
JECC: Commonsense Reasoning Tasks Derived from Interactive Fictions

Mo Yu¹ Xiaoxiao Guo¹ Yufei Feng² Yi Gu¹
Xiaodan Zhu² Michael Greenspan² Murray Campbell¹ Chuang Gan¹
¹ IBM Research ² Queens University
yum@us.ibm.com xiaoxiao.guo@ibm.com feng.yufei@queensu.ca

Abstract

1 Commonsense reasoning simulates the human ability to make presumptions about
2 our physical world, and it is an essential cornerstone in building general AI systems.
3 We propose a new commonsense reasoning dataset based on human’s Interactive
4 Fiction (IF) gameplay walkthroughs as human players demonstrate plentiful and
5 diverse commonsense reasoning. The new dataset provides a natural mixture of
6 various reasoning types and requires multi-hop reasoning. Moreover, the IF game-
7 based construction procedure requires much less human interventions than previous
8 ones. Experiments show that the introduced dataset is challenging to previous
9 machine reading models with a significant 20% performance gap compared to
10 human experts.

11 1 Introduction

12 There has been a flurry of datasets and benchmarks proposed to address natural language-based
13 commonsense reasoning [11, 27, 20, 13, 9, 15, 2, 8, 3, 16, 26]. These benchmarks usually adopt
14 a multi-choice form – with the input query and an optional short paragraph of the background
15 description, each candidate forms a statement; the task is to predict the statement that is consistent
16 with some commonsense knowledge facts.

17 These benchmarks share some limitations, as they are mostly constructed to focus on a single
18 reasoning type and require similar validation-based reasoning. First, most benchmarks concentrate
19 on a specific facet and ask human annotators to write candidate statements related to the particular
20 type of commonsense. As a result, the distribution of these datasets is unnatural and biased to a
21 specific facet. For example, most benchmarks focus on collocation, association, or other relations
22 (e.g., ConceptNet [18] relations) between words or concepts [11, 20, 13, 9]. Other examples include
23 temporal commonsense [27], physical interactions between actions and objects [3], emotions and
24 behaviors of people under the given situation [16], and cause-effects between events and states [15,
25 2, 8]. Second, most datasets require validation-based reasoning between a commonsense fact and a
26 text statement but neglect hops over multiple facts. ¹ The previous work’s limitations bias the model
27 evaluation. For example, pre-trained Language Models (LMs), such as BERT [4], well handled most
28 benchmarks. Their pre-training process may include texts on the required facts, enabling adaptation
29 to the dominating portion of commonsense validation instances. The powerful LMs with sufficient
30 capacity can fit the isolated reasoning type easily. As a result, the above limitations of previous
31 benchmarks lead to discrepancies among practical NLP tasks that require broad reasoning ability on
32 various facets.

¹Some datasets include a portion of instances that require explicit reasoning capacity, such as [2, 8, 3, 16]. But still, standalone facts can solve most such instances.

33 **Our Contribution.** We derive a new commonsense reasoning dataset from the model-based rein-
 34 *forcement learning challenge* of Interactive Fictions (IF) to address the above limitations. Recent
 35 advances [7, 1, 5] in IF games have recognized several commonsense reasoning challenges, such as
 36 detecting valid actions and predicting different actions’ effects. Figure 1 illustrates sample gameplay
 37 of the classic game *Zork1* and the required commonsense knowledge. We derive a commonsense
 38 dataset from human players’ gameplay records related to the second challenge, i.e., predicting which
 39 textual observation is most likely after applying an action or a sequence of actions to a given game
 40 state.

41 The derived dataset naturally addresses
 42 the aforementioned limitations in previous
 43 datasets. First, predicting the next obser-
 44 vation naturally requires various common-
 45 sense knowledge and reasoning types. As
 46 shown in Figure 1, a primary commonsense
 47 type is spatial reasoning, e.g., “climb the
 48 tree” makes the protagonist up on a tree.
 49 Another primary type is reasoning about ob-
 50 ject interactions. For example, keys can
 51 open locks (object relationships); “hatch
 52 egg” will reveal “things” inside the egg (ob-
 53 ject properties); “burn repellent with
 54 torch” leads to an explosion and kills the
 55 player (physical reasoning). The above in-
 56 teractions are more comprehensive than the
 57 relationships defined in ConceptNet as used
 58 in previous datasets. Second, the rich text-
 59 ual observation enables more complex rea-
 60 soning over direct commonsense validation.
 61 Due to the textual observation’s narrative
 62 nature, a large portion of the textual obser-
 63 vations are not a sole statement of the action
 64 effect, but an extended narrates about what
 65 happens because of the effect.² Third, our
 66 commonsense reasoning task formulation
 67 shares the essence of dynamics model learn-
 68 ing for model-based RL solutions related to
 69 world models and MuZero [6, 17]. As a re-
 70 sult, models developed on our benchmarks
 71 provide direct values to model-based rein-
 72 forcement learning for text-game playing.

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

81 In summary, we introduce a new commonsense reasoning dataset construction paradigm, collectively
 82 with two datasets. The larger dataset covers 29 games in multiple domains from the *Jericho Environ-*
 83 *ment* [7], named the *Jericho Environment Commonsense Comprehension task (JECC)*. The smaller
 84 dataset, aimed for the single-domain test and fast model development, includes four IF games in the
 85 *Zork Universe*, named *Zork Universe Commonsense Comprehension (ZUCC)*. We provide strong
 86 baselines to the datasets and categorize their performance gap compared to human experts.

Observation: This is a forest, with trees in all directions. To the east, there appears to be sunlight.

Action: Go east

Observation: This is a path winding through a dimly lit forest. The path heads north south here. One particularly large tree with some low branches stands at the edge of the path.

Action: Climb the tree

Observation: You are about 10 feet above the ground nestled among some large branches. The nearest branch above you is above your reach. Beside you on the branch is a small bird's nest. In the bird's nest is a large egg encrusted with precious jewels, apparently scavenged by a childless songbird.

Action: Pick up the egg +5 reward

Observation: Taken. Add [egg] to Inventory

Action: hatch egg

Observation: There is a noticeable crunch from beneath you , and inspection reveals that the egg is lying open, badly damaged. There is a golden clockwork canary nestled in the egg. it seems to have recently had a bad experience.

hundreds of steps later ...

Action: Burn repellent with torch

Observation: The can explodes and you die. A horribly smelly death: you have died. You find yourself deep within the earth in a barren prison cell.

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. Various commonsense reasoning is illustrated in the texts of observations and commands from the gameplay interaction, such as spatial relations, objective manipulation, and physical relations.

²For 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.

87 2 Related Work

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

92 **Semantic Relations between Concepts.** Most previous datasets cover the semantic relations be-
93 tween words or concepts. These relations include the concept hierarchies, such as those covered
94 by WordNet or ConceptNet, and word collocations and associations. For example, the early work
95 Winograd [11] evaluates the model’s ability to capture word collocations, associations between
96 objects, and their attributes as a pronoun resolution task. The work by [20] is one of the first datasets
97 covering the ConceptNet relational tuple validation as a question-answering task. The problem asks
98 the relation of a source object, and the model selects the target object that satisfies the relation from
99 four candidates. [13] focus on the collocations between adjectives and objects. Their task takes the
100 form of textual inference, where a premise describes an object and the corresponding hypothesis
101 consists of the object that is modified by an adjective. [9] study associations among multiple words,
102 i.e., whether a word can be associated with two or more given others (but the work does not formally
103 define the types of associations). They propose a new task format in games where the player produces
104 as many words as possible by combining existing words.

105 **Causes/Effects between Events or States.** Previous work proposes datasets that require causal
106 knowledge between events and states [15, 2, 8]. [15] takes a text generation or inference form
107 between a cause and an effect. [2] takes a similar form to ours – a sequence of two observations is
108 given, and the model selects the plausible hypothesis from multiple candidates. Their idea of data
109 construction can also be applied to include any types of knowledge. However, their dataset only
110 focuses on causal relations between events. The work of [8] utilizes multi-choice QA on a background
111 paragraph, which covers a wider range of casual knowledge for both events and statements.

112 **Other Commonsense Datasets.** [27] proposed a unique temporal commonsense dataset. The task
113 is to predict a follow-up event’s duration or frequency, given a short paragraph describing an event.
114 [3] focus on physical interactions between actions and objects, namely whether an action over an
115 object leads to a target effect in the physical world. These datasets can be solved by mostly applying
116 the correct commonsense facts; thus, they do not require reasoning. [16] propose a task of inferring
117 people’s emotions and behaviors under the given situation. Compared to the others, this task contains
118 a larger portion of instances that require reasoning beyond fact validation. The above tasks take the
119 multi-choice question-answering form.

120 **Next-Sentence Prediction.** The next sentence prediction tasks, such as SWAG [26], are also related
121 to our work. These tasks naturally cover various types of commonsense knowledge and sometimes
122 require reasoning. The issue is that the way they guarantee distractor candidates to be irrelevant
123 greatly simplified the task. In comparison, our task utilizes the IF game engine to ensure actions
124 uniquely determining the candidates, and ours has human-written texts.

125 Finally, our idea is closely related to [25], which creates a task of predicting valid actions for each IF
126 game state. [25, 24] also discussed the advantages of the supervised tasks derived from IF games for
127 natural language understanding purpose.

128 3 Dataset Construction: Commonsense Challenges from IF Games

129 We pick games supported by the *Jericho* environment [7] to construct the **JECC** dataset.³ We pick
130 games in the *Zork Universe* for the **ZUCC** dataset.⁴ We first introduce the necessary definitions in the
131 IF game domain and then describe how we construct our **ZUCC** and **JECC** datasets as the forward
132 prediction tasks based on human players’ gameplay records, followed by a summary on the improved
133 properties of our dataset compared to previous ones. The dataset will be released for public usage. It
134 can be created with our released code with MIT License.

³We 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.

⁴We pick *Zork1*, *Enchanter*, and *Sorcerer* as the training set, and the dev and sets are non-overlapping split
from *Zork3*.

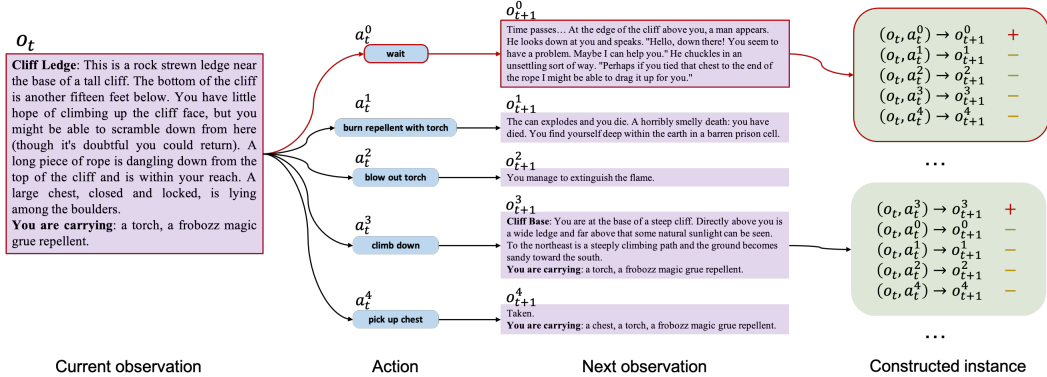


Figure 2: Illustration of our data construction process, taking an example from *Zork3*. +/−: positive/negative labels. The red colored path denotes the tuple and the resulted data instance from the human walkthrough.

Table 1: Data statistics of our **ZUCC** and **JECC** tasks. **WT** stands for walkthrough. The evaluation sets of **JECC** only consist of tuples in walkthroughs. The evaluation sets of **ZUCC** consist of all tuples after post-processing. For **JECC** the total numbers of tuples in the training games and evaluation games are close. Yet 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.

	#WT Tuples	#Tuples be- fore Proc	#Tuples af- ter 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	—	—

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 composed 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.

154 3.2 Data Construction from the Forward Prediction Task

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

164 **Data Processing.** We collect tuples from the walkthrough data provided by the Jericho environ-
165 ments. We detect the valid actions via the Jericho API and the game-specific templates. Following
166 previous work [7], we augmented the observation with the textual feedback returned by the command
167 [*inventory*] and [*look*]. The former returns the protagonist’s objects, and the latter returns the current
168 location description. When multiple actions lead to the same next-observation, we randomly keep
169 one action and next-observation in our dataset. We remove the drop *OBJ* actions since it only
170 leads to synthetic observations with minimal variety. For each step t , we keep at most 15 candidate
171 observations in O_t for the evaluation sets. When there are more than 15 candidates, we select the
172 candidate that differs most from o_t with Rouge-L measure [12].

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

184 3.3 Design Criterion and Dataset Properties

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

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

196 **Reasoning beyond verification.** A reasoning dataset should evaluate the models’ capabilities in
197 (multi-hop) reasoning with commonsense facts and background texts, beyond simple validation of
198 knowledge facts.

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

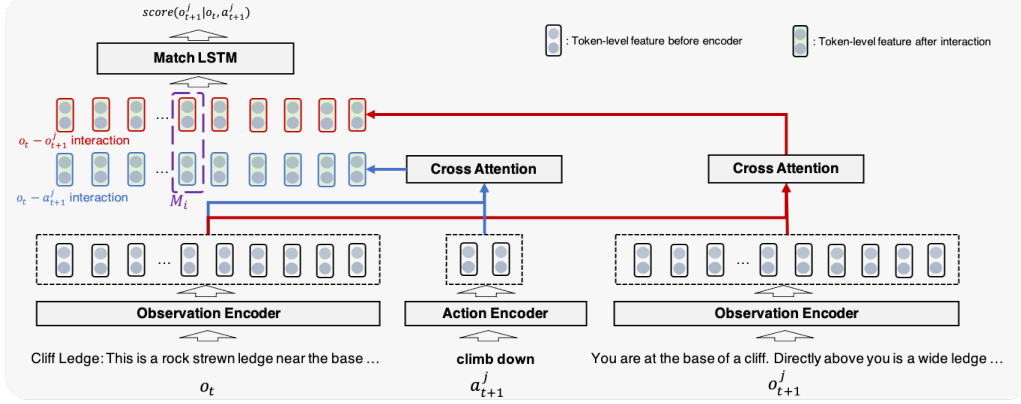


Figure 3: The co-matching architecture for our tasks.

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

209 4 Neural Inference Baselines

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

214 4.1 Neural Inference over Textual Pairs

215 • **Match LSTM** [22] represents a commonly used natural language inference model. Specifically, we
 216 concatenate o_t and a_t separated by a special split token as the premise and use the o_{t+1}^j , $j = 1, \dots, N$
 217 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
 218 premise and the hypothesis with bidirectional-LSTM model:

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

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

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

223 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.
 224 We use a scoring function $f(\cdot)$ to compute matching scores of the premise and the hypothesis via a
 225 linear transformation on the max-pooled output of M . The matching scores for all \tilde{o} are then fed to a
 226 softmax layer for the final prediction. We use the cross-entropy loss as the training objective.

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

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

230 where $[\cdot, \cdot]$ denotes concatenation. $h^{o,a}$ and $h^{\tilde{o}}$ are the last layer hidden state vectors of the “[CLS]”
 231 token. Similarly, the scoring function f computes matching scores for candidate next-observations

232 by linearly projecting the concatenated vector into a scalar. The matching scores of all \tilde{o} are grouped
 233 to a softmax layer for the final prediction.

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

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

237 The scoring function f linearly projects the last-layer hidden state of the “[CLS]” token into a scalar,
 238 and the scores are grouped to a softmax layer for final prediction. This model is much less efficient
 239 than the former two as it requires explicit combination of observation-action-next-observation as
 240 inputs. Thus this model is impractical due to the huge combinatorial space. Here we report its results
 241 for reference.

242 4.2 Neural Inference over Textual Triples

243 Existing work [10, 19, 21] has applied textual matching and entailment among triples. For example,
 244 when applying to multi-choice QA, the entailment among triples is to predict whether a question q , an
 245 answer option a can be supported by a paragraph p . In this section, we apply the most commonly used
 246 co-matching approaches [23] and its BERT variant to our task. Figure 3 illustrates our co-matching
 247 architecture.

Table 2: Evaluation on our datasets. Human performance (*) is 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.

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	–	–
<i>Textual Entailment Baselines</i>						
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†	72.16	0.6	109.48M
<i>Triple Modeling Baselines</i>						
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
Human Performance*	96.40	–	92.0	–	–	–

248 • **Co-Matching LSTM** [23] jointly encodes the question and answer with the context passage. We
 249 extend the idea to conduct the multi-hop reasoning in our setup. Specifically, similar to Equation 1,
 250 we first encode the current state observation o , the action a and the candidate next state observation
 251 \tilde{o} separately with a BiLSTM model, and use $\mathbf{H}^o, \mathbf{H}^a, \mathbf{H}^{\tilde{o}}$ to denote the output hidden vectors
 252 respectively.

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

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

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

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

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

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

262 5 Experiments

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

266 **Implementation Details.** We set learning rate of Adam to $1e^{-3}$ for LSTM-based models and $2e^{-5}$
267 for BERT-based models. The batch size varies, each corresponds to the number of valid actions
268 (up to 16 as described in data construction section). For the LSTM-based models, we use the Glove
269 embedding [14] with 100 dimensions. For both match LSTM, co-match LSTM and co-match BERT,
270 we map the final matching states M to 400 dimensional vectors, and pass these vectors to a final
271 bi-directional LSTM layer with 100-dimensional hidden states.

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

274 5.1 Overall Results

275 Table 2 summarizes the models’ accuracy on the development and test splits and the inference
276 speed on the **JECC** development set. First, all the baselines learned decent models, achieving
277 significantly better scores than a random guess. Second, the co-matching ones outperform their
278 pairwise counterparts (Co-Match BERT > BERT-Siamese/-Concat, Co-Match LSTM > Match LSTM),
279 and the co-match BERT performs consistently best on both datasets. The co-matching mechanism
280 better addressed our datasets’ underlying reasoning tasks, with a mild cost of additional inference
281 computation overhead. Third, the co-match LSTM well balances accuracy and speed. In contrast, the
282 BERT-concat, although still competitive on the accuracy, suffers from a quadratic time complexity on
283 sequence lengths, prohibiting practical model learning and inference.

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

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

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

301 5.2 Human Evaluation

302 We present to the human evaluator each time a batch of tuples starting from the same observation
303 o_t , together with its shuffled valid actions A_{t+1} and next observations O_{t+1} . For **JECC**, only the
304 walkthrough action a_{t+1} is given. The evaluators are asked to read the start observation o_t first, then
305 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

Table 3: Improvement from LSTM to BERT.

Dataset	Performance			$\frac{\Delta_{\text{BERT-LSTM}}}{\Delta_{\text{Human-LSTM}}}$
	LSTM	BERT	Human	
<i>Multi-choice QA</i>				
RACE	50.4	66.5	94.5	37%
DREAM	45.5	63.2	95.5	35%
<i>Commonsense Reasoning</i>				
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%

306 alignment, the subjects are asked to answer a secondary question: whether the provided o_t, o pair is
 307 sufficient for them to predict the action. If they believe there are not enough clues and their action
 308 prediction is based on a random guess, they are instructed to answer “UNK” to the second question.

309 We collect human predictions on 250 **ZUCC** samples and 100 **JECC** samples. The annotations are
 310 done by one of the co-authors who have experience in interactive fiction game playing (but have *not*
 311 played the development games before). The corresponding results are shown in Table 2, denoted as
 312 *Human Performance*. The human expert performs 20% higher or more compared to the machines on
 313 both datasets.

314 Finally, the annotators recognized 10.0% cases with insufficient clues in **ZUCC** and 17.0% in **JECC**,
 315 indicating an upper-bound of methods without access to history observations.⁵

316 5.3 Comparison to the Other Datasets

317 Lastly, we compare our **JECC** with the other datasets to investigate how much we can gain by
 318 merely replacing the LSTMs with pre-trained LMs like BERT for text encoding. It is to verify
 319 that the language model pre-training does not easily capture the required commonsense knowledge.
 320 When LMs contribute less, it is more likely deeper knowledge and reasoning are required so that
 321 the dataset can potentially encourage new methodology advancement. Specifically, we computed
 322 the models’ relative improvement from replacing the LSTM encoders with BERT ones to measure
 323 how much knowledge BERT has captured in pre-training. Quantitatively, we calculated the ratio
 324 between the improvement from BERT encoders to the improvement of humans to LSTM models,
 325 $\Delta_{\text{BERT-LSTM}}/\Delta_{\text{Human-LSTM}}$. The ratio measures additional information (e.g., commonsense) BERT
 326 captures, compared to the full commonsense knowledge required to fill the human-machine gap.

327 Table 3 compares the ratios on different datasets. For a fair comparison, we use all the machine
 328 performance with co-matching style architectures. We compare to related datasets with co-matching
 329 performance available, either reported in their papers or leaderboards. These include Commonsense
 330 Reasoning datasets Abductive NLI [2] and Cosmos QA [8], and the related Multi-choice QA datasets
 331 RACE [10] and DREAM [19]. Our datasets have significantly smaller ratios, indicating that much of
 332 the required knowledge in our datasets has not been captured in BERT pre-training.

333 6 Conclusion

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

⁵Humans can still make a correct prediction by first eliminating most irrelevant options then making a random guess.

References

- [1] Prithviraj Ammanabrolu and Matthew Hausknecht. Graph constrained reinforcement learning for natural language action spaces. *arXiv*, pages arXiv–2001, 2020.
- [2] Chandra Bhagavatula, Ronan Le Bras, Chaitanya Malaviya, Keisuke Sakaguchi, Ari Holtzman, Hannah Rashkin, Doug Downey, Wen-tau Yih, and Yejin Choi. Abductive commonsense reasoning. In *International Conference on Learning Representations*, 2019.
- [3] Yonatan Bisk, Rowan Zellers, Ronan LeBras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about physical commonsense in natural language. In *AAAI*, pages 7432–7439, 2020.
- [4] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 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, 2019.
- [5] Xiaoxiao Guo, Mo Yu, Yupeng Gao, Chuang Gan, Murray Campbell, and Shiyu Chang. Interactive fiction game playing as multi-paragraph reading comprehension with reinforcement learning. *arXiv preprint arXiv:2010.02386*, 2020.
- [6] David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- [7] Matthew Hausknecht, Prithviraj Ammanabrolu, Marc-Alexandre Côté, and Xingdi Yuan. Interactive fiction games: A colossal adventure. *arXiv preprint arXiv:1909.05398*, 2019.
- [8] Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 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, 2019.
- [9] Minqi Jiang, Jelena Luketina, Nantas Nardelli, Pasquale Minervini, Philip HS Torr, Shimon Whiteson, and Tim Rocktäschel. Wordcraft: An environment for benchmarking commonsense agents. *arXiv preprint arXiv:2007.09185*, 2020.
- [10] Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. Race: Large-scale reading comprehension dataset from examinations. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 785–794, 2017.
- [11] Hector Levesque, Ernest Davis, and Leora Morgenstern. The winograd schema challenge. In *Thirteenth International Conference on the Principles of Knowledge Representation and Reasoning*. Citeseer, 2012.
- [12] Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics.
- [13] James Mullenbach, Jonathan Gordon, Nanyun Peng, and Jonathan May. 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, 2019.
- [14] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.
- [15] Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A Smith, and Yejin Choi. Atomic: An atlas of machine commonsense for if-then reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3027–3035, 2019.

- 387 [16] Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. Social iqa:
388 Commonsense reasoning about social interactions. In *Proceedings of the 2019 Conference on*
389 *Empirical Methods in Natural Language Processing and the 9th International Joint Conference*
390 *on Natural Language Processing (EMNLP-IJCNLP)*, pages 4453–4463, 2019.
- 391 [17] Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Si-
392 mon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering
393 atari, go, chess and shogi by planning with a learned model. *arXiv preprint arXiv:1911.08265*,
394 2019.
- 395 [18] Robyn Speer, Joshua Chin, and Catherine Havasi. Conceptnet 5.5: an open multilingual
396 graph of general knowledge. In *Proceedings of the Thirty-First AAAI Conference on Artificial*
397 *Intelligence*, pages 4444–4451, 2017.
- 398 [19] Kai Sun, Dian Yu, Jianshu Chen, Dong Yu, Yejin Choi, and Claire Cardie. Dream: A challenge
399 data set and models for dialogue-based reading comprehension. *Transactions of the Association*
400 *for Computational Linguistics*, 7:217–231, 2019.
- 401 [20] Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. Commonsenseqa: A
402 question answering challenge targeting commonsense knowledge. In *Proceedings of the 2019*
403 *Conference of the North American Chapter of the Association for Computational Linguistics: Human*
404 *Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, 2019.
- 405 [21] Haoyu Wang, Mo Yu, Xiaoxiao Guo, Rajarshi Das, Wenhan Xiong, and Tian Gao. Do multi-hop
406 readers dream of reasoning chains? *arXiv preprint arXiv:1910.14520*, 2019.
- 407 [22] Shuohang Wang and Jing Jiang. Machine comprehension using match-lstm and answer pointer.
408 *arXiv preprint arXiv:1608.07905*, 2016.
- 409 [23] Shuohang Wang, Mo Yu, Jing Jiang, and Shiyu Chang. A co-matching model for multi-choice
410 reading comprehension. In *Proceedings of the 56th Annual Meeting of the Association for*
411 *Computational Linguistics (Volume 2: Short Papers)*, pages 746–751, 2018.
- 412 [24] Shunyu Yao, Karthik Narasimhan, and Matthew Hausknecht. Reading and acting while blind-
413 folded: The need for semantics in text game agents. In *Proceedings of the 2021 Conference*
414 *of the North American Chapter of the Association for Computational Linguistics: Human*
415 *Language Technologies*, pages 3097–3102, 2021.
- 416 [25] Shunyu Yao, Rohan Rao, Matthew Hausknecht, and Karthik Narasimhan. Keep calm and
417 explore: Language models for action generation in text-based games. In *Proceedings of the*
418 *2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages
419 8736–8754, 2020.
- 420 [26] Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. Swag: A large-scale adversarial
421 dataset for grounded commonsense inference. In *Proceedings of the 2018 Conference on*
422 *Empirical Methods in Natural Language Processing*, pages 93–104, 2018.
- 423 [27] Ben Zhou, Daniel Khashabi, Qiang Ning, and Dan Roth. “going on a vacation” takes longer
424 than “going for a walk”: A study of temporal commonsense understanding. In *Proceedings*
425 *of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th*
426 *International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages
427 3354–3360, 2019.

428 Checklist

- 429 1. For all authors...
- 430 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
431 contributions and scope? [Yes]
- 432 (b) Did you describe the limitations of your work? [Yes] At the end of Section 3.3.
- 433 (c) Did you discuss any potential negative societal impacts of your work? [Yes]

- 434 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
435 them? [Yes]
- 436 2. If you are including theoretical results...
- 437 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
438 (b) Did you include complete proofs of all theoretical results? [N/A]
- 439 3. If you ran experiments (e.g. for benchmarks)...
- 440 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
441 mental results (either in the supplemental material or as a URL)? [Yes] We include the
442 code the data in the supplemental material.
- 443 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
444 were chosen)? [Yes]
- 445 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
446 ments multiple times)? [No]
- 447 (d) Did you include the total amount of compute and the type of resources used (e.g., type
448 of GPUs, internal cluster, or cloud provider)? [Yes] See the implementation details in
449 Section 5.
- 450 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 451 (a) If your work uses existing assets, did you cite the creators? [Yes]
452 (b) Did you mention the license of the assets? [Yes]
453 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
454 (d) Did you discuss whether and how consent was obtained from people whose data you're
455 using/curating? [N/A] The data is not created via crowd-sourcing.
456 (e) Did you discuss whether the data you are using/curating contains personally identifiable
457 information or offensive content? [N/A] The data is not created via crowd-sourcing.
- 458 5. If you used crowdsourcing or conducted research with human subjects...
- 459 (a) Did you include the full text of instructions given to participants and screenshots, if
460 applicable? [Yes] See Section 5.2.
461 (b) Did you describe any potential participant risks, with links to Institutional Review
462 Board (IRB) approvals, if applicable? [N/A] Our human subjects are co-authors with
463 certain requirements, as described in Section 5.2.
464 (c) Did you include the estimated hourly wage paid to participants and the total amount
465 spent on participant compensation? [N/A]