# JECC: Commonsense Reasoning Tasks Derived from Interactive Fictions

**Mo Yu**[*1]   **Yi Gu**[*2]   **Xiaoxiao Guo**[3]   **Yufei Feng**[4]
**Xiaodan Zhu**[4]   **Michael Greenspan**[4]   **Murray Campbell**[5]   **Chuang Gan**[5]
[1] WeChat AI   [2] UC San Diego   [3] LinkedIn   [4] Queens University   [5] IBM Research
moyumyu@tencent.com   yig025@ucsd.edu

## Abstract

Commonsense reasoning simulates the human ability to make presumptions about our physical world, and it is an essential cornerstone in building general AI systems. We propose a new commonsense reasoning dataset based on human's Interactive Fiction (IF) gameplay walkthroughs as human players demonstrate plentiful and diverse commonsense reasoning. The new dataset provides a natural mixture of various reasoning types and requires multi-hop reasoning. Moreover, the IF game-based construction procedure requires much less human interventions than previous ones. Different from existing benchmarks, our dataset focuses on the assessment of functional commonsense knowledge rules rather than factual knowledge. Hence, in order to achieve higher performance on our tasks, models need to effectively utilize such functional knowledge to infer the outcomes of actions, rather than relying solely on memorizing facts. Experiments show that the introduced dataset is challenging to previous machine reading models as well as the new large language models with a significant 20% performance gap compared to human experts.[1]

## 1 Introduction

There has been a flurry of datasets and benchmarks proposed to address natural language-based commonsense reasoning (Levesque et al., 2012; Zhou et al., 2019; Talmor et al., 2019; Mullenbach et al., 2019; Jiang et al., 2020; Sap et al., 2019a; Bhagavatula et al., 2019; Huang et al., 2019; Bisk et al., 2020; Sap et al., 2019b; Zellers et al., 2018). These benchmarks usually adopt a multi-choice form – with the input query and an optional short paragraph of the background description, each candidate forms a statement; the task is to predict the

---

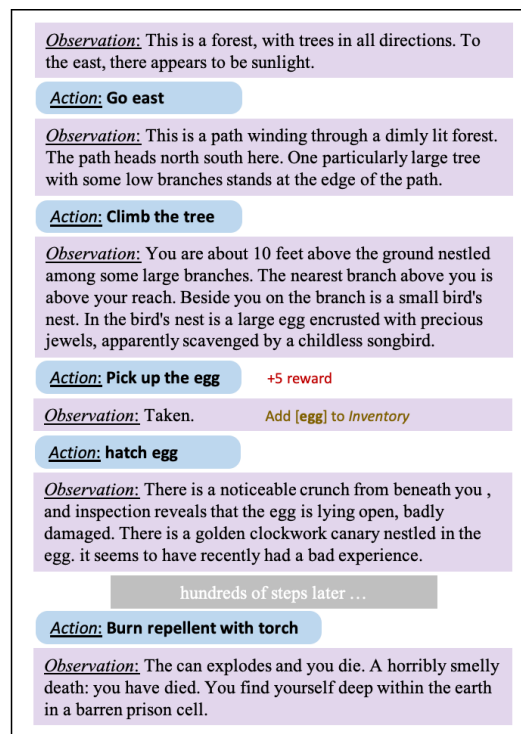[1]Our code and data are released at https://github.com/Gorov/zucc.



Figure 1: Classic dungeon game *Zork1* gameplay sample. The player receives textual observations describing the current game state and sends textual action commands to control the protagonist.

statement that is consistent with some commonsense knowledge facts.

These benchmarks share some limitations, as they are mostly constructed to focus on a single reasoning type and require similar validation-based reasoning. First, most benchmarks concentrate on a specific facet and ask human annotators to write candidate statements related to the particular type of commonsense. As a result, the distribution of these datasets is unnatural and biased to a specific facet. For example, most benchmarks focus on collocation, association, or other relations (e.g., ConceptNet (Speer et al., 2017) relations) between words or concepts (Levesque et al., 2012; Talmor et al., 2019; Mullenbach et al., 2019; Jiang et al., 2020). Other examples include temporal commonsense (Zhou et al., 2019), physical interactions

between actions and objects (Bisk et al., 2020), emotions and behaviors of people under the given situation (Sap et al., 2019b), and cause-effects between events and states (Sap et al., 2019a; Bhagavatula et al., 2019; Huang et al., 2019). Second, most datasets require validation-based reasoning between a commonsense fact and a text statement but neglect hops over multiple facts.[2] The previous work's limitations bias the model evaluation. For example, pre-trained Language Models (PLMs), such as BERT (Devlin et al., 2019), can well handle most benchmarks, because their pre-training process may include texts on the required facts thus provide shortcuts to a dominating portion of commonsense validation instances. In summary, the above limitations of previous benchmarks lead to discrepancies among practical NLP tasks that require broad reasoning ability on various facets.

**Our Contribution.** We derive *a new commonsense reasoning dataset from the model-based reinforcement learning challenge* of Interactive Fictions (IF) to address the above limitations. Recent advances (Hausknecht et al., 2019; Ammanabrolu and Hausknecht, 2020; Guo et al., 2020) in IF games have recognized several commonsense reasoning challenges, such as detecting valid actions and predicting different actions' effects. Figure 1 illustrates sample gameplay of the classic game *Zork1* and the required commonsense knowledge. We derive a commonsense dataset from human players' gameplay records related to the second challenge, i.e., predicting which textual observation is most likely after applying an action or a sequence of actions to a given game state.

The derived dataset naturally addresses the aforementioned limitations in previous datasets. First, predicting the next observation naturally requires various commonsense knowledge and reasoning types. As shown in Figure 1, a primary commonsense type is spatial reasoning, e.g., "climb the tree" makes the protagonist up on a tree. Another primary type is reasoning about object interactions. For example, keys can open locks (object relationships); "hatch egg" will reveal "things" inside the egg (object properties); "burn repellent" leads to an explosion and kills the player (physical reasoning). The above interactions are more com-

prehensive than the relationships defined in ConceptNet as used in previous datasets. Second, the rich textual observation enables more complex reasoning over direct commonsense validation. Due to the textual observation's narrative nature, a large portion of the textual observations are not a sole statement of the action effect, but an extended narrates about what happens because of the effect.[3] Third, our commonsense reasoning task formulation shares the essence of dynamics model learning for model-based RL solutions related to world models and MuZero (Ha and Schmidhuber, 2018; Schrittwieser et al., 2019). Therefore, models developed on our benchmarks provide direct values to model-based RL for text-game playing.

Finally, compared to previous works that heavily rely on human annotation, our dataset construction requires minimal human effort, providing great **expansibility**. For example, with large amounts of available IF games in dungeon crawls, Sci-Fi, mystery, comedy, and horror, it is straightforward to extend our dataset to include more data samples and cover a wide range of genres. We can also naturally increase the reasoning difficulty by increasing the prediction horizon of future observations after taking multi-step actions instead of a single one.

In summary, we introduce a new commonsense reasoning dataset construction paradigm, collectively with two datasets. The larger dataset covers 29 games in multiple domains from the *Jericho Environment* (Hausknecht et al., 2019), named the Jericho Environment Commonsense Comprehension task (**JECC**). The smaller dataset, aimed for the single-domain test and fast model development, includes four IF games in the *Zork Universe*, named Zork Universe Commonsense Comprehension (**ZUCC**). We provide strong baselines to the datasets and categorize their performance gap compared to human experts.

## 2 Related Work

Previous work has identified various types of commonsense knowledge humans master for text understanding. As discussed in the introduction section, most existing datasets cover one or a few limited types. Also, they mostly have the form of commonsense fact validation based on a text statement.

**Semantic Relations between Concepts.** Most

---

[2] Some datasets include a portion of instances that require explicit reasoning capacity, such as (Bhagavatula et al., 2019; Huang et al., 2019; Bisk et al., 2020; Sap et al., 2019b). But still, standalone facts can solve most such instances.

[3] For some actions, such as get and drop objects, the next observations are simple statements. We removed some of these actions. Details can be found in Section 3.

previous datasets cover the semantic relations between words or concepts. These relations include the concept hierarchies, such as those covered by WordNet or ConceptNet, and word collocations and associations. For example, the early work Winograd (Levesque et al., 2012) evaluates the model's ability to capture word collocations, associations between objects, and their attributes as a pronoun resolution task. The work by (Talmor et al., 2019) is one of the first datasets covering the ConceptNet relational tuple validation as a question-answering task. The problem asks the relation of a source object, and the model selects the target object that satisfies the relation from four candidates. (Mullenbach et al., 2019) focus on the collocations between adjectives and objects. Their task takes the form of textual inference, where a premise describes an object and the corresponding hypothesis consists of the object that is modified by an adjective. (Jiang et al., 2020) study associations among multiple words, i.e., whether a word can be associated with two or more given others (but the work does not formally define the types of associations). They propose a new task format in games where the player produces as many words as possible by combining existing words.

**Causes/Effects between Events or States.** Previous work proposes datasets that require causal knowledge between events and states (Sap et al., 2019a; Bhagavatula et al., 2019; Huang et al., 2019). (Sap et al., 2019a) takes a text generation or inference form between a cause and an effect. (Bhagavatula et al., 2019) takes a similar form to ours – a sequence of two observations is given, and the model selects the plausible hypothesis from multiple candidates. Their idea of data construction can also be applied to include any types of knowledge. However, their dataset only focuses on causal relations between events. The work of (Huang et al., 2019) utilizes multi-choice QA on a background paragraph, which covers a wider range of casual knowledge for both events and statements.

**Other Commonsense Datasets.** (Zhou et al., 2019) proposed a unique temporal commonsense dataset. The task is to predict a follow-up event's duration or frequency, given a short paragraph describing an event. (Bisk et al., 2020) focus on physical interactions between actions and objects, namely whether an action over an object leads to a target effect in the physical world. These datasets can be solved by mostly applying the correct com-

monsense facts; thus, they do not require reasoning. (Sap et al., 2019b) propose a task of inferring people's emotions and behaviors under the given situation. Compared to the others, this task contains a larger portion of instances that require reasoning beyond fact validation. The above tasks take the multi-choice question-answering form.

**Next-Sentence Prediction.** The next sentence prediction tasks, such as SWAG (Zellers et al., 2018), are also related to our work. These tasks naturally cover various types of commonsense knowledge and sometimes require reasoning. The issue is that the way they guarantee distractor candidates to be irrelevant greatly simplified the task. In comparison, our task utilizes the IF game engine to ensure actions uniquely determining the candidates, and ours has human-written texts.

Finally, our idea is closely related to (Yao et al., 2020), which creates a task of predicting valid actions for each IF game state. (Yao et al., 2020, 2021) also discussed the advantages of the supervised tasks derived from IF games for natural language understanding purpose.

## 3 Dataset Construction: Commonsense Challenges from IF Games

We pick games supported by the *Jericho* environment (Hausknecht et al., 2019) to construct the **JECC** dataset.[4] We pick games in the *Zork Universe* for the **ZUCC** dataset.[5] We first introduce the necessary definitions in the IF game domain and then describe how we construct our **ZUCC** and **JECC** datasets as the forward prediction tasks based on human players' gameplay records, followed by a summary on the improved properties of our dataset compared to previous ones. The dataset will be released for public usage. It can be created with our released code with MIT License.

### 3.1 Interactive Fiction Game Background

Each IF game can be defined as a Partially Observable Markov Decision Process (POMDP), namely a 7-tuple of $\langle S, A, T, O, \Omega, R, \gamma \rangle$, representing the hidden game state set, the action set, the state transition function, the set of textual observations com-

---

[4]We collect the games *905, acorncourt, advent, adventureland, afflicted, awaken, balances, deephome, dragon, enchanter, inhumane, library, moonlit, omniquest, pentari, reverb, snacktime, sorcerer, zork1* for training, *zork3, detective, ztuu, jewel, zork2* as the development set, *temple, gold, karn, zenon, wishbringer* as the test set.

[5]We pick *Zork1, Enchanter*, and *Sorcerer* as the training set, and the dev and sets are non-overlapping split from *Zork3*.
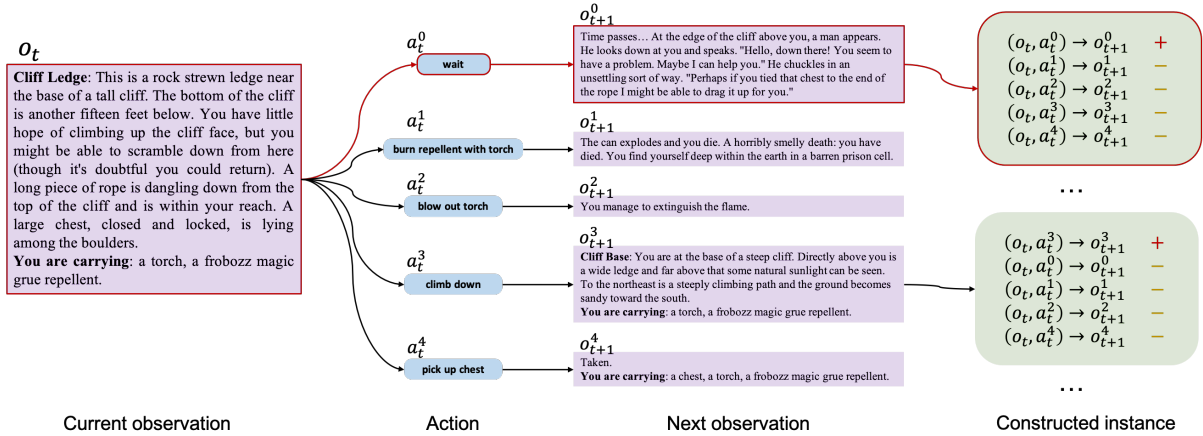
Figure 2: Illustration of our data construction process, taking an example from *Zork3*. $+/-$: positive/negative labels. The red colored path denotes the tuple and the resulted data instance from the human walkthrough.

posed from vocabulary words, the textual observation function, the reward function, and the discount factor respectively. The game playing agent interacts with the game engine in multiple turns until the game is over or the maximum number of steps is reached. At the $t$-th turn, the agent receives a textual observation describing the current game state $o_t \in O$ and sends a textual action command $a_t \in A$ back. The agent receives additional reward scalar $r_t$ which encodes the game designers' objective of game progress. Thus the task of the game playing can be formulated to generate a textual action command per step as to maximize the expected cumulative discounted rewards $\mathbf{E}\left[\sum_{t=0}^{\infty} \gamma^t r_t\right]$. Most IF games have a deterministic dynamics, and the next textual observation is uniquely determined by an action choice. Unlike most previous work on IF games that design autonomous learning agents, we utilize human players' gameplay records that achieve the highest possible game scores.

**Trajectories and Walkthroughs.** A *trajectory* in text game playing is a sequence of tuples $\{(o_t, a_t, r_t, o_{t+1})\}_{t=0}^{T-1}$, starting with the initial textual observation $o_0$ and the game terminates at time step $t = T$, i.e., the last textual observation $o_T$ describes the game termination scenario. We define the *walkthrough* of a text game as a trajectory that completes the game progress and achieves the highest possible game scores.

### 3.2 Data Construction from the Forward Prediction Task

**The Forward Prediction Task.** We represent our commonsense reasoning benchmark as a next-

|  | #WT Tuples | #Tuples before Proc | #Tuples after Proc |
|---|---|---|---|
| **ZUCC** | | | |
| Train | 913 | 17,741 | 10,498 |
| All Eval | 271 | 4,069 | 2,098 |
| Dev | – | – | 1,276 |
| Test | – | – | 822 |
| **JECC** | | | |
| Train | 2,526 | 48,843 | 24,801 |
| All Eval | 2,063 | 53,160 | 25,891 |
| Dev | 917 | – | – |
| Test | 1,146 | – | – |

Table 1: Data statistics of our **ZUCC** and **JECC** tasks. **WT** stands for walkthrough. The evaluation sets of **ZUCC** consist of all tuples after post-processing. The evaluation sets of **JECC** only consist of tuples in walkthroughs. As discussed in the dataset construction criteria (Section 3.3), we only evaluate the models with tuples from the walkthroughs to ensure a representative distribution of required knowledge.

observation prediction task, given the current observation and action. The benchmark construction starts with all the tuples in a walkthrough trajectory, and we then extend the tuple set by including all valid actions and their corresponding next-observations conditioned on the current observations in the walkthrough. Specifically, for a walkthrough tuple $(o_t, a_t, r_t, o_{t+1})$, we first obtain the complete valid action set $A_t$ for $o_t$. We sample and collect one next observation $o_{t+1}^j$ after executing the corresponding action $a_t^j \in A_t$. The next-observation prediction task is thus to select the next observation $o_{t+1}^j$ given $(o_t, a_t^j)$ from the complete set of next observations $O_{t+1} = \{o_{t+1}^k, \forall k\}$. Figure 2 illustrates our data construction process.

**Data Processing.** We collect tuples from the walkthrough data provided by the Jericho environments. We detect the valid actions via the Jericho API and the game-specific templates. Following previous work (Hausknecht et al., 2019), we augmented the observation with the textual feedback returned by the command [*inventory*] and [*look*]. The former returns the protagonist's objects, and the latter returns the current location description. When multiple actions lead to the same next-observation, we randomly keep one action and next-observation in our dataset. We remove the `drop OBJ` actions since it only leads to synthetic observations with minimal variety. For each step $t$, we keep at most 15 candidate observations in $O_t$ for the evaluation sets. When there are more than 15 candidates, we select the candidate that differs most from $o_t$ with Rouge-L measure (Lin, 2004).

During evaluation, for **JECC**, we only evaluate on the tuples on walkthroughs. As will be discussed in 3.3, this helps our evaluation reflects a natural distribution of commonsense knowledge required, which is an important criterion pointed out by our introduction. However for **ZUCC** the walkthough data is too small, therefore we consider all the tuples during evaluation. This leads to some problems. First, there are actions that do not have the form of `drop OBJ` but have the actual effects of dropping objects. Through the game playing process, more objects will be collected in the inventory at the later stages. These cases become much easier as long as these non-standard drop actions have been recognized. A similar problem happens to actions like `burn repellent` that can be performed at every step once the object is in the inventory. To deal with such problems, we down-sample these biased actions to achieve similar distributions in development and test sets. Table 1 summarizes statistics of the resulted **JECC** and **ZUCC** datasets.

### 3.3 Design Criterion and Dataset Properties

**Knowledge coverage and distribution.** As discussed in the introduction, an ideal commonsense reasoning dataset needs to cover various commonsense knowledge types, especially useful ones for understanding language. A closely related criterion is that the required commonsense knowledge and reasoning types should reflect a natural distribution in real-world human language activities.

Our **JECC** and **ZUCC** datasets naturally meet

| Dimension | Count | Dimension | Count |
|---|---|---|---|
| similarity | 3 | utility | 6 |
| distinctness | 1 | desire/goal | 1 |
| taxonomic | 0 | quality | 15 |
| part-whole | 5 | comparative | 1 |
| spatial | 16 | temporal | 56 |
| creation | 0 | relational-other | 6 |

Table 2: The covered commonsense knowledge dimensions by our dataset. All the examples require *temporal* knowledge because the knowledge cause-effect is categorized into this type in the schema.

these two criteria. The various IF games cover diverse domains, and human players demonstrate plentiful and diverse commonsense reasoning in finishing the games. The commonsense background information and interventions are recorded in human-written texts (by the game designers and the players, respectively). With the improved coverage of commonsense knowledge following a natural distribution, our datasets have the potential of better evaluating reasoning models, alleviating the biases from previous datasets on a specific knowledge reasoning type.

**Reasoning beyond verification.** We hope to evaluate the models' capabilities in (multi-hop) reasoning with commonsense facts and background texts, beyond simple validation of knowledge facts.

By design, our datasets depart from simple commonsense validation. Neither the input (current observation and action) nor the output (next observation) directly describes a knowledge fact. Instead, they are narratives that form a whole story. Moreover, our task formulation explicitly requires using commonsense knowledge to understand how the action impacts the current state, then reason the effects, and finally verifies whether the next observation coheres with the action effects. These solution steps form a multi-step reasoning process.

### 3.4 The Coverage of Knowledge Dimensions

We conducted human annotation to investigate the range of commonsense knowledge types covered by our datasets. We employed the knowledge type schema from (Ilievski et al., 2021) and manually examined and categorized a total of 56 examples that could be resolved using various types of commonsense knowledge. Despite the small sample size, Table 2 illustrates that our task encompasses a wide array of commonsense types.

Importantly, unlike (Ilievski et al., 2021) and

many other datasets designed for commonsense assessments, our datasets focus on evaluating functional commonsense knowledge, such as rules, rather than factual knowledge. For example, both our datasets and previous work may cover the *spatial* knowledge. However, instead of assessing the static fact "the Great Door is to the south of the Royal Hall", we require an understanding of the functional knowledge that "moving to south make the original position to the north of the new position".

Similarly, instead of simply knowing the *property* fact that "magic grue repellent is explosible", we require the knowledge of the functional rule that "gas in a can may explode when heated". Thus, in conjunction with the knowledge rule that "burning a thing with a torch can heats it", we can infer that the can explodes, resulting in the player's death. Both the *property* and the *causal* (categorized under the *temporal* type) knowledge in this example, required by our task, are functional knowledge rules rather than static facts.

Among all the dimensions, we do not cover the *creation* dimension, as it typically pertains to entity-specific facts rather than general rules. Additionally, the *taxonomic* dimension was not observed in the samples we studied from *Zork3*.

## 4 Neural Inference Baselines

We formulate our task as a textual entailment task that the models infer the next state $o_{t+1}$ given $o_t$ and $a_t$. We provide strong textual entailment-based baselines for our benchmark. We categorize the baselines into two types, namely pairwise textual inference methods and the triplewise inference methods. The notations $o_t$, $a_t$ of observations and actions represent their word sequences.

### 4.1 Neural Inference over Textual Pairs

• **Match LSTM** (Wang and Jiang, 2016) represents a commonly used natural language inference model. Specifically, we concatenate $o_t$ and $a_t$ separated by a special split token as the premise and use the $o_{t+1}^j, j = 1, ...N$ as the hypothesis. For simplicity *we denote $o_t$, $a_t$ and a candidate $o_{t+1}^j$ as $o, a, \tilde{o}$.* We encode the premise and the hypothesis with bidirectional-LSTM model:

$$\boldsymbol{H}^{o,a} = \text{BiLSTM}([o,a]), \boldsymbol{H}^{\tilde{o}} = \text{BiLSTM}(\tilde{o}), \quad (1)$$

where $\boldsymbol{H}^{o,a}$ and $\boldsymbol{H}^{\tilde{o}}$ are the sequences of BiLSTM output $d$-dimensional hidden vectors that correspond to the premise and hypothesis respectively. We apply the bi-attention model to compute the match between the premise and the hypothesis, which is followed by a Bi-LSTM model to get the final hidden sequence for prediction:

$$\bar{\boldsymbol{H}}^{\tilde{o}} = \boldsymbol{H}^{\tilde{o}}\boldsymbol{G}^{\tilde{o}}, \boldsymbol{G}^{\tilde{o}} = \text{SoftMax}((W^g\boldsymbol{H}^{\tilde{o}} + b^g \otimes e)^T\boldsymbol{H}^{o,a})$$
$$\boldsymbol{M} = \text{BiLSTM}([\boldsymbol{H}^{o,a}, \bar{\boldsymbol{H}}^{\tilde{o}}, \boldsymbol{H}^{o,a} - \bar{\boldsymbol{H}}^{\tilde{o}}, \boldsymbol{H}^{o,a} \odot \bar{\boldsymbol{H}}^{\tilde{o}}]).$$

Here $W^g \in \mathbb{R}^{d \times d}$ and $b^g \in \mathbb{R}^d$ are learnable parameters and $e \in \mathbb{R}^{|\tilde{o}|}$ denotes a vector of all 1s. We use a scoring function $f(\cdot)$ to compute matching scores of the premise and the hypothesis via a linear transformation on the max-pooled output of $\boldsymbol{M}$. The matching scores for all $\tilde{o}$ are then fed to a softmax layer for the final prediction. We use the cross-entropy loss as the training objective.

• **BERT Siamese** uses a pre-trained BERT model to separately encode the current observation-action pair $(o_t, a_t)$ and candidate observations $\tilde{o}$. All inputs to BERT start with the "[CLS]" token, and we concatenate $o_t$ and $a_t$ with a "[SEP]" token:

$$\boldsymbol{h}^{o,a} = \text{BERT}([o,a]), \quad \boldsymbol{h}^{\tilde{o}} = \text{BERT}(\tilde{o}),$$
$$l_j = f([\boldsymbol{h}^{o,a}, \boldsymbol{h}^{\tilde{o}}, \boldsymbol{h}^{o,a} - \boldsymbol{h}^{\tilde{o}}, \boldsymbol{h}^{o,a} \odot \boldsymbol{h}^{\tilde{o}}]),$$

where $[\cdot, \cdot]$ denotes concatenation. $\boldsymbol{h}^{o,a}$ and $\boldsymbol{h}^{\tilde{o}}$ are the last layer hidden state vectors of the "[CLS]" token. Similarly, the scoring function $f$ computes matching scores for candidate next-observations by linearly projecting the concatenated vector into a scalar. The matching scores of all $\tilde{o}$ are grouped to a softmax layer for the final prediction.

• **BERT Concat** represents the standard pairwise prediction mode of BERT. We concatenate $o$ and $a$ with a special split token as the first segment and treat $\tilde{o}$ as the second. We then concatenate the two with the "[SEP]" token:

$$l_j = f(\text{BERT}([o, a, \tilde{o}])).$$

The scoring function $f$ linearly projects the last-layer hidden state of the "[CLS]" token into a scalar, and the scores are grouped to a softmax layer for final prediction. This model is much less efficient than the former two as it requires explicit combination of observation-action-next-observation as inputs. Thus this model is impractical due to the huge combinatorial space. Here we report its results for reference.
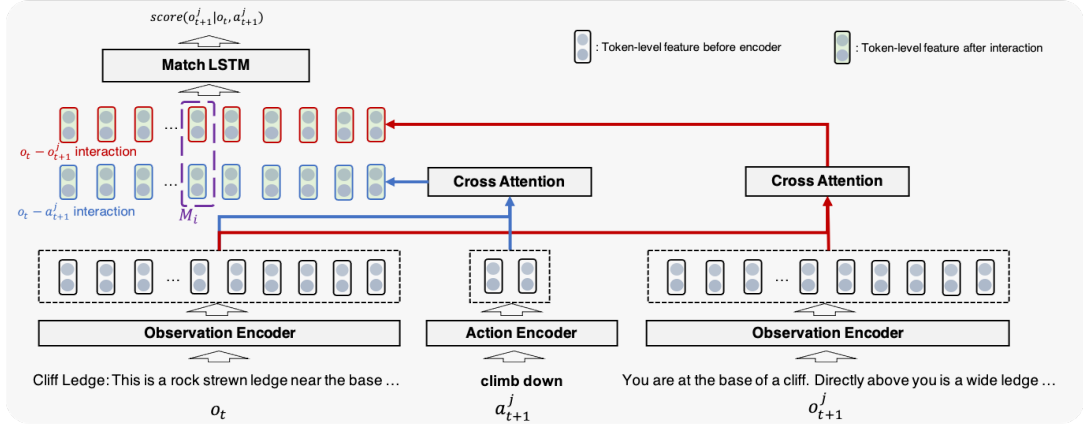
Figure 3: The co-matching architecture for our tasks.

## 4.2 Neural Inference over Textual Triples

Existing work (Lai et al., 2017; Sun et al., 2019; Wang et al., 2019) has applied textual matching and entailment among triples. For example, when applying to multi-choice QA, the entailment among triples is to predict whether a question $q$, an answer option $a$ can be supported by a paragraph $p$. In this section, we apply the most commonly used co-matching approaches (Wang et al., 2018) and its BERT variant to our task. Figure 3 illustrates our co-matching architecture.

• **Co-Matching LSTM** (Wang et al., 2018) jointly encodes the question and answer with the context passage. We extend the idea to conduct the multi-hop reasoning in our setup. Specifically, similar to Equation 1, we first encode the current state observation $o$, the action $a$ and the candidate next state observation $\tilde{o}$ separately with a BiLSTM model, and use $\boldsymbol{H}^o, \boldsymbol{H}^a, \boldsymbol{H}^{\tilde{o}}$ to denote the output hidden vectors respectively.

We then integrate the co-matching to the baseline readers by applying bi-attention described in Equation 2 on $(\boldsymbol{H}^o, \boldsymbol{H}^{\tilde{o}})$, and $(\boldsymbol{H}^a, \boldsymbol{H}^{\tilde{o}})$ using the same set of parameters:

$$\bar{\boldsymbol{H}}^o = \boldsymbol{H}^o \boldsymbol{G}^o, \boldsymbol{G}^o = \text{SoftMax}((W^g \boldsymbol{H}^o + b^g \otimes e_o)^T \boldsymbol{H}^{\tilde{o}})$$
$$\bar{\boldsymbol{H}}^a = \boldsymbol{H}^a \boldsymbol{G}^a, \boldsymbol{G}^a = \text{SoftMax}((W^g \boldsymbol{H}^a + b^g \otimes e_a)^T \boldsymbol{H}^{\tilde{o}}),$$

where $W^g \in \mathbb{R}^{d \times d}$ and $b^g \in \mathbb{R}^d$ are learnable parameters and $e_o \in \mathbb{R}^{|o|}, e_a \in \mathbb{R}^{|a|}$ denote vectors of all 1s. We further concatenate the two output hidden sequences $\bar{\boldsymbol{H}}^o$ and $\bar{\boldsymbol{H}}^a$, followed by a BiLSTM model to get the final sequence representation:

$$M = \text{BiLSTM} \left( \begin{bmatrix} \boldsymbol{H}^{\tilde{o}}, \bar{\boldsymbol{H}}^o, \boldsymbol{H}^{\tilde{o}} - \bar{\boldsymbol{H}}^o, \boldsymbol{H}^{\tilde{o}} \odot \bar{\boldsymbol{H}}^o \\ \boldsymbol{H}^{\tilde{o}}, \bar{\boldsymbol{H}}^a, \boldsymbol{H}^{\tilde{o}} - \bar{\boldsymbol{H}}^a, \boldsymbol{H}^{\tilde{o}} \odot \bar{\boldsymbol{H}}^a \end{bmatrix} \right)$$
(2)

A scoring function $f$ linearly projects the max-pooled output of $M$ into a scalar.

• **Co-Matching BERT** replaces the LSTM encoders with BERT encoders. Specifically, it separately encodes $o, a, \tilde{o}$ with BERT. Given the encoded hidden vector sequences $\boldsymbol{H}^o, \boldsymbol{H}^a$ and $\boldsymbol{H}^{\tilde{o}}$, it follows Co-Matching LSTM's bi-attention and scoring function to compute the matching score.

## 4.3 Large Language Models

Finally, we test the performance of the recent large language models on our task, in order to verify whether the assessed commonsense knowledge and the inference skills can be well handled by these models. Specifically, we use ChatGPT in a zero-shot setting.

## 5 Experiments

We first evaluate all the proposed baselines on our datasets. Then we conduct a human study on a subset of our development data to investigate how human experts perform and the performance gap between machines and humans.

**Implementation Details.** We set learning rate of Adam to 1e$^{-3}$ for LSTM-based models and 2e$^{-5}$ for BERT-based models. The batch size various, each corresponds to the number of valid actions (up to 16 as described in data construction section). For the LSTM-based models, we use the Glove embedding (Pennington et al., 2014) with 100 dimensions. For both match LSTM, co-match LSTM and co-match BERT, we map the final matching states $M$ to 400 dimensional vectors, and pass these vectors to a final bi-directional LSTM layer with 100-dimensional hidden states.

| Method | ZUCC | | JECC | | Inference Speed (#states/sec) | #Parameters |
|---|---|---|---|---|---|---|
| | Dev Acc | Test Acc | Dev Acc | Test Acc | | |
| Random Guess | 10.66 | 16.42 | 7.92 | 8.01 | – | – |
| *Textual Entailment Baselines* | | | | | | |
| Match LSTM | 57.52 | 62.17 | 64.99 | 66.14 | 33.8 | 1.43M |
| BERT-siamese | 49.29 | 53.77 | 61.94 | 63.87 | 9.1 | 109.49M |
| BERT-concat | 64.73 | 64.48 | $67.39^{\dagger}$ | 72.16 | 0.6 | 109.48M |
| *Triple Modeling Baselines* | | | | | | |
| Co-Match LSTM | 72.34 | 75.91 | 70.01 | 71.64 | 25.8 | 1.47M |
| Co-Match BERT | 72.79 | 75.56 | 74.37 | 75.48 | 7.0 | 110.23M |
| ChatGPT* | – | – | 51.0 | – | – | – |
| Human Performance* | 96.40 | – | 92.0 | – | – | – |

Table 3: Evaluation on our datasets. ChatGPT and human performance (*) are computed on subsets of our data. BERT-concat (†) performs not well on JECC dev set, because the dev instances are longer on average. The concatenated inputs are more likely beyond BERT's length limit. **Inference speeds** of models are evaluated on the development set of our **JECC** dataset with a single V100 GPU.

All the experiments run on servers using a single Tesla V100 GPU with 32G memory for both training and evaluation. We use Pytorch 1.4.0; CUDA 10.2; Transformer 3.0.2; and Jericho 2.4.3.

## 5.1 Overall Results

Table 3 summarizes the models' accuracy on the development and test splits and the inference speed on the **JECC** development set. First, all the baselines learned decent models, achieving significantly better scores than a random guess. Second, the co-matching ones outperform their pairwise counterparts (Co-Match BERT > BERT-Siamese/-Concat, Co-Match LSTM > Match LSTM), and the co-match BERT performs consistently best on both datasets. The co-matching mechanism better addressed our datasets' underlying reasoning tasks, with a mild cost of additional inference computation overhead. Third, the co-match LSTM well balances accuracy and speed. In contrast, the BERT-concat, although still competitive on the accuracy, suffers from a quadratic time complexity on sequence lengths, prohibiting practical model learning and inference.

BERT-Concat represents recent general approaches to commonsense reasoning tasks. We manually examined the incorrect predictions and identified two error sources. First, it is challenging for the models to distinguish the structures of current/next observations and actions, especially when directly taking as input complicated concatenated strings of multiple types of elements. For example, it may not learn which parts of the inputs correspond to inventories. Second, the concatenation often makes the texts too long for BERT.

Albeit the models consistently outperform random guesses, the best development results on both datasets are still far below human-level performance. Compared to the human expert's near-perfect performance, the substantial performance gaps confirm that our datasets require important commonsense that humans always possess.

Finally, ChatGPT demonstrates a poor performance on the same subset for the human study. Given the wide range of commonsense knowledge types addressed by our **JECC**, we attribute this challenge primarily to the necessity of reasoning beyond mere knowledge facts. Consequently, we believe that leveraging more advanced prompting techniques such as Chain-of-Thought (Wei et al., 2022) may yield better results, and we leave this for future work.

**Remark on the Performance Consistency.** It seems that the BERT-Concat and co-match LSTM/BERT models achieve inconsistent results on the **ZUCC** and **JECC**. We point out that this inconsistency is mainly due to the different distributions – for the **JECC** we hope to keep a natural distribution of commonsense challenges, so we only evaluate on walkthrough tuples. To clarify, we also evaluate the three models on *all tuples* from **JECC** development games. The resulted accuracies are 59.84 (BERT-Concat), 68.58 (co-match LSTM), and 68.96 (co-match BERT), consistent with their ranks on **ZUCC**.

## 5.2 Human Evaluation

We present to the human evaluator each time a batch of tuples starting from the same observation $o_t$, together with its shuffled valid actions $A_{t+1}$ and

11233

| Dataset | Performance | | | $\frac{\triangle_{\text{BERT-LSTM}}}{\triangle_{\text{Human-LSTM}}}$ |
|---------|------|------|-------|------|
|         | LSTM | BERT | Human |      |
| *Multi-choice QA* | | | | |
| RACE | 50.4 | 66.5 | 94.5 | 37% |
| DREAM | 45.5 | 63.2 | 95.5 | 35% |
| *Commonsense Reasoning* | | | | |
| Abductive NLI | 50.8 | 68.6 | 91.4 | 44% |
| Cosmos QA | 44.7 | 67.6 | 94.0 | 46% |
| Our **ZUCC** | 72.3 | 72.8 | 96.4 | 2% |
| Our **JECC** | 70.0 | 74.4 | 92.0 | 20% |

Table 4: Improvement from LSTM to BERT.

next observations $O_{t+1}$. For **JECC**, only the walk-through action $a_{t+1}$ is given. The evaluators are asked to read the start observation $o_t$ first, then to align each $o \in O_{t+1}$ with an action $a \in A_{t+1}$. For each observation $o$, besides labeling the action's alignment, the subjects are asked to answer a secondary question: whether the provided $o_t, o$ pair is sufficient for them to predict the action. If they believe there are not enough clues and their action prediction is based on a random guess, they are instructed to answer "UNK" to the second question.

We collect human predictions on 250 **ZUCC** samples and 100 **JECC** samples. The annotations are done by one of the co-authors who have experience in interactive fiction game playing (but have *not* played the development games before). The corresponding results are shown in Table 3, denoted as *Human Performance*. The human expert performs 20% higher or more compared to the machines on both datasets.

Finally, the annotators recognized 10.0% cases with insufficient clues in **ZUCC** and 17.0% in **JECC**, indicating an upper-bound of methods without access to history observations.[6]

### 5.3 Comparison to the Other Datasets

Lastly, we compare our **JECC** with the other datasets to investigate how much we can gain by merely replacing the LSTMs with pre-trained LMs like BERT for text encoding. It is to verify that the language model pre-training does not easily capture the required commonsense knowledge. When LMs contribute less, it is more likely deeper knowledge and reasoning are required so that the dataset can potentially encourage new methodology advancement. Specifically, we computed the models' relative improvement from replacing the LSTM encoders with BERT ones to measure

how much knowledge BERT has captured in pre-training. Quantitatively, we calculated the ratio between the improvement from BERT encoders to the improvement of humans to LSTM models, $\triangle_{\text{BERT-LSTM}}/\triangle_{\text{Human-LSTM}}$. The ratio measures additional information (e.g., commonsense) BERT captures, compared to the full commonsense knowledge required to fill the human-machine gap.

Table 4 compares the ratios on different datasets. For a fair comparison, we use all the machine performance with co-matching style architectures. We compare to related datasets with co-matching performance available, either reported in their papers or leaderboards. These include Commonsense Reasoning datasets Abductive NLI (Bhagavatula et al., 2019) and Cosmos QA (Huang et al., 2019), and the related Multi-choice QA datasets RACE (Lai et al., 2017) and DREAM (Sun et al., 2019). Our datasets have significantly smaller ratios, indicating that much of the required knowledge in our datasets has not been captured in BERT pre-training.

## 6 Conclusion

Interactive Fiction (IF) games encode plentiful and diverse commonsense knowledge of the physical world. In this work, we derive commonsense reasoning benchmarks **JECC** and **ZUCC** from IF games in the *Jericho Environment*. Taking the form of predicting the most likely observation when applying an action to a game state, our automatically generated benchmark covers comprehensive commonsense reasoning types such as spatial reasoning and object interaction, etc. Our experiments show that current popular neural models have limited performance compared to humans. To our best knowledge, we do not identify significant negative impacts on society resulting from this work.

## Limitations

Our dataset construction method has certain limitations. One important limitation is that it is difficult to get the distribution of the required commonsense knowledge types. This can be addressed in future work with human designed commonsense knowledge schema and human annotation.

One potential risk of our work is that the text games may be limited by the time of writing, thus raise fairness considerations. However, our dataset construction strategy is not limited to these specific games, better sampling games can help to reduce such biases.

---

[6]Humans can still make a correct prediction by first eliminating most irrelevant options then making a random guess.

# References

Prithviraj Ammanabrolu and Matthew Hausknecht. 2020. Graph constrained reinforcement learning for natural language action spaces. *arXiv*, pages arXiv–2001.

Chandra Bhagavatula, Ronan Le Bras, Chaitanya Malaviya, Keisuke Sakaguchi, Ari Holtzman, Hannah Rashkin, Doug Downey, Wen-tau Yih, and Yejin Choi. 2019. Abductive commonsense reasoning. In *International Conference on Learning Representations*.

Yonatan Bisk, Rowan Zellers, Ronan LeBras, Jianfeng Gao, and Yejin Choi. 2020. Piqa: Reasoning about physical commonsense in natural language. In *AAAI*, pages 7432–7439.

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

Xiaoxiao Guo, Mo Yu, Yupeng Gao, Chuang Gan, Murray Campbell, and Shiyu Chang. 2020. Interactive fiction game playing as multi-paragraph reading comprehension with reinforcement learning. *arXiv preprint arXiv:2010.02386*.

David Ha and Jürgen Schmidhuber. 2018. World models. *arXiv preprint arXiv:1803.10122*.

Matthew Hausknecht, Prithviraj Ammanabrolu, Marc-Alexandre Côté, and Xingdi Yuan. 2019. Interactive fiction games: A colossal adventure. *arXiv preprint arXiv:1909.05398*.

Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. Cosmos qa: Machine reading comprehension with contextual commonsense reasoning. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2391–2401.

Filip Ilievski, Alessandro Oltramari, Kaixin Ma, Bin Zhang, Deborah L McGuinness, and Pedro Szekely. 2021. Dimensions of commonsense knowledge. *Knowledge-Based Systems*, 229:107347.

Minqi Jiang, Jelena Luketina, Nantas Nardelli, Pasquale Minervini, Philip HS Torr, Shimon Whiteson, and Tim Rocktäschel. 2020. Wordcraft: An environment for benchmarking commonsense agents. *arXiv preprint arXiv:2007.09185*.

Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. Race: Large-scale reading comprehension dataset from examinations. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 785–794.

Hector Levesque, Ernest Davis, and Leora Morgenstern. 2012. The winograd schema challenge. In *Thirteenth International Conference on the Principles of Knowledge Representation and Reasoning*. Citeseer.

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

James Mullenbach, Jonathan Gordon, Nanyun Peng, and Jonathan May. 2019. Do nuclear submarines have nuclear captains? a challenge dataset for commonsense reasoning over adjectives and objects. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6054–6060.

Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.

Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A Smith, and Yejin Choi. 2019a. Atomic: An atlas of machine commonsense for if-then reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3027–3035.

Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019b. Social iqa: Commonsense reasoning about social interactions. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4453–4463.

Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. 2019. Mastering atari, go, chess and shogi by planning with a learned model. *arXiv preprint arXiv:1911.08265*.

Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: an open multilingual graph of general knowledge. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, pages 4444–4451.

Kai Sun, Dian Yu, Jianshu Chen, Dong Yu, Yejin Choi, and Claire Cardie. 2019. Dream: A challenge data set and models for dialogue-based reading comprehension. *Transactions of the Association for Computational Linguistics*, 7:217–231.

Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. Commonsenseqa: A question answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for*

*Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158.

Haoyu Wang, Mo Yu, Xiaoxiao Guo, Rajarshi Das, Wenhan Xiong, and Tian Gao. 2019. Do multi-hop readers dream of reasoning chains? *arXiv preprint arXiv:1910.14520*.

Shuohang Wang and Jing Jiang. 2016. Machine comprehension using match-lstm and answer pointer. *arXiv preprint arXiv:1608.07905*.

Shuohang Wang, Mo Yu, Jing Jiang, and Shiyu Chang. 2018. A co-matching model for multi-choice reading comprehension. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 746–751.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*.

Shunyu Yao, Karthik Narasimhan, and Matthew Hausknecht. 2021. Reading and acting while blindfolded: The need for semantics in text game agents. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3097–3102.

Shunyu Yao, Rohan Rao, Matthew Hausknecht, and Karthik Narasimhan. 2020. Keep calm and explore: Language models for action generation in text-based games. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 8736–8754.

Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. 2018. Swag: A large-scale adversarial dataset for grounded commonsense inference. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 93–104.

Ben Zhou, Daniel Khashabi, Qiang Ning, and Dan Roth. 2019. "going on a vacation" takes longer than "going for a walk": A study of temporal commonsense understanding. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3354–3360.

## A   For every submission:

☑ A1. Did you describe the limitations of your work?
*In the limitation section.*

☑ A2. Did you discuss any potential risks of your work?
*In the limitation section.*

☑ A3. Do the abstract and introduction summarize the paper's main claims?
*Left blank.*

☒ A4. Have you used AI writing assistants when working on this paper?
*Left blank.*

## B   ☑ Did you use or create scientific artifacts?

*Section 3.*

☑ B1. Did you cite the creators of artifacts you used?
*Section 3.*

☑ B2. Did you discuss the license or terms for use and / or distribution of any artifacts?
*Section 3.*

☑ B3. Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)?
*Our dataset will be released for public research under CC-BY 4.0.*

☐ B4. Did you discuss the steps taken to check whether the data that was collected / used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect / anonymize it?
*Not applicable. We build our dataset on top of text games collected in Jericho, which do not have personal information.*

☑ B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?
*The first paragraph in Section 3 provide the list games.*

☑ B6. Did you report relevant statistics like the number of examples, details of train / test / dev splits, etc. for the data that you used / created? Even for commonly-used benchmark datasets, include the number of examples in train / validation / test splits, as these provide necessary context for a reader to understand experimental results. For example, small differences in accuracy on large test sets may be significant, while on small test sets they may not be.
*Table 1.*

## C   ☑ Did you run computational experiments?

*Section 5.*

☑ C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used?
*Section 5 and Table 2.*

---

*The Responsible NLP Checklist used at ACL 2023 is adopted from NAACL 2022, with the addition of a question on AI writing assistance.*

☑ C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?
*Section 5.*

☒ C3. Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run?
*According to the numbers in Table 2. Different methods have clear performance gaps between them (and between a model and humans).*

☑ C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation), did you report the implementation, model, and parameter settings used (e.g., NLTK, Spacy, ROUGE, etc.)?
*Section 3.2.*

**D  ☑ Did you use human annotators (e.g., crowdworkers) or research with human participants?**

*Section 5.2.*

☐ D1. Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.?
*Not applicable. We only use human study to compute the human accuracy. The annotators are the paper authors who have not seen the data before.*

☐ D2. Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)?
*Not applicable. We only use human study to compute the human accuracy. The annotators are the paper authors who have not seen the data before.*

☐ D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating? For example, if you collected data via crowdsourcing, did your instructions to crowdworkers explain how the data would be used?
*Not applicable. We only use human study to compute the human accuracy. The annotators are the paper authors who have not seen the data before.*

☐ D4. Was the data collection protocol approved (or determined exempt) by an ethics review board?
*Not applicable. We only use human study to compute the human accuracy.*

☐ D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?
*Not applicable. We only use human study to compute the human accuracy. The annotators are the paper authors who have not seen the data before.*