January 8th, there will be scattered clouds with temperatures ranging from -5.19°C to -3.67°C. Throughout this period, wind speeds will vary between 0.63 m/s and 8.12 m/s, and the pressure will range from <u>991<sup>I</sup></u> hPa to 1034 hPa.	Over the next five days, the weather in Niemegek, Germany is expected to be quite varied. On January 4th, there will be light rain with temperatures ranging from 2.17°C to <u>9.55°C<sup>I</sup></u> . The following day, January 5th, there will be light snow with temperatures dropping to as low as <u>-0.89°C<sup>I</sup></u> . On January 6th, the weather will be overcast with temperatures ranging from -0.32°C to <u>1.69°C<sup>I</sup></u> . The trend continues on January 7th with overcast clouds and temperatures dropping to as low as <u>-2.51°C<sup>I</sup></u> . Finally, on January 8th, there will be scattered clouds with temperatures ranging from -5.19°C to <u>-3.67°C<sup>I</sup></u> . Throughout this period, wind speeds will vary between 0.63 m/s and <u>8.12 m/s<sup>I</sup></u> , and the pressure will range from <u>991 hPa<sup>I</sup></u> to <u>1034 hPa<sup>I</sup></u> .

Table 8: An example of the model outputs on openweather (see the input in Table 8).

## Gigabyte GSmart Mika M3

name: GSmart Mika M3, id: gigabyte\_gsmart\_mika\_m3-6797

quick specifications

<b>Display size</b>	5.0"
<b>Display resolution</b>	720x1280 pixels
<b>Camera pixels</b>	13MP
<b>Video pixels</b>	Video recorder
<b>RAM size</b>	1GB RAM
<b>Chipset</b>	MT6582
<b>Battery size</b>	1900mAh
<b>Battery type</b>	Li-Ion

detailed specifications

<b>Network</b>	<b>Technology</b>	GSM / HSPA
	<b>2G bands</b>	GSM 850 / 900 / 1800 / 1900 - SIM 1 & SIM 2
	<b>3G bands</b>	HSDPA 900 / 2100
	<b>Speed</b>	HSPA 21.1/5.76 Mbps
<b>Launch</b>	<b>Announced</b>	2014, October. Released 2014, November
	<b>Status</b>	Discontinued
<b>Body</b>	<b>Dimensions</b>	144 x 70.5 x 8.3 mm (5.67 x 2.78 x 0.33 in)
	<b>Weight</b>	146 g (5.15 oz)
	<b>SIM</b>	Dual SIM (Mini-SIM/ Micro-SIM)
<b>Display</b>	<b>Type</b>	IPS LCD
	<b>Size</b>	5.0 inches, 68.9 cm <sup>2</sup> (~67.9% screen-to-body ratio)
	<b>Resolution</b>	720 x 1280 pixels, 16:9 ratio (~294 ppi density)
<b>Platform</b>	<b>OS</b>	Android 4.4.2 (KitKat)
	<b>Chipset</b>	Mediatek MT6582 (28 nm)
	<b>CPU</b>	Quad-core 1.3 GHz Cortex-A7
	<b>GPU</b>	Mali-400MP2
<b>Memory</b>	<b>Card slot</b>	microSDHC (dedicated slot)
	<b>Internal</b>	8GB 1GB RAM
<b>Main Camera</b>	<b>Single</b>	13 MP, AF
	<b>Features</b>	LED flash
	<b>Video</b>	Yes
<b>Selfie camera</b>	<b>Single</b>	8 MP
	<b>Video</b>	
<b>Sound</b>	<b>Loudspeaker</b>	Yes
	<b>3.5mm jack</b>	Yes
<b>Comms</b>	<b>WLAN</b>	Wi-Fi 802.11 b/g/n, hotspot
	<b>Bluetooth</b>	3.0
	<b>Positioning</b>	GPS
	<b>NFC</b>	No
	<b>Radio</b>	FM radio
	<b>USB</b>	microUSB 2.0
<b>Features</b>	<b>Sensors</b>	Accelerometer, proximity
<b>Battery</b>	<b>Type</b>	Li-Ion 1900 mAh, removable
	<b>Stand-by</b>	Up to 150 h
	<b>Talk time</b>	Up to 4 h 30 min
<b>Misc</b>	<b>Colors</b>	Navy Blue

Figure 8: An example of an input from gsmarena (see the outputs in Table 9).

Model	Human annotations ( $\mathcal{E}_{\text{hum}}$ )	GPT-4 annotations ( $\mathcal{E}_{\text{gpt}}$ )
Llama2	The Gigabyte GSmart Mika M3 is a smartphone that was released in 2014. It has a 5-inch IPS LCD display with a resolution of 720x1280 pixels and a 13MP main camera. The device also has an 8MP front camera and is powered by a Mediatek MT6582 chipset with 1GB of RAM and 8GB of internal storage. It has dual SIM support and is available in Navy Blue.	The Gigabyte GSmart Mika M3 is a smartphone that was released in 2014. It has a <b>5-inch<sup>I</sup></b> IPS LCD display with a resolution of 720x1280 pixels and a 13MP main camera. The device also has an <b>8MP front camera<sup>NC</sup></b> and is powered by a Mediatek MT6582 chipset with 1GB of RAM and 8GB of internal storage. It has dual SIM support and <b>is available in Navy Blue<sup>NC</sup></b> .
Mistral	The GSmart Mika M3 is a smartphone developed by Gigabyte. It was released in November 2014 and is powered by the Android 4.4.2 operating system. The device features a 5.0-inch IPS LCD display with a resolution of 720 x 1280 pixels, which translates to 294 pixels per inch. The smartphone is equipped with a 13-megapixel rear camera <b>with autofocus<sup>NC</sup></b> and LED flash, as well as an 8-megapixel front-facing camera. It is powered by a Mediatek MT6582 chipset with a quad-core 1.3 GHz Cortex-A7 CPU and a Mali-400MP2 GPU. The device has 8GB of internal storage and 1GB of RAM, with the option to expand storage via a microSDHC card. The GSmart Mika M3 also features Wi-Fi, Bluetooth, GPS, and <b>NFC connectivity<sup>I</sup></b> , as well as a 3.5mm headphone jack and a microUSB 2.0 port for charging and data transfer. The device is available in Navy Blue and is powered by a 1900mAh Li-Ion battery, which provides up to 150 hours of stand-by time and up to 4 hours and 30 minutes of talk time.	The GSmart Mika M3 is a smartphone developed by Gigabyte. It was <b>released in November 2014<sup>I</sup></b> and is powered by the Android 4.4.2 operating system. The device features a 5.0-inch IPS LCD display with a resolution of 720 x 1280 pixels, which translates to 294 pixels per inch. The smartphone is equipped with a 13-megapixel rear camera with autofocus and LED flash, as well as an 8-megapixel front-facing camera. It is powered by a Mediatek MT6582 chipset with a quad-core 1.3 GHz Cortex-A7 CPU and a Mali-400MP2 GPU. The device has 8GB of internal storage and 1GB of RAM, with the option to expand storage via a microSDHC card. The GSmart Mika M3 also features Wi-Fi, Bluetooth, GPS, and <b>NFC connectivity<sup>I</sup></b> , as well as a 3.5mm headphone jack and a microUSB 2.0 port for charging and data transfer. The device is <b>available in Navy Blue<sup>M</sup></b> and is powered by a 1900mAh Li-Ion battery, which provides up to 150 hours of stand-by time and up to 4 hours and 30 minutes of talk time.
Zephyr	The Gigabyte GSmart Mika M3 is a smartphone that boasts a 5.0-inch IPS LCD display with a resolution of 720 x 1280 pixels. It runs on Android 4.4.2 (KitKat) and is powered by a quad-core 1.3 GHz Cortex-A7 processor with a Mali-400MP2 GPU. The device comes with 8GB of internal storage and 1GB of RAM, which can be expanded via a dedicated microSDHC card slot. The primary camera is a 13-megapixel <b>autofocus lens<sup>NC</sup></b> with an LED flash, while the secondary camera is an 8-megapixel lens. Connectivity options include Wi-Fi, Bluetooth, GPS, FM radio, and microUSB 2.0. The smartphone is equipped with a 1900mAh removable Li-Ion battery that provides up to <b>4 hours<sup>I</sup></b> of talk time and up to 150 hours of standby time. The device is available in Navy Blue and supports dual SIM cards (Mini-SIM/Micro-SIM).	The Gigabyte GSmart Mika M3 is a smartphone that boasts a 5.0-inch IPS LCD display with a resolution of 720 x 1280 pixels. It runs on Android 4.4.2 (KitKat) and is powered by a quad-core 1.3 GHz Cortex-A7 processor with a Mali-400MP2 GPU. The device comes with 8GB of internal storage and 1GB of RAM, which can be expanded via a dedicated microSDHC card slot. The primary camera is a 13-megapixel autofocus lens with an LED flash, while the secondary camera is an 8-megapixel lens. Connectivity options include Wi-Fi, Bluetooth, GPS, FM radio, and microUSB 2.0. The smartphone is equipped with a 1900mAh removable Li-Ion battery that provides <b>up to 4 hours of talk time<sup>I</sup></b> and up to 150 hours of standby time. The device is available in Navy Blue and supports dual SIM cards (Mini-SIM/Micro-SIM).
GPT-3.5	The GSmart Mika M3, also known as the Gigabyte GSmart Mika M3, is a smartphone that was announced in October 2014 and released in November of the same year. It features a 5.0-inch IPS LCD display with a resolution of 720 x 1280 pixels and is powered by a Quad-core 1.3 GHz Cortex-A7 processor. The phone comes with 1GB of RAM and 8GB of internal storage, which can be expanded via a microSDHC card slot. The main camera is 13MP with LED flash, while the selfie camera is 8MP. The phone runs on Android 4.4.2 (KitKat) and is equipped with a 1900mAh removable battery. Other features include dual SIM capability, Wi-Fi, Bluetooth, GPS, FM radio, and microUSB 2.0. The GSmart Mika M3 is available in the color Navy Blue and has a sleek and compact design, measuring 144 x 70.5 x 8.3 mm and weighing 146g.	The GSmart Mika M3, also known as the <b>Gigabyte GSmart Mika M3<sup>I</sup></b> , is a smartphone that was announced in October 2014 and released in November of the same year. It features a 5.0-inch IPS LCD display with a resolution of 720 x 1280 pixels and is powered by a Quad-core 1.3 GHz Cortex-A7 processor. The phone comes with 1GB of RAM and 8GB of internal storage, which can be expanded via a microSDHC card slot. <b>The main camera is 13MP with LED flash<sup>NC</sup>, while the selfie camera is 8MP<sup>NC</sup></b> . The phone runs on Android 4.4.2 (KitKat) <b>and is equipped with a 1900mAh removable battery<sup>NC</sup></b> . Other features include dual SIM capability, Wi-Fi, Bluetooth, GPS, FM radio, and microUSB 2.0. <b>The GSmart Mika M3 is available in the color Navy Blue and has a sleek and compact design<sup>NC</sup></b> , measuring 144 x 70.5 x 8.3 mm and weighing 146g.

Table 9: An example of the model outputs on gsmarena (see the input in Table 9).

## Kristall Saratov – HC Norilsk

### metadata

<b>id</b>	11482971
<b>startDatetime</b>	2023-11-28 14:30:00
<b>startTimestamp</b>	1701181800

### awayScore

<b>display</b>	5
<b>period1</b>	2
<b>period2</b>	3
<b>period3</b>	0
<b>normaltime</b>	5

### season

<b>name</b>	VHL 23/24
<b>year</b>	23/24
<b>id</b>	53511

### awayTeam

<b>name</b>	HC Norilsk
<b>shortName</b>	HC Norilsk
<b>nameCode</b>	NOR
<b>id</b>	485576
<b>country</b>	Russia

### homeTeam

<b>name</b>	Kristall Saratov
<b>shortName</b>	Kristall Saratov
<b>nameCode</b>	KRI
<b>id</b>	78185
<b>country</b>	Russia

### homeScore

<b>display</b>	1
<b>period1</b>	0
<b>period2</b>	0
<b>period3</b>	1
<b>normaltime</b>	1

Figure 9: An example of an input from ice\_hockey (see the outputs in Table 10).



Model	Human annotations ( $\mathcal{E}_{\text{hum}}$ )	GPT-4 annotations ( $\mathcal{E}_{\text{gpt}}$ )
Llama2	In the VHL 23/24 season, Kristall Saratov hosted HC Norilsk at their home arena. Kristall Saratov <b>won the game<sup>1</sup></b> with a final score of 1-5. Kristall Saratov's home team <b>scored one goal in the first period<sup>1</sup>, while HC Norilsk scored two goals in the second and third periods<sup>1</sup>. The game lasted for 60 minutes<sup>NC</sup></b> , with Kristall Saratov's home team <b>leading the game for the entire duration<sup>1</sup></b> .	In the VHL 23/24 season, Kristall Saratov hosted HC Norilsk at their home arena. <b>Kristall Saratov won the game with a final score of 1-5.<sup>1</sup> Kristall Saratov's home team scored one goal in the first period<sup>1</sup>, while HC Norilsk scored two goals in the second and third periods.<sup>1</sup> The game lasted for 60 minutes<sup>NC</sup>, with Kristall Saratov's home team leading the game for the entire duration.<sup>1</sup></b>
Mistral	On November 28th, 2023, the VHL 23/24 season game between Kristall Saratov and HC Norilsk took place. The game was played at the home arena of Kristall Saratov, with the home team <b>winning 1-0<sup>1</sup></b> . The game was played in three periods, with the first period <b>ending in a 0-0 tie<sup>1</sup></b> , the second period ending in <b>a 2-3 score in favor of HC Norilsk<sup>1</sup></b> , and the third period ending in a 1-0 score in favor of Kristall Saratov. The game ended in normal time, with a <b>total score of 3-1 in favor of HC Norilsk<sup>1</sup></b> .	On November 28th, 2023, the VHL 23/24 season game between Kristall Saratov and HC Norilsk took place. The game was played at the home arena of Kristall Saratov, <b>with the home team winning 1-0<sup>1</sup></b> . The game was played in three periods, with the first period ending in a 0-0 tie, <b>the second period ending in a 2-3 score in favor of HC Norilsk<sup>1</sup></b> , and the third period ending in a 1-0 score in favor of Kristall Saratov. The game ended in normal time, <b>with a total score of 3-1 in favor of HC Norilsk<sup>1</sup></b> .
Zephyr	In the VHL 23/24 season, Kristall Saratov hosted HC Norilsk. The first period saw two <b>quick<sup>NC</sup> goals</b> from Norilsk, with <b>one in the second and third periods as well<sup>1</sup></b> , resulting in a dominant <b>5-0<sup>1</sup> victory</b> . Kristall Saratov managed to score a <b>single goal</b> in the third period, but it was not enough to overcome Norilsk's lead. The final score was 1-5 in favor of HC Norilsk.	In the VHL 23/24 season, Kristall Saratov hosted HC Norilsk. The first period saw two quick goals from Norilsk, <b>with one in the second<sup>1</sup> and third periods as well, resulting in a dominant 5-0 victory<sup>1</sup></b> . Kristall Saratov managed to score a <b>single goal</b> in the third period, but it was not enough to overcome Norilsk's lead. The final score was 1-5 in favor of HC Norilsk.
GPT-3.5	In a match from the VHL 23/24 season, Kristall Saratov faced off against HC Norilsk. The game took place on November 28, 2023, at <b>15:30<sup>1</sup></b> . Despite a goal from Kristall Saratov in the third period, HC Norilsk emerged victorious with a final score of 5-1. The away team dominated the game, scoring 2 goals in the first period and 3 goals in the second period. Kristall Saratov struggled to keep up, resulting in a clear win for HC Norilsk.	In a match from the VHL 23/24 season, Kristall Saratov faced off against HC Norilsk. The game took place on November 28, 2023, at 15:30. Despite a goal from Kristall Saratov in the third period, HC Norilsk emerged victorious with a final score of 5-1. The away team dominated the game, scoring 2 goals in the first period and 3 goals in the second period. Kristall Saratov struggled to keep up, resulting in a clear win for HC Norilsk.

Table 10: An example of the model outputs on ice\_hockey (see the input in Table 10).

## Djibouti

Life expectancy at birth. The average number of years that a newborn could expect to live, if he or she were to pass through life exposed to the sex- and age-specific death rates prevailing at the time of his or her birth, for a specific year, in a given country, territory, or geographic area.

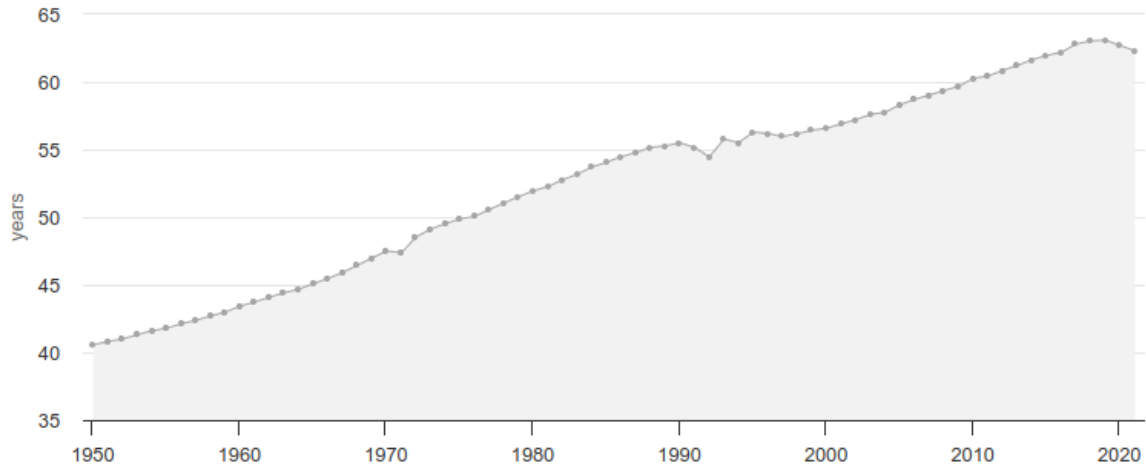


Figure 10: An example of an input from owid (see the outputs in Table 11).

Model	Human annotations ( $\mathcal{E}_{\text{hum}}$ )	GPT-4 annotations ( $\mathcal{E}_{\text{gpt}}$ )
<b>Llama2</b>	Life expectancy at birth in Djibouti has increased significantly over the past few decades, with a notable increase from 40.5665 years in 1950 to 62.3049 years in 2021.	Life expectancy at birth in Djibouti has increased significantly over the past few decades, with a notable increase from 40.5665 years in 1950 to 62.3049 years in 2021.
<b>Mistral</b>	This chart displays the life expectancy at birth in Djibouti from 1950 to 2021. The data shows a steady increase in life expectancy over the years, with a <b>significant jump in the 1970s and 1980s<sup>I</sup></b> . In 2021, the average life expectancy at birth in Djibouti was 62.3 years.	This chart displays the life expectancy at birth in Djibouti from 1950 to 2021. The data shows a <b>steady increase<sup>I</sup></b> in life expectancy over the years, with a <b>significant jump in the 1970s and 1980s<sup>I</sup></b> . In 2021, the average life expectancy at birth in Djibouti was <b>62.3 years<sup>I</sup></b> .
<b>Zephyr</b>	This line graph illustrates the evolution of life expectancy at birth in Djibouti from 1950 to 2021. Initially, the average number of years that a newborn could expect to live was around 40.5 years in 1950, and it gradually increased over the following decades, reaching 62.3 years in 2021. The trend shows a steady improvement, with some fluctuations in the <b>1980s<sup>I</sup></b> and 1990s. The data presented here is based on the sex- and age-specific death rates prevailing at the time of birth for a specific year in a given country, territory, or geographic area.	This line graph illustrates the evolution of life expectancy at birth in Djibouti from 1950 to 2021. Initially, the average number of years that a newborn could expect to live was <b>around 40.5 years in 1950<sup>I</sup></b> , and it gradually increased over the following decades, reaching <b>62.3 years in 2021<sup>I</sup></b> . <b>The trend shows a steady improvement<sup>M</sup></b> , with some fluctuations in the 1980s and 1990s. The data presented here is based on the sex- and age-specific death rates prevailing at the time of birth for a specific year in a given country, territory, or geographic area.
<b>GPT-3.5</b>	The chart shows the life expectancy at birth in Djibouti from 1950 to 2021. The life expectancy has steadily increased over the years, from 40.57 years in 1950 to 62.30 years in 2021, with some fluctuations in between. The data reflects the average number of years a newborn could expect to live if exposed to the prevailing death rates at the time of birth.	The chart shows the life expectancy at birth in Djibouti from 1950 to 2021. <b>The life expectancy has steadily increased over the years<sup>I</sup></b> , from <b>40.57 years in 1950<sup>I</sup></b> to <b>62.30 years in 2021<sup>I</sup></b> , with some fluctuations in between. The data reflects the average number of years a newborn could expect to live if exposed to the prevailing death rates at the time of birth.

Table 11: An example of the model outputs on owid (see the input in Table 11).

## Driven Brands

<b>country</b>	United States of America
<b>industry</b>	automotive services industry
<b>instance of</b>	business

Figure 11: An example of an input from wikidata (see the outputs in Table 12).

Model	Human annotations ( $\mathcal{E}_{\text{hum}}$ )	GPT-4 annotations ( $\mathcal{E}_{\text{gpt}}$ )
<b>Llama2</b>	Driven Brands is an American <u>franchise holding<sup>NC</sup></u> company that <u>owns, operates, and franchises<sup>NC</sup></u> automotive services businesses in the United States. <u>The company was founded in 1996 and is headquartered in Charlotte, North Carolina. It operates several brands, including Maaco, Meineke, Econo Lube, and CARSTAR<sup>NC</sup>.</u>	Driven Brands is an American franchise holding company that owns, operates, and franchises automotive services businesses in the United States. The company was <u>founded in 1996<sup>NC</sup></u> and is <u>headquartered in Charlotte, North Carolina<sup>NC</sup></u> . It operates several brands, including Maaco, Meineke, Econo Lube, and CARSTAR.
<b>Mistral</b>	Driven Brands is a business operating in the automotive services industry in the United States of America.	Driven Brands is a business operating in the automotive services industry in the United States of America.
<b>Zephyr</b>	Driven Brands is a business operating in the automotive services industry based in the United States of America. As an instance of the broader category of businesses, Driven Brands <u>provides goods or services to customers in exchange for financial remuneration<sup>NC</sup></u> .	Driven Brands is a business operating in the automotive services industry based in the United States of America. <u>As an instance of the broader category of businesses,<sup>O</sup></u> Driven Brands provides <u>goods or services<sup>I</sup></u> to customers in exchange for financial remuneration.
<b>GPT-3.5</b>	Driven Brands is a business based in the United States of America, operating within the automotive services industry.	Driven Brands is a business based in the United States of America, operating within the automotive services industry.

Table 12: An example of the model outputs on wikidata (see the input in Table 12).

		Incorrect		Not Checkable		Misleading		Other		All categories		# Tok.
		$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}}$	$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}}$	$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}}$	$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}}$	$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}}$	
<b>Llama 2</b>	openweather	4.27	6.99	2.23	1.93	0.58	0.20	0.16	0.26	7.24	9.38	166.7
	gsmarena	0.67	1.96	0.87	0.31	0.20	0.24	0.09	0.10	1.83	2.61	90.3
	ice_hockey	1.47	1.98	0.99	0.90	0.09	0.00	0.12	0.05	2.67	2.93	63.3
	owid	1.30	2.26	0.58	0.15	0.32	0.13	0.07	0.00	2.27	2.54	60.3
	wikidata	0.15	0.74	1.59	1.24	0.08	0.03	0.07	0.04	1.89	2.05	38.5
<b>Mistral</b>	openweather	5.65	8.76	2.47	0.41	1.36	0.50	0.53	0.04	10.01	9.71	193.5
	gsmarena	0.51	2.02	1.10	0.94	0.30	0.25	0.14	0.02	2.05	3.23	146.3
	ice_hockey	1.47	1.79	0.92	0.76	0.14	0.08	0.13	0.10	2.66	2.73	92.4
	owid	2.42	3.11	0.40	0.18	0.29	0.27	0.07	0.02	3.18	3.58	91.1
	wikidata	0.10	0.45	0.72	0.40	0.13	0.19	0.38	0.31	1.33	1.35	51.0
<b>Zephyr</b>	openweather	4.22	7.39	1.01	0.37	0.34	0.41	0.16	0.00	5.73	8.17	130.9
	gsmarena	0.33	2.24	1.04	1.00	0.23	0.48	0.11	0.01	1.71	3.73	142.8
	ice_hockey	0.89	1.57	0.61	0.20	0.11	0.04	0.07	0.01	1.68	1.82	83.1
	owid	1.68	2.64	0.49	0.16	0.22	0.30	0.10	0.01	2.49	3.11	85.2
	wikidata	0.09	0.36	0.72	0.29	0.12	0.21	0.36	0.21	1.29	1.07	48.1
<b>GPT-3.5</b>	openweather	1.57	4.34	0.57	0.15	0.38	0.53	0.05	0.00	2.57	5.02	112.8
	gsmarena	0.20	1.64	0.80	1.42	0.21	0.34	0.17	0.01	1.38	3.41	129.5
	ice_hockey	0.81	0.76	0.46	0.17	0.07	0.04	0.09	0.01	1.43	0.98	84.4
	owid	0.64	1.87	0.25	0.02	0.17	0.27	0.01	0.00	1.07	2.16	62.2
	wikidata	0.05	0.18	0.36	0.16	0.06	0.10	0.05	0.10	0.52	0.54	35.7

Table 13: The average *numbers of errors per output* for each domain (lower is better). We also include the average number of tokens per output in the rightmost column. See Table 3 for aggregated results.

		Incorrect		Not Checkable		Misleading		Other		All categories	
		$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}}$	$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}}$	$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}}$	$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}}$	$\mathcal{E}_{\text{hum}}$	$\mathcal{E}_{\text{gpt}}$
<b>Llama 2</b>	openweather	76%	98%	71%	68%	30%	14%	12%	22%	90%	100%
	gsmarena	38%	86%	49%	24%	18%	16%	5%	7%	74%	94%
	ice_hockey	70%	77%	51%	55%	7%	0%	10%	5%	89%	94%
	owid	70%	92%	37%	14%	24%	12%	6%	0%	88%	93%
	wikidata	12%	47%	79%	63%	8%	2%	5%	4%	87%	89%
<b>Mistral</b>	openweather	78%	99%	61%	16%	38%	18%	17%	4%	92%	100%
	gsmarena	33%	89%	54%	54%	19%	17%	12%	2%	73%	99%
	ice_hockey	74%	83%	62%	52%	12%	8%	9%	7%	88%	97%
	owid	75%	96%	27%	8%	24%	25%	6%	2%	86%	97%
	wikidata	8%	34%	44%	29%	10%	17%	24%	27%	67%	72%
<b>Zephyr</b>	openweather	82%	99%	45%	25%	23%	25%	14%	0%	90%	100%
	gsmarena	22%	91%	47%	52%	19%	30%	9%	1%	61%	99%
	ice_hockey	54%	78%	39%	17%	11%	4%	6%	1%	76%	87%
	owid	70%	93%	31%	9%	20%	26%	7%	1%	85%	95%
	wikidata	6%	29%	49%	22%	8%	18%	22%	18%	66%	66%
<b>GPT-3.5</b>	openweather	64%	98%	32%	12%	23%	25%	5%	0%	75%	99%
	gsmarena	15%	74%	42%	58%	13%	19%	12%	1%	57%	97%
	ice_hockey	64%	50%	30%	12%	7%	4%	8%	1%	76%	57%
	owid	43%	89%	13%	2%	14%	24%	1%	0%	57%	90%
	wikidata	4%	14%	27%	14%	6%	9%	5%	9%	38%	36%

Table 14: The percentage of *outputs containing at least one error* for each domain (lower is better). See Table 4 for aggregated results.