

# CORECODE: A Common Sense Annotated Dialogue Dataset with Benchmark Tasks for Chinese Large Language Models

Dan Shi<sup>1</sup>, Chaobin You<sup>1</sup>, Jiantao Huang<sup>2</sup>, Taihao Li<sup>2</sup>, Deyi Xiong<sup>1\*</sup>

<sup>1</sup>College of Intelligence and Computing, Tianjin University, Tianjin, China

<sup>2</sup>Zhejiang Lab, Hangzhou, China

{shidan, chaobinyou, dyxiong}@tju.edu.cn, {jthuang, lith}@zhejianglab.com

## Abstract

As an indispensable ingredient of intelligence, commonsense reasoning is crucial for large language models (LLMs) in real-world scenarios. In this paper, we propose CORECODE, a dataset that contains abundant commonsense knowledge manually annotated on dyadic dialogues, to evaluate the commonsense reasoning and commonsense conflict detection capabilities of Chinese LLMs. We categorize commonsense knowledge in everyday conversations into three dimensions: entity, event, and social interaction. For easy and consistent annotation, we standardize the form of commonsense knowledge annotation in open-domain dialogues as “domain: slot = value”. A total of 9 domains and 37 slots are defined to capture diverse commonsense knowledge. With these pre-defined domains and slots, we collect 76,787 commonsense knowledge annotations from 19,700 dialogues through crowdsourcing. To evaluate and enhance the commonsense reasoning capability for LLMs on the curated dataset, we establish a series of dialogue-level reasoning and detection tasks, including commonsense knowledge filling, commonsense knowledge generation, commonsense conflict phrase detection, domain identification, slot identification, and event causal inference. A wide variety of existing open-source Chinese LLMs are evaluated with these tasks on our dataset. Experimental results demonstrate that these models are not competent to predict CORECODE’s plentiful reasoning content, and even ChatGPT could only achieve 0.275 and 0.084 accuracy on the domain identification and slot identification tasks under the zero-shot setting. We release the data and codes of CORECODE at <https://github.com/danshi777/CORECODE> to promote commonsense reasoning evaluation and study of LLMs in the context of daily conversations.

## 1 Introduction

Commonsense reasoning is a crucial component of intelligence (Liu and Singh 2004; Cambria et al. 2011; Storks, Gao, and Chai 2019), which involves the ability to make logical deductions, infer implicit information and apply background knowledge to solve problems as well as understand the world. In recent years, exploring and improving the ability of NLP models for the acquisition and application of

commonsense knowledge has been attracting growing interest, leading to extensive research in this field (Lin et al. 2019; Bauer, Wang, and Bansal 2018; Lv et al. 2020; Wang et al. 2020; Liu et al. 2021; Jiang et al. 2021; Liu et al. 2022).

It is widely acknowledged that LLMs, trained on a huge amount of data, are able to obtain broad knowledge covering a wide range of domains (Rae et al. 2021; Hoffmann et al. 2022; Touvron et al. 2023; Du et al. 2022a; Guo et al. 2023), including commonsense knowledge (West et al. 2022; Bian et al. 2023; Bang et al. 2023). However, commonsense reasoning is still regarded as a major challenge for LLMs (Zhou et al. 2020; Bhargava and Ng 2022). Studies disclose that LLMs fall short in performing adequate commonsense reasoning (Wei et al. 2022). For example, ChatGPT<sup>1</sup> does not precisely know what the needed commonsense knowledge for answering a specific question is (e.g., questions in social and temporal domains) (Bian et al. 2023).

To mitigate this issue, we propose CORECODE (Commonsense Reasoning and Conflict Detection in dialogues), a dataset that contains abundant commonsense knowledge manually annotated on Chinese dyadic dialogues, to assess how much commonsense knowledge the LLMs have gained and how well they can be improved in commonsense reasoning and conflict detection with the annotated knowledge in CORECODE.

Specifically, we focus on annotating fine-grained commonsense knowledge in multi-turn dyadic dialogues. The knowledge annotated in a dialogue is context-sensitive and grounded exclusively in that particular dialogue. Inspired by the annotation convention used in task-oriented dialogue, in which dialogue states are denoted in the form of “domain: slot = value”, e.g. “hotel: price range = moderate” (Budzianowski et al. 2018; Zhu et al. 2020; Quan et al. 2020), we standardize the representation of commonsense knowledge in open-domain dialogues also in the form of “domain: slot = value”. We categorize commonsense knowledge into three dimensions, namely entity, event, and social interaction, and then construct an ontology over these dimensions, which defines all possible domains for each dimension and all possible slots for each domain. Thanks to the guidance of this ontology, crowdsourcing annotators are able to conveniently annotate fine-grained commonsense

\* Corresponding author

<sup>1</sup><https://openai.com/blog/chatgpt>

knowledge in a consistent way.

Over the curated dataset, we develop six benchmark tasks: commonsense knowledge filling, commonsense knowledge generation, commonsense conflict phrase detection, domain identification, slot identification and event causal inference. These tasks, organized in different forms (e.g., multiple-choice questions, span extraction, text generation), facilitate the evaluation and enhancement of commonsense reasoning in LLMs.

We conduct numerous experiments on CORECODE, attempting to explore two main research questions: (1) Can LLMs master and apply commonsense knowledge well enough to achieve good performance on these tasks? (2) How much further improvements can be obtained by LLMs if they are fine-tuned on CORECODE? Extensive experiments demonstrate that our benchmark tasks are challenging for existing Chinese LLMs, as all evaluated LLMs perform poorly on most tasks. We also show that although the performance of LLMs improves after being fine-tuned on CORECODE, they fail to obtain robust commonsense reasoning ability. When perturbations are introduced, the fine-tuning performance has significantly dropped.

## 2 Related Work

A variety of datasets and benchmarks focusing on different aspects of commonsense knowledge over textual inputs have been proposed, including science common sense datasets ARC (Clark et al. 2018) and QASC (Khot et al. 2020), temporal common sense dataset MC-TACO (Zhou et al. 2019), numerical common sense dataset NumerSense (Lin et al. 2020), event common sense dataset HellaSWAG (Zellers et al. 2019), physical common sense dataset PIQA (Bisk et al. 2020), social common sense dataset Social IQA (Sap et al. 2019b) and general common sense datasets CommonsenseQA (Talmor et al. 2019), OpenBookQA (Mihaylov et al. 2018), and WSC (Levesque, Davis, and Morgenstern 2012). These datasets only examine the model’s knowledge and ability in a certain commonsense aspect in the form of multiple-choice questions.

Meanwhile, there have also been many studies devoted to annotating commonsense knowledge involved in utterances in dialogues. ATOMIC (Sap et al. 2019a; Hwang et al. 2021) is one such dataset that consists of a large set of inference types. However, ATOMIC is context-insensitive, as its commonsense reasoning operates on phrases taken out of context, disregarding whether an event is performed by the same individual. TIMEDIAL (Qin et al. 2021) focuses on the time reasoning ability of language models in dialogues, while CICERO (Ghosal et al. 2022) provides cause, subsequent events, prerequisites, motivations, and emotional reactions for utterances in dialogues, focusing on these five event-related reasoning types. Both datasets cover only a specific aspect of commonsense knowledge. CIDER (Ghosal et al. 2021) extracts knowledge in dialogues into knowledge triplets, which covers fewer commonsense knowledge types than us. For example, *subsequent event*, *subsequent emotional reaction*, *frequency* are beyond the scope of CIDER.

To the best of our knowledge, CORECODE is the first large-scale Chinese dialogue-oriented commonsense knowl-

edge annotation dataset involving comprehensive commonsense knowledge in three dimensions: entity, event, and social interaction, covering a large number of perspectives such as attributes, time, space, and causality. Yet another feature that must be emphasized is that within CORECODE, we manually provide phrases corresponding to the phrases in an original dialogue, which are against common sense in that context. This aims to probe the model’s capacity to detect and locate such phrases that are inconsistent with the context in terms of commonsense reasoning.

## 3 Dataset Creation

The raw data of CORECODE is derived from NaturalConv (Wang et al. 2021) and DuLeMon (Xu et al. 2022) datasets, both of which contain multi-turn dialogues between two people. Dialogues in NaturalConv involve a variety of topics (including but not limited to sports, entertainment, and technology). We first take an automatic screening method to identify dialogues that are rich in commonsense knowledge, following Zhou et al. (2021).

Specifically, we first identify candidate concepts (nouns, verbs, adjectives) in each turn of a dialogue using part-of-speech tagging. We then query the ConceptNet using the identified concepts in each utterance to obtain a list of one-hop commonsense triples of the form  $(e_1, r, e_2)$ . Next, we check if entity  $e_2$  from the triple appears in the concept set of the succeeding utterance. If there is a match, it indicates a potential commonsense link between the two utterances.

Unlike Zhou et al. (2021) who retain dialogues with only one commonsense triple match, we employ a stricter criterion by retaining dialogues where more than three commonsense triple matches are detected. This selection ensures that the kept dialogues possess a substantial amount of commonsense reasoning. The statistics of the screening results on NaturalConv and DuLeMon are shown in Appendix A.

Moreover, to differentiate between the two sides of the conversation, we employ the notation “A: ” or “B: ” preceding each utterance to denote the respective speaker.

## Data Annotation

Over the selected dialogues, we perform commonsense knowledge annotation. To guarantee the consistency of annotations across multiple crowd-sourced workers, we adopt a standardized annotation procedure.

We categorize commonsense knowledge in everyday conversations into three dimensions: *entity*, *event*, and *social interaction*. Crowd-sourced workers first need to identify specific instances under these three dimensions from dialogues. Then, with the assistance of linguists, we divide each of these three dimensions into multiple domains to which their commonsense knowledge belongs, and define different slots for each domain, forming a two-level hierarchical taxonomy. Such design is guided by three fundamental principles: *coverage*, *exclusivity*, and *easiness*. The *coverage* rule ensures that the commonsense knowledge system encompasses nearly all conceivable types of commonsense knowledge in dialogues. *Exclusivity* mandates that each commonsense knowledge type remains distinct, devoid

of any overlap with other types. Lastly, the *easiness* principle indicates that the commonsense knowledge system is straightforward for annotators to employ. With this convention, crowd-sourced workers are instructed to annotate the identified instances with commonsense knowledge in the form of “domain: slot = value”. In addition to such annotations, they are also required to provide phrases that, in terms of common sense, conflict with the original textual context. Below, we describe each step in detail.

**Entity, Event, Social Interaction Recognition.** The first step of the annotation process is to identify specific instances of entity, event, and social interaction that exist in dialogues, according to the following definitions.

- **Entities** refer to objectively existing and distinguishable physical objects in the real world, either representing a general category of people or things, such as “*cats*”, “*movies*”, or referring to specific individuals or objects, such as “*Yao Ming*”, “*Wolf Warrior*”, etc.
- **Events** are typically text spans in the form of “subject + predicate” or “subject + predicate + object”. They are fine-grained semantic units that describe the state of entities and their actions (Zhou et al. 2022). For example, “*He looks very excited*” describes the state of the subject, and “*He broke his toy*” illustrates an action where the subject interacts with the object.
- **Social interactions** refer to the set of rules and guidelines that constrain people’s behavior when interacting with others. It encompasses a collection of social norms and customs that people are expected to adhere to (Bian et al. 2023). For instance, “*It is customary to knock on the door before entering someone else’s room*”.

**Annotation of Involved Commonsense Knowledge.** Under each of the three dimensions, we define domains and slots. For entities, we divide the relevant commonsense knowledge into three corresponding domains: attribute, comparison, and space. These domains capture specific properties of the object itself, relationships between the object and other objects, and relationships between the object and the spatial environment in which it is located, respectively. Under each domain, there are further divisions into different slots. For example, under the attribute domain, there are slots “Is”, “Is A”, “Has”, “Is Made Of”, and so on. For events, the relevant commonsense knowledge includes the prerequisite, cause, and consequence of an event, as well as the temporal and spatial factors associated with the event. For social interactions, we focus on the social norms that humans follow. Instead of subdividing into multiple domains, we divide seven slots under the social norms domain. There are 9 domains and 37 slots included in the three dimensions in total. The full inventory of all domains and slots is listed in Appendix B.

The second step of the annotation process is to label each entity, event or social interaction instance with its commonsense knowledge in the form of “domain: slot = value”. The annotated “value” does not necessarily need to be an exact span extracted from the original dialogue, but can be a grammatically correct and semantically fluent clause summarized

from the dialogue, ensuring that the event and its “domain: slot = value” in isolation is informationally complete and logically consistent. It has been emphasized to annotators that for the “event cause” slot in the “cause” domain and the “subsequent event” slot in the “consequence” domain, they should be annotated in the form of an event, i.e., in the “subject + predicate” form or “subject + predicate + object” form.

In addition, the annotators need to indicate which phrases or clauses in the original dialogue led to the identification of this commonsense knowledge, so as to provide a basis for the next step.

**Rewriting of Commonsense Conflict Phrases.** Finally, for each set of phrases from the original dialogue indicated in the previous step, annotators are required to choose one phrase and provide it with the following two commonsense conflict phrases:

(1) **Commonsense Conflict Phrase 1:** This phrase should be obtained by conforming to the minimal modification principle, i.e., modifying only one or two words in the original phrase. There should be a commonsense conflict or error after using this phrase to replace the original phrase in the dialogue.

(2) **Commonsense Conflict Phrase 2:** This phrase should be created by modifying as many words as possible in the original phrase in compliance with the maximum modification principle. When constructing this phrase, annotators can include words that appear in the dialogue to maintain consistency with the dialogue’s context. However, it is crucial to ensure as much as possible that the meaning of this phrase differs from the Commonsense Conflict Phrase 1.

The purpose of this annotation step is to explore whether LLMs are able to detect the location of phrases that conflict with the dialogue context in terms of common sense. Therefore, annotators must ensure that after replacing the original phrase in the dialogue with the annotated conflict phrase, there should be only a commonsense error while the dialogue maintains grammatically correct and fluent.

To comprehensively evaluate the commonsense reasoning ability of LLMs, we propose two distinct annotated subsets with varying difficulty levels. During the annotation procedure on 9.7K dialogues, we represent the subject and object of events using the speaker indicators “A” or “B” from the dialogue and group these annotated instances as an EASY set. A HARD set is annotated on another 10K dialogues, where “x” is uniformly employed to denote the subject of all events, while “y” is used to represent the predicate of all events, regardless of the dialogue participant to whom the event pertains. Significant challenges in reasoning through events are provided in the HARD set, as LLMs are required to first deduce and locate the event initiator before reasoning.

## Annotation Quality Control

In order to standardize the annotation form and control the quality of common sense annotations, we design and develop a knowledge acquisition platform where crowd-sourced workers need to properly click on the appropriate

buttons and fill in the corresponding values given the dialogue history.

We adopt a very strict quality control protocol to ensure the quality of annotations. First, we train two reviewers with 200 dialogues. The annotation consistency of the two reviewers is high, with an average Cohen’s Kappa (McHugh 2012) of 80.7% across the annotation tasks. We only hire annotators who have relevant experience in text annotation, e.g., those who have participated in annotation tasks such as Chinese multi-turn dialogue writing and correction, entity extraction or syntactic structure annotation in Chinese texts.

Second, 200 candidate workers participate in a pre-annotation stage. They adhere to the prescribed rules to annotate dialogues. The two reviewers will review annotations of these participants to distinguish whether the annotations meet the requirements. The process has an elimination rate of roughly 80%, with 43 labelers passing this stage.

Third, we proceed to the training phase. We divided the participants into groups of 5 people each. We train 1-2 quality inspectors within each group, who in turn are responsible for the instruction of the annotators. During this progression, quality inspectors evaluate the rule comprehension and error correction capabilities of the annotators. Those who do not meet the criteria are subjected to further training or eliminated from the process.

At last, 6 quality inspectors with an average Cohen’s Kappa of 59.4%, as well as 15 annotators, proceed to the formal annotation stage. We take iterative verification and revision during this stage. Any data deemed unsatisfactory will be returned for revision until they are qualified.

## Overall Statistics

The overall statistics of the annotated dataset are shown in Table 1. After annotating on 19.7K dialogues, we obtain 76,787 annotations, each consisting of the original dialogue, an entity/event/social interaction instance, a commonsense knowledge represented by a domain-slot-value triplet, the involved phrase from the original dialogue, and two commonsense knowledge conflict phrases. The average number of turns and tokens per dialogue is 19.40 and 501.58, indicating that the dialogues annotated are quite long and informative. Since the knowledge in the social interaction dimension is mainly used to constrain our behavior, but is rarely mentioned in our dialogues, there is little commonsense knowledge annotated for this dimension. The three dimensions of entity, event and social interaction annotations account for 58.42%, 41.54% and 0.03% of overall annotations, respectively.

## 4 Benchmark Tasks

We use our dataset as a testbed and define 6 tasks in different forms, attempting to evaluate dialogue-level commonsense reasoning capabilities of Chinese LLMs. For each task, we provide both its definition and associated prompt that is constructed to allow LLMs to complete the task in the continuation to the prompt.

	HARD	EASY	Total
# dialogues	10,000	9,700	19,700
Max. turns per dialogue	26	26	26
Min. turns per dialogue	14	16	15
Avg. turns per dialogue	18.69	20.10	19.40
Max. # tokens per dialogue	1,002	953	977.5
Min. # tokens per dialogue	194	231	212.5
Avg. # tokens per dialogue	464.18	538.98	501.58
Avg. # tokens per turn	24.83	26.81	25.82
# annotated instances	37,777	39,010	76,787
# annotated entities	21,320	23,541	44,861
# annotated events	16,439	15,461	31,900
# annotated social interactions	18	8	26
# domain-slot-value triplets	37,777	39,010	76,787
# commonsense conflict phrases	75,554	78,020	153,574
Avg. # tokens of conflict phrase 1	4.81	5.31	5.06
Avg. # tokens of conflict phrase 2	4.74	4.91	4.83

Table 1: Overall statistics of the CORECODE dataset.

### Commonsense Knowledge Filling

**Task definition.** This task is to fill desirable commonsense knowledge into a masked dialogue where a commonsense phrase is replaced with [MASK]. In order to automatically assess the performance of the task, we formulate the task in the form of multiple-choice questions.

**Prompt.** The input prompt to LLMs for this task consists of the question, masked dialogue, answer choices, and suffix: question \n masked dialogue \n (a) *phrase*<sub>1</sub> (b) *phrase*<sub>2</sub> (c) *phrase*<sub>3</sub> \n “answer: the correct option is”. The three phrases are the corresponding masked commonsense phrase and two manually composed commonsense conflict phrases. See Appendix C for examples of all tasks.

### Commonsense Knowledge Generation

**Task definition.** We frame this task as a generative task that takes the annotated commonsense knowledge values as ground truth and asks LLMs to generate the values according to the dialogue context.

**Prompt.** The input prompt is formatted as: dialogue \n question \n “answer:”, where the question is formed by the entity/event/social interaction and its annotated slot through a predefined template and some explanatory text.

### Commonsense Conflict Phrase Detection

**Task definition.** We define this task as a span extraction task. We replace the corresponding phrases in the original dialogue with the annotated commonsense conflict phrases, and ask LLMs to extract the commonsense conflict phrases.

**Prompt.** The input prompt is in the form of: dialogue with replaced commonsense conflict phrases \n question \n “answer:”.

### Domain Identification

**Task definition.** This task is also defined as a multiple-choice question task. Take the entity dimension as an example, LLMs are required to select the domain to which the

relationship between an entity and its annotated value belongs, based on the given dialogue context. Since the social interaction dimension includes a single domain, this task is performed on the entity and event dimensions.

**Prompt.** The prompt is formatted like: question \n entity or event \n annotated value \n dialogue \n (a) *domain*<sub>1</sub> (b) *domain*<sub>2</sub> ··· (x) *domain*<sub>n</sub> \n “answer: the correct domain is”.

### Slot Identification

**Task definition.** This task is similar to the Domain Identification task, except that this task is chosen from more fine-grained slot options and executed on all the three dimensions.

**Prompt.** The format of the prompt is: question \n entity or event or social interaction \n annotated value \n dialogue \n (a) *slot*<sub>1</sub> (b) *slot*<sub>2</sub> ··· (x) *slot*<sub>n</sub> \n “answer: the correct option is”.

### Event Causal Inference

Causal inference is one of the crucial reasoning abilities of human intelligence, which involves establishing the correct cause-and-consequence relationships between events. These relationships are captured in the “cause: event cause” slot and the “consequence: subsequent event” slot of our taxonomy. We specially design three generative event causal inference tasks that utilize the annotated knowledge involved in these two slots.

- **Subtask 1: Event Cause Inference.** Given the dialogue and event, LLMs are required to generate the cause of the event.
- **Subtask 2: Subsequent Event Inference.** Given the dialogue and event, the consequence of the event is generated by LLMs.
- **Subtask 3: Clipped Subsequent Event Inference.** Given the event and the truncated dialogue where the context succeeding the event is discarded, we require LLMs to generate the consequence of the event.

## 5 Experiments

### Evaluated LLMs

We evaluated a diverse list of Chinese LLMs that cover a variety of training processes and scales<sup>2</sup>: (1) LLMs only being pre-trained on large-scale training corpora, including GLM-10B (Du et al. 2022b) and BLOOM-7.1B (Scao et al. 2022), (2) LLMs being both pre-trained and instruction-tuned, including ChatGLM-6B<sup>3</sup>, ChatGLM2-6B<sup>4</sup>, MOSS-SFT-16B<sup>5</sup>, Baichuan-7B<sup>6</sup>, BLOOMZ-1.7B,

<sup>2</sup>All the experiments in the main paper were conducted on the HARD set. Experimental results on the EASY set are available in Appendix F.

<sup>3</sup><https://github.com/THUDM/ChatGLM-6B>

<sup>4</sup><https://github.com/THUDM/ChatGLM2-6B>

<sup>5</sup><https://huggingface.co/chnlp/moss-moon-003-sft>

<sup>6</sup><https://github.com/baichuan-inc/baichuan-7B>

BLOOMZ-7.1B, BLOOMZ-7.1B-MT (Muennighoff et al. 2022), and BELLE-7B, which is the SFT version based on BLOOMZ-7.1B-MT. We used two variants of BELLE fine-tuned on 200K and 2M instructions separately, i.e., BELLE-7B-0.2M<sup>7</sup> and BELLE-7B-2M<sup>8</sup>. We also evaluated two variants Chinese-Alpaca-Plus-7B and Chinese-Alpaca-Plus-13B of Chinese-Alpaca-Plus (Cui, Yang, and Yao 2023). We conducted experiments using the recommended hyperparameter settings for all LLMs. We also evaluated ChatGPT<sup>9</sup> (i.e., GPT-3.5-turbo) from OpenAI as a reference.

Furthermore, to explore the impact of in-context learning (ICL) on model performance, we also carried out experiments on ChatGLM-6B under the few-shot settings, including 1-shot, 3-shot and 5-shot settings.

### Evaluation Metrics

For the commonsense knowledge filling, domain identification and slot identification tasks (we refer these three tasks to the selection tasks), we used the accuracy of selecting the correct answer as the evaluation metric. During inference, we have found that even if we explicitly state in the prompt that models should output only the answer option indicator (i.e. a, b, c, etc.), not all models follow this instruction. There is no uniformity in the form of answers generated by each model. Moreover, sometimes models output answers with rationales attached. In order to avoid the underestimation of the model performance due to the varying output formats, we adopted a series of filtering measures to find the correct answer in the output as much as possible. For example, in the case where the ground-truth is “(a) premise”, the generated answers “a”, “A”, “(a)”, “(A)”, “a)”, “A)”, “(a)premise”, “(a) premise”, “premise” are all counted as correct. We conducted a quantitative analysis experiment to showcase the effectiveness of the filtering measures. Detailed results of this experiment can be found in Appendix D.

For the span extraction task, i.e., the commonsense conflict phrase detection task, we used F1 and exact match (EM) scores calculated by comparing model outputs to ground-truth answers.

For the two generation tasks, namely the commonsense knowledge generation task and the event causal inference task, we evaluated LLMs with F1 and EM scores together with reference based metrics: BLEU (Papineni et al. 2002), METEOR (Banerjee and Lavie 2005), ROUGE (Lin 2004) and CIDEr (Vedantam, Lawrence Zitnick, and Parikh 2015).

### Performance of LLMs without Fine-tuning

We report the performance of the selection tasks in Table 3, and the performance of the span extraction task and generation tasks in Table 2. We can see that CoRECODE is a very challenging benchmark for all evaluated LLMs.

From Table 3, we observe that models which are instruction-tuned with SFT significantly outperform models being only pre-trained. The best-performing models across the three tasks are ChatGLM2-6B, BLOOMZ-7.1B, and

<sup>7</sup><https://huggingface.co/BelleGroup/BELLE-7B-0.2M>

<sup>8</sup><https://huggingface.co/BelleGroup/BELLE-7B-2M>

<sup>9</sup><https://openai.com/chatgpt>

Model	Commonsense Knowledge Generation								Commonsense Conflict Phrase Detection	
	F1	EM	BLEU1	BLEU2	METEOR	ROUGE-L	CIDEr	ChatGPT Score	F1	EM
GLM-10B	0.023	0.000	0.000	0.000	0.032	0.000	0.001	3.190	0.011	0.000
BLOOM-7.1B	0.071	0.000	0.017	0.000	0.115	0.004	0.017	3.455	0.024	0.000
ChatGLM2-6B	0.160	0.004	0.001	0.000	0.145	0.001	0.002	3.940	0.029	0.001
BELLE-7B-0.2M	0.090	0.019	0.015	0.000	0.105	0.010	0.041	3.265	0.024	0.000
BELLE-7B-2M	0.111	0.008	0.004	0.000	0.140	0.003	0.010	3.555	0.007	0.000
BLOOMZ-1.7B	0.388	0.234	0.234	0.000	0.164	0.234	0.585	4.060	0.004	0.000
BLOOMZ-7.1B	0.438	0.284	0.282	0.000	0.199	0.283	0.707	3.980	0.041	0.003
BLOOMZ-7.1B-MT	0.435	0.300	0.300	0.000	0.184	0.300	0.750	4.030	0.047	0.010
MOSS-SFT-16B	0.199	0.071	0.066	0.000	0.147	0.049	0.174	3.780	0.038	0.001
Baichuan-7B	0.071	0.000	0.000	0.000	0.072	0.000	0.001	3.445	0.002	0.000
Chinese-Alpaca-Plus-7B	0.129	0.015	0.014	0.000	0.099	0.015	0.039	3.375	0.021	0.000
Chinese-Alpaca-Plus-13B	0.133	0.021	0.018	0.000	0.104	0.020	0.051	3.490	0.021	0.000
ChatGLM-6B	0.147	0.000	0.000	0.000	0.166	0.000	0.000	3.745	0.044	0.001
ChatGLM-6B 1-shot	0.202	0.061	0.035	0.000	0.154	0.048	0.120	3.770	0.038	0.002
ChatGLM-6B 3-shot	0.274	0.115	0.091	0.000	0.175	0.111	0.277	3.885	0.060	0.006
ChatGLM-6B 5-shot	0.215	0.097	0.095	0.000	0.147	0.096	0.240	3.685	0.052	0.007
ChatGPT	0.296	0.071	0.044	0.000	0.258	0.045	0.111	-	0.104	0.021

Table 2: Overall performance of evaluated LLMs on the commonsense knowledge generation and commonsense conflict phrase detection task.

Model	CKF	DI	SI
GLM-10B	0.157	0.060	0.051
BLOOM-7.1B	0.329	0.108	0.039
ChatGLM-6B	0.788	0.246	0.113
ChatGLM2-6B	0.818	0.286	0.153
BELLE-7B-0.2M	0.392	0.208	0.212
BELLE-7B-2M	0.599	0.169	0.109
BLOOMZ-1.7B	0.709	0.248	0.044
BLOOMZ-7.1B	0.758	0.444	0.165
BLOOMZ-7.1B-MT	0.695	0.341	0.168
MOSS-SFT-16B	0.445	0.353	0.110
Baichuan-7B	0.416	0.071	0.055
Chinese-Alpaca-Plus-7B	0.584	0.385	0.060
Chinese-Alpaca-Plus-13B	0.510	0.126	0.449
ChatGPT	0.896	0.275	0.084

Table 3: Overall performance of evaluated LLMs on the three selection tasks. CKF: Commonsense Knowledge Filling. DI: Domain Identification. SI: Slot Identification.

Chinese-Alpaca-Plus-13B, respectively. Notably, on the slot identification task, Chinese-Alpaca-Plus-13B achieves an outstanding and unparalleled score.

On the commonsense knowledge generation task, BLOOMZ family achieves very high scores, as shown in Table 2. After checking the outputs of each model, we have found that models like ChatGLM and BELLE usually generate leading sentences or explanatory reasons in their responses, despite our prompt explicitly instructing them not to do so. In contrast, BLOOMZ-1.7B and BLOOMZ-7.1B typically generate relatively short phrases as answers, which is consistent with the form of our annotations. They hence achieve higher scores than other evaluated LLMs.

To exclude the effect of answer form and answer length on the performance, we handed over the outputs of evaluated LLMs to ChatGPT for scoring, the average results of which are also reported in Table 2. We described the task to ChatGPT and asked it to score the answers according to our pre-defined scoring criteria (see in Appendix E). The average

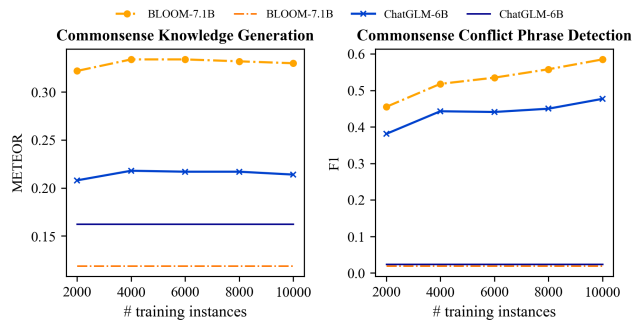


Figure 1: Performance of fine-tuned LLMs on the commonsense knowledge generation and commonsense conflict phrase detection task. The horizontal lines show the performance of LLMs without fine-tuning.

scores obtained by these models vary from 3 to 5. According to our criteria, this suggests that the answers generated by LLMs are more likely to be “answers that fit the context of the dialogue but are not a specific answer to the question” or “answers that are semantically inconsistent with the ground-truth answer but are also correct”.

From Tabel 2 we also observe that the model performance improves under the few-shot settings. However, the performance under the 5-shot setting is worse than that under the 3-shot setting. This might be due to the long length of our dialogues (as shown in Table 1, the average number of tokens in a dialogue is 501). The excessive length of model inputs under the 5-shot setting might lead to a decline in performance.

### Performance of LLMs Being Fine-tuned on CORECODE

We further evaluated LLMs after they were fine-tuned on CORECODE.

Specifically, we fine-tuned BLOOM-7.1B and ChatGLM-6B on 2K, 4K, 6K, 8K, and 10K examples respectively in the

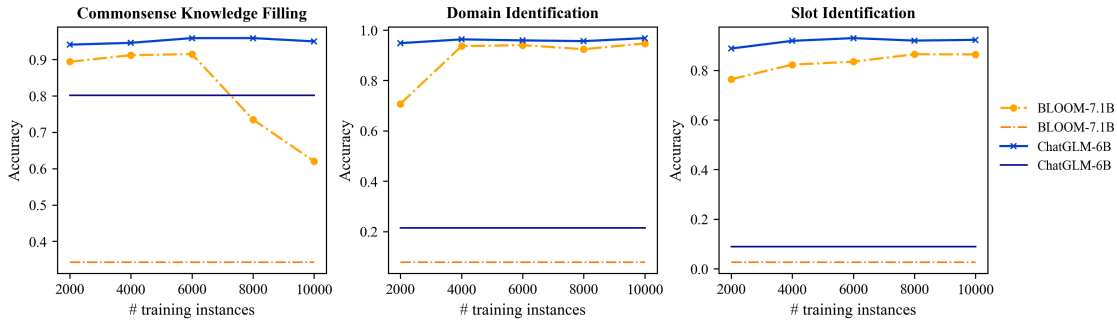


Figure 2: Results of fine-tuned LLMs on the three selection tasks. The horizontal lines show the performance of LLMs without fine-tuning.

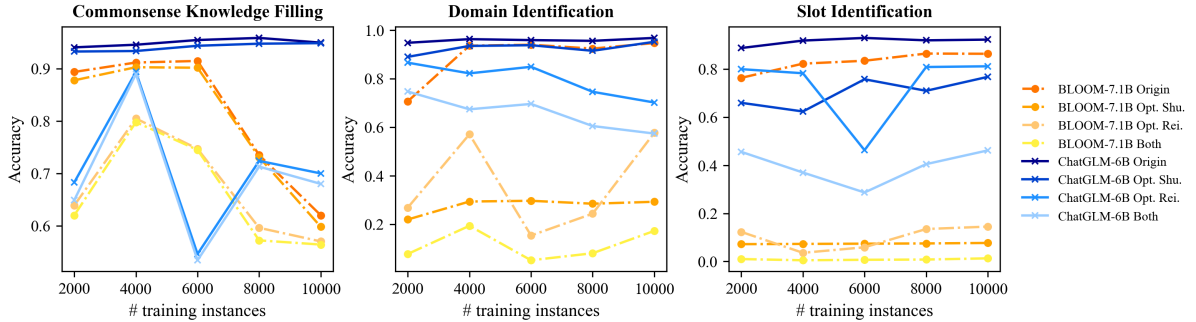


Figure 3: Results of fine-tuned LLMs on the perturbed test sets of the three selection tasks, by option re-indicating (Opt. Rei.), option shuffling (Opt. Shu.) and both.

LoRA (Hu et al. 2022) manner, and tested these fine-tuned models on another 2K data. Results on the commonsense knowledge generation and commonsense conflict phrase detection task are shown in Figure 1. Fine-tuning on different sizes of data results in large performance gains for both models. On the commonsense conflict phrase detection task, the performance in terms of F1 rises as the size of training data increases. In contrast, on the commonsense knowledge generation task, the performance rises first and then falls as the number of training instances increases, indicating that approximately 4K training instances are sufficient for this task. Training with the same amount of training data for the same duration on both tasks brings more performance gains for BLOOM-7.1B than for ChatGLM-6B. The reason could be that it is easier for BLOOM-7.1B without SFT to acquire such knowledge than ChatGLM-6B with SFT.

For the three selection tasks, as shown in Figure 2, there is a positive correlation between model performance and training data size on most tasks. Both models obtain a substantial improvement after fine-tuning.

### Robustness Analysis

Although fine-tuning on CORECODE significantly improves LLMs in commonsense reasoning, is the commonsense reasoning ability that LLMs obtained through fine-tuning robust?

To investigate this question, we conducted three robustness tests on the three selection tasks: (1) option re-

indicating, (2) option shuffling, and (3) both. For (1) option re-indicating, we change the option indicators from a, b, c to 1, 2, 3 in the process of forming the prompt. For (2) option shuffling, we shuffle the candidate options and then re-form the input prompt. For (3) both, we implement both option re-indicating and option shuffling.

Experiment results are shown in Figure 3. We find a decrease in accuracy for both models. Generally, the two LLMs are especially sensitive to option re-indicating, demonstrating larger drops. However, they are more robust to option shuffling, maintaining relatively higher accuracy. The largest performance degradation occurs when both perturbations are executed.

Perturbation causes a dramatic drop to BLOOM-7.1B. When we fine-tune LLMs on CORECODE, we use option indicators, e.g., “b”, as labels to be learned/predicted. ChatGLM-6B with SFT is better capable of understanding and following instructions than BLOOM-7.1B. It is able to align the indicator to the corresponding answer option during training and combine it with the task instruction to master the involved commonsense reasoning ability. BLOOM-7.1B, however, prefers to learn to answer by memorizing the corresponding input-label mappings. After re-indicating and shuffling answer options, BLOOM-7.1B struggles to answer correctly. For instance, on the slot identification task, our training data has a large number of examples with the label “b”. BLOOM-7.1B seems to learn such a shortcut incorrectly (i.e., mapping questions to label “b”). After shuf-

fling answer options (the correct answer indicators are now mostly not “b”), the model still outputs plenty of “b”, resulting in very low accuracies, i.e., 0.072, 0.073, 0.074, 0.075, and 0.077 for models trained on 2K, 4K, 6K, 8K, and 10K data, respectively.

## 6 Conclusion

In this paper, we have presented CORECODE, a large-scale commonsense knowledge annotated dialogue dataset with over 76K annotations, and defined 6 benchmark tasks in the form of selection, extraction and generation, to assess the capability of LLMs in learning and applying commonsense knowledge. A diverse list of Chinese LLMs have been evaluated, which achieve poor performance on all tasks, demonstrating the difficulty and utility of the proposed dataset. We have further revealed the robustness issue of LLM commonsense knowledge acquisition via fine-tuning. We hope this work could be used to track and facilitate future advances in context-sensitive LLM commonsense reasoning.

## Acknowledgments

The present research was supported by Zhejiang Lab (No. 2022KH0AB01). We would like to thank the anonymous reviewers for their insightful comments.

## References

Banerjee, S.; and Lavie, A. 2005. METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, 65–72. Ann Arbor, Michigan: Association for Computational Linguistics.

Bang, Y.; Cahyawijaya, S.; Lee, N.; Dai, W.; Su, D.; Wilie, B.; Lovénia, H.; Ji, Z.; Yu, T.; Chung, W.; et al. 2023. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *arXiv preprint arXiv:2302.04023*.

Bauer, L.; Wang, Y.; and Bansal, M. 2018. Commonsense for generative multi-hop question answering tasks.

Bhargava, P.; and Ng, V. 2022. Commonsense knowledge reasoning and generation with pre-trained language models: A survey. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, 12317–12325.

Bian, N.; Han, X.; Sun, L.; Lin, H.; Lu, Y.; and He, B. 2023. Chatgpt is a knowledgeable but inexperienced solver: An investigation of commonsense problem in large language models. *arXiv preprint arXiv:2303.16421*.

Bisk, Y.; Zellers, R.; Gao, J.; and Choi, Y. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, 7432–7439. ISBN 2374-3468. Issue: 05.

Budzianowski, P.; Wen, T.-H.; Tseng, B.-H.; Casanueva, I.; Ultes, S.; Ramadan, O.; and Gašić, M. 2018. MultiWOZ - A Large-Scale Multi-Domain Wizard-of-Oz Dataset for Task-Oriented Dialogue Modelling. In *Proceedings of the*

*2018 Conference on Empirical Methods in Natural Language Processing*, 5016–5026. Association for Computational Linguistics. EMNLP 2018.

Cambria, E.; Song, Y.; Wang, H.; and Hussain, A. 2011. Isanette: A Common and Common Sense Knowledge Base for Opinion Mining. *2011 IEEE 11th International Conference on Data Mining Workshops*, 315–322.

Clark, P.; Cowhey, I.; Etzioni, O.; Khot, T.; Sabharwal, A.; Schoenick, C.; and Tafjord, O. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge.

Cui, Y.; Yang, Z.; and Yao, X. 2023. Efficient and effective text encoding for chinese llama and alpaca. *arXiv preprint arXiv:2304.08177*.

Du, N.; Huang, Y.; Dai, A. M.; Tong, S.; Lepikhin, D.; Xu, Y.; Krikun, M.; Zhou, Y.; Yu, A. W.; and Firat, O. 2022a. Glam: Efficient scaling of language models with mixture-of-experts. In *International Conference on Machine Learning*, 5547–5569. PMLR. ISBN 2640-3498.

Du, Z.; Qian, Y.; Liu, X.; Ding, M.; Qiu, J.; Yang, Z.; and Tang, J. 2022b. GLM: General Language Model Pretraining with Autoregressive Blank Infilling. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 320–335. Dublin, Ireland: Association for Computational Linguistics.

Ghosal, D.; Hong, P.; Shen, S.; Majumder, N.; Mihalcea, R.; and Poria, S. 2021. CIDER: Commonsense Inference for Dialogue Explanation and Reasoning. In *Proceedings of the 22nd Annual Meeting of the Special Interest Group on Discourse and Dialogue*, 301–313. Singapore and Online: Association for Computational Linguistics.

Ghosal, D.; Shen, S.; Majumder, N.; Mihalcea, R.; and Poria, S. 2022. CICERO: A Dataset for Contextualized Commonsense Inference in Dialogues. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 5010–5028. Dublin, Ireland: Association for Computational Linguistics.

Guo, Z.; Jin, R.; Liu, C.; Huang, Y.; Shi, D.; Yu, L.; Liu, Y.; Li, J.; Xiong, B.; Xiong, D.; et al. 2023. Evaluating large language models: A comprehensive survey. *arXiv preprint arXiv:2310.19736*.

Hoffmann, J.; Borgeaud, S.; Mensch, A.; Buchatskaya, E.; Cai, T.; Rutherford, E.; de Las Casas, D.; Hendricks, L. A.; Welbl, J.; and Clark, A. 2022. An empirical analysis of compute-optimal large language model training. 35: 30016–30030.

Hu, E. J.; Shen, Y.; Wallis, P.; Allen-Zhu, Z.; Li, Y.; Wang, S.; Wang, L.; and Chen, W. 2022. LoRA: Low-Rank Adaptation of Large Language Models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.

Hwang, J. D.; Bhagavatula, C.; Le Bras, R.; Da, J.; Sakaguchi, K.; Bosselut, A.; and Choi, Y. 2021. (Comet-)atomic 2020: on symbolic and neural commonsense knowledge graphs. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, 6384–6392. ISBN 2374-3468. Issue: 7.



- Jiang, L.; Bosselut, A.; Bhagavatula, C.; and Choi, Y. 2021. “I’m Not Mad”: Commonsense Implications of Negation and Contradiction. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 4380–4397. Association for Computational Linguistics.
- Khot, T.; Clark, P.; Guerquin, M.; Jansen, P.; and Sabharwal, A. 2020. Qasc: A dataset for question answering via sentence composition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, 8082–8090. ISBN 2374-3468. Issue: 05.
- Levesque, H.; Davis, E.; and Morgenstern, L. 2012. The winograd schema challenge. In *Thirteenth international conference on the principles of knowledge representation and reasoning*.
- Lin, B. Y.; Chen, X.; Chen, J.; and Ren, X. 2019. KagNet: Knowledge-Aware Graph Networks for 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)*, 2829–2839. Hong Kong, China: Association for Computational Linguistics.
- Lin, B. Y.; Lee, S.; Khanna, R.; and Ren, X. 2020. Birds have four legs?! NumerSense: Probing Numerical Commonsense Knowledge of Pre-Trained Language Models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 6862–6868. Online: Association for Computational Linguistics.
- Lin, C.-Y. 2004. ROUGE: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out*, 74–81. Barcelona, Spain: Association for Computational Linguistics.
- Liu, H.; and Singh, P. 2004. ConceptNet—a practical commonsense reasoning tool-kit. *BT technology journal*, 22(4): 211–226.
- Liu, J.; Liu, A.; Lu, X.; Welleck, S.; West, P.; Le Bras, R.; Choi, Y.; and Hajishirzi, H. 2022. Generated Knowledge Prompting for Commonsense Reasoning. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 3154–3169. Association for Computational Linguistics.
- Liu, Y.; Wan, Y.; He, L.; Peng, H.; and Philip, S. Y. 2021. Kg-bart: Knowledge graph-augmented bart for generative commonsense reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, 6418–6425. ISBN 2374-3468. Issue: 7.
- Lv, S.; Guo, D.; Xu, J.; Tang, D.; Duan, N.; Gong, M.; Shou, L.; Jiang, D.; Cao, G.; and Hu, S. 2020. Graph-Based Reasoning over Heterogeneous External Knowledge for Commonsense Question Answering. 34(5): 8449–8456.
- McHugh, M. L. 2012. Interrater reliability: the kappa statistic. *Biochemia medica*, 22(3): 276–282.
- Mihaylov, T.; Clark, P.; Khot, T.; and Sabharwal, A. 2018. Can a Suit of Armor Conduct Electricity? A New Dataset for Open Book Question Answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 2381–2391. Brussels, Belgium: Association for Computational Linguistics.
- Muennighoff, N.; Wang, T.; Sutawika, L.; Roberts, A.; Biderman, S.; Scao, T. L.; Bari, M. S.; Shen, S.; Yong, Z. X.; Schoelkopf, H.; Tang, X.; Radev, D.; Aji, A. F.; Alzubair, K.; Albanie, S.; Alyafeai, Z.; Webson, A.; Raff, E.; and Raffel, C. 2022. Crosslingual Generalization through Multitask Finetuning. *CoRR*, abs/2211.01786.
- Papineni, K.; Roukos, S.; Ward, T.; and Zhu, W.-J. 2002. Bleu: a Method for Automatic Evaluation of Machine Translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, 311–318. Philadelphia, Pennsylvania, USA: Association for Computational Linguistics.
- Qin, L.; Gupta, A.; Upadhyay, S.; He, L.; Choi, Y.; and Faruqi, M. 2021. TIMEDIAL: Temporal Commonsense Reasoning in Dialog. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 7066–7076. Online: Association for Computational Linguistics.
- Quan, J.; Zhang, S.; Cao, Q.; Li, Z.; and Xiong, D. 2020. RiSAWOZ: A Large-Scale Multi-Domain Wizard-of-Oz Dataset with Rich Semantic Annotations for Task-Oriented Dialogue Modeling. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 930–940. Association for Computational Linguistics.
- Rae, J. W.; Borgeaud, S.; Cai, T.; Millican, K.; Hoffmann, J.; Song, F.; Aslanides, J.; Henderson, S.; Ring, R.; and Young, S. 2021. Scaling language models: Methods, analysis & insights from training gopher.
- Sap, M.; Le Bras, R.; Allaway, E.; Bhagavatula, C.; Lourie, N.; Rashkin, H.; Roof, B.; Smith, N. A.; and Choi, Y. 2019a. Atomic: An atlas of machine commonsense for if-then reasoning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, 3027–3035. ISBN 2374-3468. Issue: 01.
- Sap, M.; Rashkin, H.; Chen, D.; Le Bras, R.; and Choi, Y. 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)*, 4463–4473. Hong Kong, China: Association for Computational Linguistics.
- Scao, T. L.; Fan, A.; Akiki, C.; Pavlick, E.; Ilić, S.; Hesslow, D.; Castagné, R.; Luccioni, A. S.; Yvon, F.; and Gallé, M. 2022. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*.
- Storks, S.; Gao, Q.; and Chai, J. Y. 2019. Recent Advances in Natural Language Inference: A Survey of Benchmarks, Resources, and Approaches. *arXiv: Computation and Language*.
- Talmor, A.; Herzig, J.; Lourie, N.; and Berant, J. 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)*, 4149–4158. Minneapolis, Minnesota: Association for Computational Linguistics.
- Touvron, H.; Lavril, T.; Izacard, G.; Martinet, X.; Lachaux, M.-A.; Lacroix, T.; Rozière, B.; Goyal, N.; Hambro, E.; and Azhar, F. 2023. Llama: Open and efficient foundation language models.
- Vedantam, R.; Lawrence Zitnick, C.; and Parikh, D. 2015. CIDEr: Consensus-Based Image Description Evaluation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Wang, P.; Peng, N.; Ilievski, F.; Szekely, P.; and Ren, X. 2020. Connecting the Dots: A Knowledgeable Path Generator for Commonsense Question Answering. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, 4129–4140. Association for Computational Linguistics.
- Wang, X.; Li, C.; Zhao, J.; and Yu, D. 2021. Naturalconv: A chinese dialogue dataset towards multi-turn topic-driven conversation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, 14006–14014.
- Wei, J.; Wang, X.; Schuurmans, D.; Bosma, M.; Xia, F.; Chi, E.; Le, Q. V.; and Zhou, D. 2022. Chain-of-thought prompting elicits reasoning in large language models. 35: 24824–24837.
- West, P.; Bhagavatula, C.; Hessel, J.; Hwang, J.; Jiang, L.; Le Bras, R.; Lu, X.; Welleck, S.; and Choi, Y. 2022. Symbolic Knowledge Distillation: from General Language Models to Commonsense Models. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 4602–4625. Seattle, United States: Association for Computational Linguistics.
- Xu, X.; Gou, Z.; Wu, W.; Niu, Z.-Y.; Wu, H.; Wang, H.; and Wang, S. 2022. Long time no see! open-domain conversation with long-term persona memory. *arXiv preprint arXiv:2203.05797*.
- Zellers, R.; Holtzman, A.; Bisk, Y.; Farhadi, A.; and Choi, Y. 2019. HellaSwag: Can a Machine Really Finish Your Sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 4791–4800. Florence, Italy: Association for Computational Linguistics.
- Zhou, B.; Khashabi, D.; Ning, Q.; and Roth, D. 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)*, 3363–3369. Hong Kong, China: Association for Computational Linguistics.
- Zhou, P.; Gopalakrishnan, K.; Hedayatnia, B.; Kim, S.; Pujara, J.; Ren, X.; Liu, Y.; and Hakkani-Tür, D. Z. 2021. Commonsense-Focused Dialogues for Response Generation: An Empirical Study. In *SIGDIAL Conferences*.
- Zhou, X.; Zhang, Y.; Cui, L.; and Huang, D. 2020. Evaluating commonsense in pre-trained language models. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, 9733–9740.
- Zhou, Y.; Shen, T.; Geng, X.; Long, G.; and Jiang, D. 2022. ClarET: Pre-training a Correlation-Aware Context-To-Event Transformer for Event-Centric Generation and Classification. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2559–2575. Dublin, Ireland: Association for Computational Linguistics.
- Zhu, Q.; Huang, K.; Zhang, Z.; Zhu, X.; and Huang, M. 2020. CrossWOZ: A Large-Scale Chinese Cross-Domain Task-Oriented Dialogue Dataset. *Transactions of the Association for Computational Linguistics*, 8: 281–295.

## A Data Statistics

The automated filtering results on NaturalConv and DuLeMon are demonstrated in Table 8.

## B Inventory of All Domains and Slots

See Tabel 5 for the full inventory of all domains and slots.

## C Examples of All Tasks

Examples of input prompts for all tasks are shown in Figure 4 and Figure 5, where the Chinese text in the black box is the prompt we input to LLMs, and the text in the gray box is the corresponding English translation.

## D Quantitative Analysis on the Reliability of the Automatic Evaluation

As mentioned in Section 5, we employed an automatic metric, i.e., accuracy, to assess the performance of LLMs on the three selection tasks. Given the highly uncontrollable outputs of LLMs, we applied filtering measures to their outputs to prevent underestimation of LLM performance. To determine the extent to which these filtering measures improve evaluation, we conducted a quantitative analysis on the reliability of the auto evaluation, which is displayed in Table 4. We manually examined 100 outputs from each LLM and recorded the count of correct answers in various scenarios, including: (1) the raw outputs that are counted as correct (“a”, “(a)”, “(a)premise”, and “premise” are all treated as permissible formats), (2) the filtered outputs that are counted as correct, and (3) the actually correct outputs, which contain outputs that cannot be automatically evaluated as correct due to irregular formats.

The results reveal that the number of answers counted as correct significantly increases when filtering measures are applied, aligning more closely with the actual situation (i.e., the actually correct cases indicated by # correct). This indicates the effectiveness of the filtering measures, substantially enhancing the evaluation process and mitigating the underestimation of the LLMs’ performance.

## E Prompt for ChatGPT Scoring

Figure 6 demonstrates the complete prompt input provided to ChatGPT, which contains detailed scoring criteria, requiring ChatGPT to score the outputs of LLMs on the commonsense knowledge generation task.

## F Experimental Results on the EASY Set

As described before, our EASY set and HARD set differ significantly in how subjects and predicates are represented in event annotations. In the EASY set, they are indicated by speaker indicators “A” and “B”, whereas in the HARD set, they are uniformly denoted as “x” and “y”. Consequently, we conducted experiments to assess and compare the performance of the LLMs across the event-related tasks: the commonsense knowledge generation and event causal inference tasks.

The performance of the commonsense knowledge generation task is reported in Table 6. We observe that LLMs

Models	# raw correct	# filtered correct	# correct
GLM-10B	0	18	21
BLOOM-7.1B	2	29	32
ChatGLM-6B	37	85	85
BELLE-7B-2M	22	69	69
BLOOMZ-7.1B	49	76	78
MOSS-SFT-16B	40	40	40
Baichuan-7B	1	40	48
Chinese-Alpaca-Plus-7B	0	62	62
Chinese-Alpaca-Plus-13B	0	53	53

Table 4: Results of quantitative analysis of automatical evaluation.

exhibit slightly better performance on the EASY set compared to the HARD set. This variance can be attributed to the fact that LLMs are required to engage in more intricate reasoning when handling the HARD set. In this scenario, they must first deduce and identify the initiator of the event, followed by subsequent inferences. For instance, when posed with the query “*What is the occurrence time of the event ‘Person A doesn’t go out much?’*”, ChatGLM-6B responds with “*Based on the context of the event, the occurrence time of ‘Person A doesn’t go out much’ is ‘Sunday’*”. Conversely, when confronted with the question “*What is the occurrence time of the event ‘Person x doesn’t go out much?’*”, it gave an erroneous response “*It is not possible to ascertain the occurrence time of the event ‘Person x doesn’t go out much’*”. It is noteworthy that instances of successful reasoning do exist. For example, when prompted with the question “*What is the cause of the event ‘Person x is not in a good mood?’*”, ChatGLM-6B adeptly deduced the initiator of the event and subsequently reasoned to deliver the accurate answer: “*The event ‘Person x is not in a good mood’ is caused by the event ‘Person B has recently lost his job’*”.

The results on the three subtasks of the event causal inference task (namely, the event cause inference task, the subsequent event inference task, and the clipped subsequent event inference task) are demonstrated in Table 7. On these three subtasks, the results on the HARD set also exhibit some degree of degradation compared to the EASY set. Overall, the notably low performance in these three subtasks underscores the formidable nature of the event causal inference problem for existing Chinese LLMs.

	Domain	Slot	Description	Example	
Entity	Attribution	<b>Is</b>	States that S presents. Its value is an adjective.	dog	obedient
		<b>Is A</b>	S is a specific instance or a subtype of O. Its value is a noun.	dog	animal
		<b>Has</b>	S contains O. O is part of S.	bird	wings
		<b>Has Type</b>	O is a specific instance or a subtype of S. Contrary to "Is A".	movie	romance movie
		<b>Is Made Of</b> <b>Is Part Of</b> <b>Be Used To</b> <b>Capable Of</b>	S is made of O. S is part of O. O contains S. Contrary to "Has". What S can be used to do. What S can do.	noodles watermelon rind knife programmer	flour watermelon cut coding
Comparison	<b>Equivalent Entity</b>	Equivalent entity	movie	film	
	<b>Similar Entity</b>	Similar entity	pomeranian	bichon frise	
	<b>Opposite Entity</b>	Opposite entity	spear	shield	
Space	<b>At Location</b>	O is a typical location for S.	traffic police	road	
	<b>Located Near</b>	S and O are typically near each other.	school	mall	
	<b>Spatially Contains</b>	S spatially contains O.	school	canteen	
	<b>Is Spatially Contained</b>	S is spatially contained by O.	cinema	mall	
Event	Prerequisite	<b>Prerequisite</b>	Prerequisite of the event.	Mary is a good student	Mary usually studies diligently and treats people politely.
	Cause	<b>Event Cause</b>	An event that directly causes this result.	Mary fell down on her bicycle.	A car rushed out.
		<b>Emotional Cause</b>	Motivations, expectations, intentions of the person(s) leading to this result.	Mary fell down on her bicycle.	Mary wanted to be safe.
		<b>Temporal Cause</b>	The event occurring as a result of being at a certain time.	Mary fell down on her bicycle.	It's dark.
		<b>Spatial Cause</b>	The event occurring as a result of being at a certain spatial location.	Mary fell down on her bicycle.	Mary was close to a car.
	Consequence	<b>Subsequent Event</b>	Subsequent event resulting from this event.	Mary won the first prize in the competition	Mary's dad rewarded Mary with a new computer.
		<b>Subsequent Emotional Reaction</b>	Subsequent emotional reaction resulting from this event.	Mary won the first prize in the competition	Mary's happy.
		<b>Subsequent Time Change</b> <b>Subsequent Location Change</b>	Subsequent time change resulting from this event. Subsequent spatial location change resulting from this event.	Mary has been studying all day. Mary has been studying all day.	It's time for dinner. Mary's back home.
	Time	<b>Occurrence Time</b>	Time when the event occurs.	Mary fell down on her bicycle.	last month
		<b>Start Time</b>	Time when the event starts.	Mary's sleeping.	8 p.m.
<b>End Time</b>		Time when the event ends.	Mary's sleeping.	8 a.m.	
<b>Duration</b> <b>Frequency</b>		Duration of the event. Frequency of the event.	Mary's sleeping. Mary's sleeping.	eight hours every day	
Space	<b>Location</b>	Location where the event occurs.	Mary got perfect marks on the exam.	school	
Social interaction	Social norms	<b>xAttr</b>	The attribute of x.	x passed judgment on y.	impolite
		<b>xIntent</b>	The intention of x.	x paid for their meal.	Gets on with y.
		<b>xNeedTo</b>	What x need to do.	x passed judgment on y.	Apologize to y.
		<b>xReact</b>	The reaction of x.	x paid for their meal.	Feels good.
		<b>xEffect</b>	The effect on x.	x paid for their meal.	Feels further connected to y
		<b>yReact</b>	The reaction of y.	x paid for their meal.	Feels uncomfortable.
		<b>yEffect</b>	The effect on y.	x passed judgment on y.	Feels angry.

Table 5: Annotated domains and slots in our dataset. S and O denote subject and object, respectively.

Model	Commonsense Knowledge Generation (EASY)							Commonsense Knowledge Generation (HARD)						
	F1	EM	BLEU1	BLEU2	METEOR	ROUGE-L	CIDEr	F1	EM	BLEU1	BLEU2	METEOR	ROUGE-L	CIDEr
GLM-10B	0.027	0.001	0.000	0.000	0.038	0.001	0.002	0.023	0.000	0.000	0.000	0.032	0.000	0.001
BLOOM-7.1B	0.083	0.000	0.018	0.000	0.131	0.003	0.018	0.071	0.000	0.017	0.000	0.115	0.004	0.017
ChatGLM-6B	0.147	0.001	0.000	0.000	0.173	0.000	0.000	0.147	0.000	0.000	0.000	0.166	0.000	0.000
ChatGLM2-6B	0.166	0.006	0.001	0.000	0.181	0.001	0.001	0.160	0.004	0.001	0.000	0.145	0.001	0.002
BELLE-7B-0.2M	0.112	0.032	0.027	0.000	0.130	0.020	0.072	0.090	0.019	0.015	0.000	0.105	0.010	0.041
BELLE-7B-2M	0.110	0.019	0.008	0.000	0.140	0.006	0.021	0.111	0.008	0.004	0.000	0.140	0.003	0.010
BLOOMZ-1.7B	0.422	0.241	0.239	0.000	0.200	0.240	0.599	0.388	0.234	0.234	0.000	0.164	0.234	0.585
BLOOMZ-7.1B	0.488	0.297	0.294	0.000	0.251	0.296	0.741	0.438	0.284	0.282	0.000	0.199	0.283	0.707
BLOOMZ-7.1B-MT	0.474	0.306	0.305	0.000	0.221	0.306	0.764	0.435	0.300	0.300	0.000	0.184	0.300	0.750
MOSS-SFT-16B	0.208	0.074	0.058	0.000	0.178	0.110	0.208	0.199	0.071	0.066	0.000	0.147	0.049	0.174
Baichuan-7B	0.075	0.002	0.001	0.000	0.069	0.002	0.004	0.071	0.000	0.000	0.000	0.072	0.000	0.001
Chinese-Alpaca-Plus-7B	0.126	0.018	0.016	0.000	0.107	0.018	0.045	0.129	0.015	0.014	0.000	0.099	0.015	0.039
Chinese-Alpaca-Plus-13B	0.126	0.022	0.019	0.000	0.105	0.020	0.051	0.133	0.021	0.018	0.000	0.104	0.020	0.051
Average	0.197	0.078	0.076	0.000	0.148	0.079	0.194	0.184	0.074	0.073	0.000	0.129	0.071	0.183

Table 6: Experimental Results on the commonsense knowledge generation task on the EASY set and HARD set.

Model	F1	EM	BLEU1	METEOR	ROUGE-L	CIDEr	F1	EM	BLEU1	METEOR	ROUGE-L	CIDEr
	Event Cause Inference (EASY)						Event Cause Inference (HARD)					
GLM-10B	0.039	0.000	0.000	0.064	0.000	0.000	0.033	0.000	0.000	0.057	0.000	0.000
BLOOM-7.1B	0.113	0.000	0.000	0.158	0.000	0.001	0.112	0.000	0.000	0.139	0.000	0.000
ChatGLM-6B	0.142	0.000	0.000	0.222	0.000	0.000	0.141	0.000	0.000	0.214	0.000	0.000
ChatGLM2-6B	0.132	0.000	0.000	0.202	0.000	0.000	0.125	0.000	0.000	0.189	0.000	0.000
BELLE-7B-0.2M	0.101	0.000	0.000	0.149	0.000	0.000	0.101	0.001	0.000	0.140	0.000	0.000
BELLE-7B-2M	0.127	0.003	0.000	0.181	0.000	0.000	0.123	0.004	0.000	0.169	0.000	0.000
BLOOMZ-1.7B	0.244	0.020	0.017	0.145	0.017	0.042	0.214	0.013	0.011	0.124	0.011	0.028
BLOOMZ-7.1B	0.320	0.033	0.026	0.205	0.026	0.066	0.310	0.021	0.019	0.192	0.019	0.047
BLOOMZ-7.1B-MT	0.278	0.027	0.027	0.159	0.027	0.068	0.286	0.020	0.020	0.154	0.020	0.049
MOSS-SFT-16B	0.184	0.007	0.005	0.151	0.004	0.013	0.191	0.011	0.009	0.147	0.007	0.024
Baichuan-7B	0.099	0.000	0.000	0.092	0.000	0.000	0.100	0.000	0.000	0.089	0.000	0.000
Chinese-Alpaca-Plus-7B	0.113	0.000	0.000	0.114	0.000	0.001	0.109	0.000	0.000	0.104	0.000	0.000
Chinese-Alpaca-Plus-13B	0.138	0.001	0.001	0.128	0.001	0.002	0.113	0.000	0.001	0.130	0.001	0.001
Average	0.156	0.007	0.006	0.151	0.006	0.015	0.150	0.005	0.005	0.142	0.004	0.012
Model	Subsequent Event Inference (EASY)						Subsequent Event Inference (HARD)					
GLM-10B	0.045	0.000	0.000	0.068	0.000	0.000	0.040	0.000	0.000	0.053	0.000	0.000
BLOOM-7.1B	0.121	0.000	0.000	0.164	0.000	0.000	0.106	0.000	0.001	0.148	0.000	0.002
ChatGLM-6B	0.112	0.001	0.000	0.177	0.000	0.000	0.102	0.000	0.000	0.166	0.000	0.000
ChatGLM2-6B	0.142	0.002	0.000	0.186	0.000	0.000	0.118	0.000	0.000	0.159	0.000	0.000
BELLE-7B-0.2M	0.083	0.000	0.001	0.121	0.000	0.000	0.079	0.002	0.000	0.115	0.000	0.000
BELLE-7B-2M	0.117	0.004	0.001	0.150	0.001	0.002	0.108	0.005	0.000	0.136	0.000	0.000
BLOOMZ-1.7B	0.228	0.032	0.028	0.144	0.030	0.074	0.168	0.015	0.013	0.108	0.013	0.034
BLOOMZ-7.1B	0.311	0.062	0.057	0.196	0.057	0.144	0.268	0.037	0.030	0.174	0.033	0.082
BLOOMZ-7.1B-MT	0.288	0.050	0.049	0.174	0.049	0.122	0.233	0.023	0.023	0.142	0.023	0.058
MOSS-SFT-16B	0.164	0.021	0.016	0.124	0.012	0.041	0.132	0.005	0.004	0.106	0.003	0.009
Baichuan-7B	0.072	0.000	0.000	0.085	0.000	0.000	0.065	0.000	0.000	0.074	0.000	0.000
Chinese-Alpaca-Plus-7B	0.107	0.001	0.001	0.106	0.001	0.002	0.092	0.001	0.001	0.093	0.001	0.002
Chinese-Alpaca-Plus-13B	0.123	0.002	0.002	0.109	0.002	0.004	0.109	0.000	0.000	0.099	0.000	0.001
Average	0.147	0.013	0.012	0.139	0.012	0.030	0.125	0.007	0.005	0.121	0.006	0.014
Model	Clipped Subsequent Event Inference (EASY)						Clipped Subsequent Event Inference (HARD)					
GLM-10B	0.048	0.000	0.000	0.037	0.000	0.000	0.043	0.000	0.000	0.029	0.000	0.000
BLOOM-7.1B	0.089	0.000	0.000	0.115	0.000	0.000	0.084	0.000	0.000	0.108	0.000	0.000
ChatGLM-6B	0.069	0.000	0.000	0.111	0.000	0.000	0.067	0.000	0.000	0.108	0.000	0.000
ChatGLM2-6B	0.083	0.000	0.000	0.113	0.000	0.000	0.073	0.000	0.000	0.104	0.000	0.000
BELLE-7B-0.2M	0.070	0.000	0.000	0.091	0.000	0.000	0.065	0.000	0.000	0.084	0.000	0.000
BELLE-7B-2M	0.045	0.000	0.000	0.054	0.000	0.000	0.056	0.000	0.000	0.062	0.000	0.000
BLOOMZ-1.7B	0.118	0.000	0.000	0.070	0.000	0.000	0.096	0.001	0.001	0.057	0.001	0.003
BLOOMZ-7.1B	0.164	0.008	0.007	0.090	0.008	0.020	0.145	0.004	0.004	0.084	0.004	0.010
BLOOMZ-7.1B-MT	0.133	0.004	0.004	0.069	0.004	0.011	0.116	0.002	0.002	0.063	0.002	0.006
MOSS-SFT-16B	0.096	0.002	0.001	0.077	0.001	0.004	0.092	0.001	0.001	0.072	0.001	0.003
Baichuan-7B	0.059	0.000	0.000	0.074	0.000	0.000	0.052	0.000	0.000	0.063	0.000	0.000
Chinese-Alpaca-Plus-7B	0.069	0.000	0.000	0.064	0.000	0.000	0.067	0.000	0.000	0.057	0.000	0.000
Chinese-Alpaca-Plus-13B	0.055	0.000	0.000	0.048	0.000	0.000	0.061	0.000	0.000	0.051	0.000	0.000
Average	0.084	0.001	0.001	0.078	0.001	0.003	0.078	0.001	0.001	0.072	0.001	0.002

Table 7: Experimental Results on the three subtasks of the event causal inference task on the EASY set and HARD set.

	# dialogues	# 1 triple matches	# 2 triple matches	# 3 triple matches	# dialogues in CORECODE
NaturalConv	19919	17015	12938	9145	9145
DuLeMon	27501	24362	21926	18346	10555
Total	47420	41377	34864	27491	19700

Table 8: Statistics after screening on NaturalConv and DuLeMon. # dialogues: total number of dialogues in the dataset. # 1 triple matches: number of dialogues where more than one commonsense triple matches occurred.

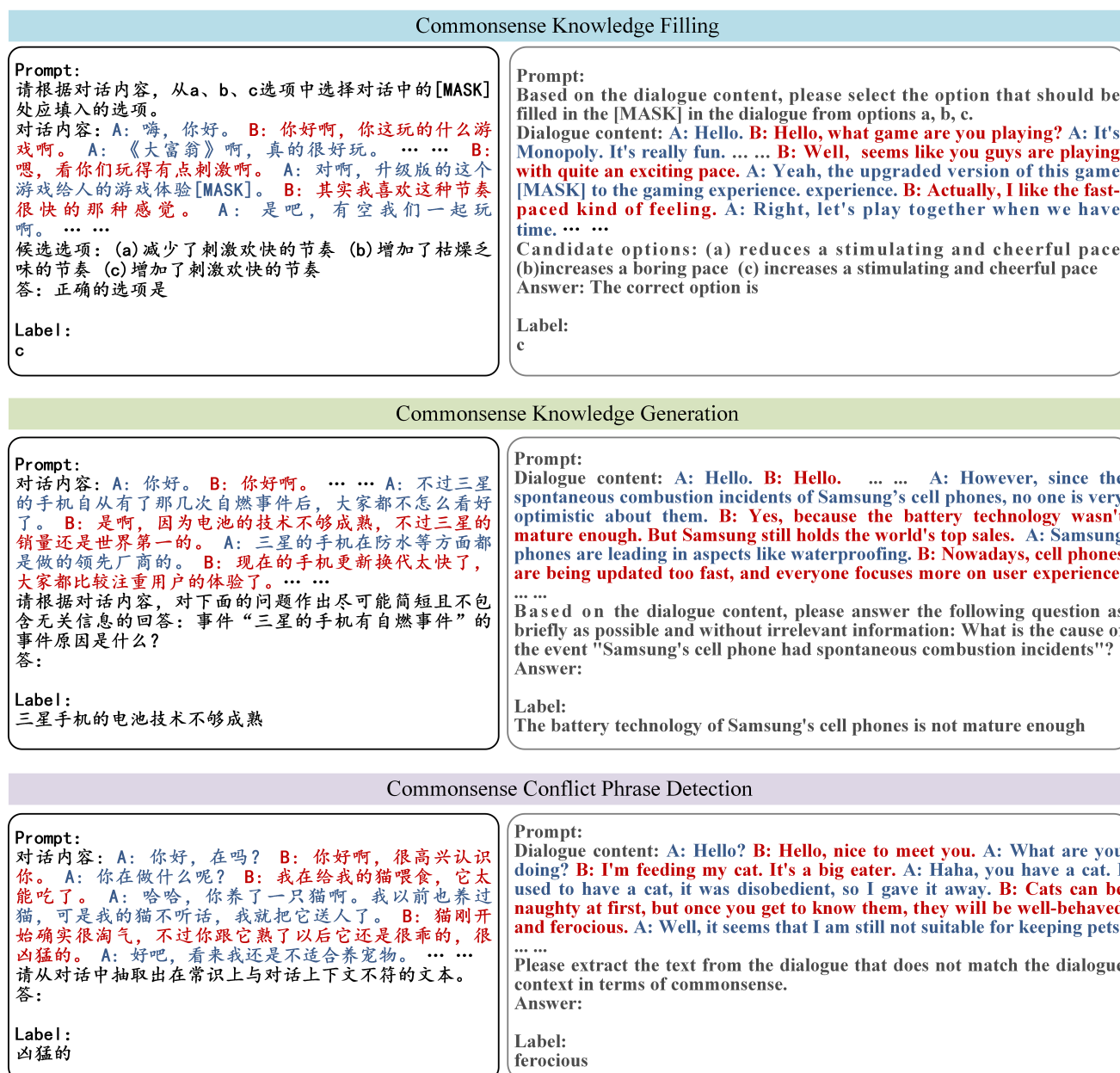


Figure 4: Examples of the commonsense knowledge filling, commonsense knowledge generation and commonsense conflict phrase detection tasks. Due to the extensive length of the dialogues, we only display a relevant section of each dialogue that pertains to the reasoning required for the task, rather than presenting the entire dialogue.

### Domain Identification

**Prompt:**

请根据对话内容，从a、b、c等候选领域中选择下面两个短语之间的关系所属的领域。

短语1: 郝伟对队员很严格 短语2: 郝伟希望以最好的状态迎接比赛

对话内容: ... .. A: 没错，对于即将到来的比赛郝伟都亲自上阵示范了。 B: 看来他对队员的要求还是很严格的呀。 A: 那肯定的，要以最好的状态迎接比赛。

B: 加油啊！希望看到国奥队出色的成绩。 A: 嗯嗯。

B: 哎呀，我得先走了。 A: 好，再见。 B: 拜拜！

候选领域: (a)前提 (b)原因 (c)后续 (d)时间 (e)空间

答: 正确的领域是

Label:

b

**Prompt:**

Based on the dialogue content, please choose the domain to which the relationship between the following two phrases belongs from the candidates a, b, c, and so on.

Phrase 1: Hao Wei is very strict with his team members Phrase 2: Hao Wei wants to be in the best condition for the game

Dialogue content: ... .. A: That's right, Hao Wei personally demonstrated for the upcoming match. B: He seems to be very strict with his team members. A: Definitely, he wants to be in the best condition for the game.

B: Keep it up! Hoping to see an excellent performance from the national team. A: Me too. B: Oh well, I have to go now. A: Alright, goodbye. B: Bye!

Candidate domains: (a) Prerequisite (b) Cause (c) Consequence (d) Time (e) Space

Answer: The correct domain is

Label:

b

### Slot Identification

**Prompt:**

请根据对话内容，从a、b、c等选项中选择下面两个短语之间的关系。

短语1: 军队 短语2: 食堂

对话内容: A: 中午好呀！在吗？ B: 你好，在呢。 A: 吃饭了吗？ B: 还没有开饭呢，你一问，还真饿了。 A: 等着别人做好了叫你吗？ B: 食堂还没有到开饭的点儿呢。 A: 是学生吗？ B: 不是，我在军队呢，中午不想回去吃了，就在食堂里面吃点儿。 A: 军队的伙食一定很好。 B: 还好吧，你是做什么的？ ... ..

候选选项: (a)是 (b)是一个 (c)有 (d)包括...类型 (e)由...制成 (f)是...的一部分 (g)用途 (h)能力 (i)相同物 (j)相似物 (k)相反物 (l)在某个位置 (m)空间上相邻 (n)空间上包含 (o)空间上包含于

答: 正确的选项是

Label:

n

**Prompt:**

Based on the dialogue content, please choose the relationship between the following two phrases from the options a, b, c, and so on.

Phrase 1: Army Phrase 2: Canteen

Dialogue content: A: Good noon! Are you there? B: Hello, yes, I'm here. A: Have you eaten yet? B: Not yet, but now that you've asked, I suddenly feel hungry. A: Are you waiting for someone to call you when the food is ready? B: It's not mealtime at the canteen yet. A: Are you a student? B: No, I'm in the army. I don't want to go back for lunch, so I'll eat in the canteen. A: The food in the army must be very good. B: It's okay. What do you do? ... ..

Candidate options: (a) Is (b) Is A (c) Has (d) Has Type (e) Is Made Of (f) Is Part of (g) Be Used To (h) Capable Of (i) Equivalent Entity (j) Similar Entity (k) Opposite Entity (l) At Location (m) Located Near (n) Spatially Contains (o) Is Spatially Contained

Answer: The correct option is

Label:

n

Figure 5: Examples of the domain identification and slot identification tasks. The event causal inference task can be regarded as a special case of the commonsense knowledge generation task, hence no illustrative examples are shown here.

[对话内容开始]

A: 哈喽, 你好  
 B: 你好呀  
 A: 你喜欢看足球比赛吗?  
 B: 当然喜欢啊, 我只要一有空闲时间我就会看足球比赛  
 A: 最近的荷甲你看了吗, 阿贾克斯对战海伦芬这场比赛真的好精彩  
 B: 哈哈, 我昨天才看完, 我现在脑袋都还回想这场比赛呢, 阿贾克斯对战海伦芬的比赛好像是阿贾克斯大获全胜呢  
 A: 嗯嗯, 是的, 阿贾克斯简直可以说是秒杀海伦芬了, 真的太厉害了  
 B: 对啊, 在这场比赛中阿贾克斯可以说是神勇无敌了, 很多球感觉都进得莫名其妙的  
 A: 哈哈, 就是, 我记得在比赛开始后的十几分钟的时候本来双方都在僵局中, 突然阿贾克斯一个意外的进球打破了这个僵局  
 B: 就是, 我看到这儿的时候我都懵了, 这个球是怎么进的都不知道, 也不知道是谁踢的呢  
 A: 哈哈, 别说你知道了, 我敢说连阿贾克斯的球员都不知道这个球是怎么进的呢  
 B: 嘻嘻, 是的呢, 不过呢他们赢了还是不错的, 虽然赢得有点莫名其妙  
 A: 哈哈, 人家可是光明正大踢的球, 也没有作弊, 赢了这场比赛也是正常的嘛  
 B: 嗯嗯, 也是。对了, 你喜欢哪个足球明星?  
 A: 我很喜欢贝克汉姆, 我超级喜欢他踢球时的状态, 真的很帅  
 B: 哈哈, 看来我们很有缘分了, 我也很喜欢贝克汉姆, 我觉得他的球技是真的很棒  
 A: 对啊, 尤其是他赢了比赛之后无比激动, 全场奔跑的样子真的就像刚放出的野兽一样  
 B: 哈哈, 你把他比作成野兽真的厉害了, 很形象  
 A: 嘻嘻, 他真的很像嘛, 而且他在球场上的状态就是像野兽盯着自己的猎物一样, 真的很霸气  
 B: 我要休息了, 再见

[对话内容结束]

[问题开始]

请根据对话内容, 直接回答下面问题的答案, 不要重述问题或解释原因: “贝克汉姆”是一个或一种什么?

[问题结束]

[打分标准开始]

0: 生成答案的语句不通顺或意思不流畅;  
 1: 尽管生成答案的语句通顺, 意思流畅, 但不是问题的正确答案, 且与对话情境相差较远;  
 2: 尽管生成答案的语句通顺, 意思流畅, 符合对话情境, 但不是问题的正确答案;  
 3: 尽管生成答案的语句通顺, 意思流畅, 符合对话情境, 但语义比较宽泛, 没有具体地回答问题;  
 4: 生成答案的语句通顺, 意思流畅, 尽管与标准答案的语义不一致, 但也是一个正确的答案;  
 5: 生成答案的语句通顺, 意思流畅, 且与标准答案的语义一致。

[打分标准结束]

[标准答案开始]

足球明星

[标准答案结束]

[候选答案开始]

“贝克汉姆”是一个足球明星。

[候选答案结束]

请你参照以上标准答案, 为以上问题的候选答案进行评分, 请首先输出一行, 该行仅包含一个值, 表示得分。在接下来的行中, 请对你的评判进行解释, 避免任何潜在的偏见和后果。

[Dialogue Content Begins]

A: Hello!  
 B: Hello!  
 A: Do you like watching football games?  
 B: Of course I do. I watch football games whenever I have free time.  
 A: Have you watched the recent Eredivisie? The match between Ajax and Heerenveen was really exciting.  
 B: Haha, I just finished watching it yesterday, and the game is still replaying in my mind. The game between Ajax and Heerenveen seemed to be a complete victory for Ajax.  
 A: Yeah, indeed. Ajax completely dominated Heerenveen. They were so impressive.  
 B: Absolutely. Ajax seemed invincible in this game. Many goals were scored inexplicably.  
 A: Haha, yes. I remember in the first fifteen minutes or so, both teams were at a deadlock, and suddenly, Ajax scored unexpectedly, breaking that deadlock.  
 B: Exactly! I was stunned when I saw that. I had no idea how that goal went in or who scored it.  
 A: Haha, I think it wasn't just you who was confused. I dare say that even the Ajax players didn't know how this goal was scored.  
 B: Maybe that's true. But still, their win was impressive, even if it was a bit unexpected.  
 A: They played fair and square, no cheating involved. It is normal for them to win this game.  
 B: Yeah, that's true. By the way, which football star do you like?  
 A: I really like David Beckham. I admire how he looks when he plays. He looks so handsome on the field.  
 B: Haha, seems like we have something in common. I'm a big fan of Beckham too. His football skills are truly fantastic.  
 A: I especially like how he was so excited after winning the game and ran all over the field, really like a released beast.  
 B: Haha, likening him to a beast is quite something, very vivid.  
 A: He really does! His demeanor on the field is like a beast eyeing its prey, so commanding.  
 B: I need to take a break now. Goodbye.

[Dialogue Content Ends]

[Question Begins]

Based on the dialogue content, please answer the following question directly without restating the question or explaining the reason: What is "Beckham"??

[Question ends]

[Scoring Criteria begins]

0: The generated answer's sentences are not coherent or lack fluency.  
 1: Although the generated answer has a coherent statement and fluent meaning, it is not the correct answer to the question and is far from the dialogue context;  
 2: Although the generated answer has a coherent statement and fluent meaning, and fit to the dialogue context, it is not the correct answer to the question;  
 3: Although the generated answer has a coherent statement and fluent meaning, and fit to the dialogue context, it is semantically broad and does not specifically answer the question;  
 4: The generated answer has a coherent statement and fluent meaning. Although it is semantically inconsistent with the ground-truth answer, it is also a correct answer;  
 5: The generated answer has a coherent statement and fluent meaning, and is semantically consistent with the ground-truth answer.

[Scoring Criteria ends]

[Ground-truth Answer Begins]

Football star

[Ground-truth Answer Ends]

[Candidate Answer Begins]

"David Beckham" is a football star.

[Candidate Answer Ends]

Please rate the candidate answer for the above question based on the provided ground-truth answer. First, output a line containing only one value, indicating the score. In the next line, explain your judgment avoiding any potential bias and consequences.

Figure 6: Prompt for ChatGPT to score the outputs of LLMs on the commonsense knowledge generation task. The Chinese text in the black box (left) is the prompt we input to the ChatGPT, while the text in the gray box (right) is the corresponding English translation.