

# CKBP v2: Better Annotation and Reasoning for Commonsense Knowledge Base Population

Anonymous ACL submission

## Abstract

Commonsense Knowledge Bases (CSKB) Population, which aims at automatically expanding knowledge in CSKBs with external resources, is an important yet hard task in NLP. Fang et al. (2021a) proposed a CSKB Population (CKBP) framework with an evaluation set CKBP v1. However, CKBP v1 relies on crowdsourced annotations that suffer from a considerable number of mislabeled answers, and the evaluation set lacks alignment with the external knowledge source due to random sampling. In this paper, we introduce CKBP v2, a new high-quality CSKB Population evaluation set that addresses the two aforementioned issues by employing domain experts as annotators and incorporating diversified adversarial samples to make the evaluation data more representative. We show that CKBP v2 serves as a challenging and representative evaluation dataset for the CSKB Population task, while its development set aids in selecting a population model that leads to improved knowledge acquisition for downstream commonsense reasoning. A better population model can also help acquire more informative commonsense knowledge as additional supervision signals for both generative commonsense inference and zero-shot commonsense question answering. Specifically, the question-answering model based on DeBERTa-v3-large (He et al., 2023b) even outperforms powerful large language models in a zero-shot setting, including ChatGPT and GPT-3.5.

## 1 Introduction

Recently introduced LLMs have shown a remarkable performance on many reasoning benchmarks (Hoffmann et al., 2022; Chowdhery et al., 2022; Bang et al., 2023; Chan et al., 2023), yet there still exists a need to ensure the alignment between the generation of LLMs with external knowledge at the inference time to avoid hallucination and for safer use (Kim et al., 2022a; He et al., 2023a; Peng et al., 2023). The source of

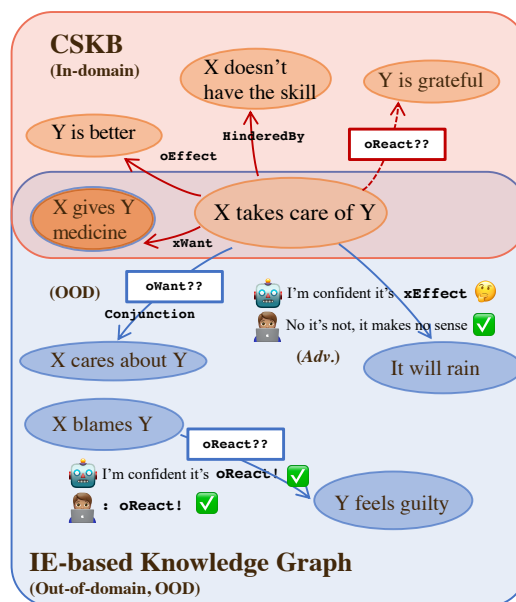


Figure 1: An example of CSKB Population. The coral part indicates the conventional case of CSKB Completion, and the blue part is the population on external knowledge graphs. We include an adversarially constructed sample set in our CKBP v2 by re-annotating the confident predictions by language models.

external knowledge, which can be commonsense, factual, or domain knowledge, should be selected and processed carefully depending on the purpose of generation. However, existing (high-quality) human-annotated knowledge bases are usually far from complete to serve as the source of external knowledge for LLMs.

Regarding commonsense knowledge bases, to extend limited human annotations, CSKB Population (Fang et al., 2021a) stands as a way to acquire missing knowledge, thereby enriching and expanding the existing CSKBs. Unlike CSKB Completion (Li et al., 2016; Saito et al., 2018; Malaviya et al., 2020), which adopts a close-world assumption and only deals with entities and events within CSKBs, the Population task deals with both existing and unseen entities and events, thus requiring a

more generalized reasoning ability.

Several works have been conducted on CSKB Population. Fang et al. (2021a) studied a framework that links four CSKBs, ConceptNet (Speer et al., 2017), ATOMIC (Sap et al., 2019a), ATOMIC<sub>20</sub> (Hwang et al., 2021), and GLUCOSE (Mostafazadeh et al., 2020), to a large-scale discourse knowledge base, ASER (Zhang et al., 2020, 2022). The resulting knowledge base not only served as the unified source of commonsense knowledge but also was used as the training set to train population models in order to identify missing commonsense knowledge. To evaluate models, the authors created an evaluation set (denoted as CKBP v1), in which they applied fine-grained rules to select candidate commonsense knowledge from ASER and enlisted human annotators to manually annotate these candidates.

However, there are two major limitations in CKBP v1. First, the quality of CKBP v1 is limited. CKBP v1 instances are randomly sampled from the whole population space, resulting in a low recall of plausible commonsense knowledge due to the noise in candidate discourse knowledge. Moreover, as pointed out by Davis (2023), current crowd-sourced commonsense benchmarks often contain a substantial fraction of incorrect answers, we also find it true for CKBP v1 after manual inspection. For example, annotators frequently make mistakes on some subtle relations such as *xIntent*, which should describe an *intention* instead of a *consequence*. Second, it’s unclear how to leverage populated or expanded commonsense knowledge in CKBP to further improve downstream commonsense reasoning. All previous investigations into CKBP stay within the population task itself without generalizing to actual downstream applications.

Therefore, to address the two limitations, this work presents a more high-quality and adversarially constructed evaluation set by expert annotation, and a comprehensive pipeline for conducting a series of downstream experiments. The aim is to leverage the new CKBP benchmark effectively and facilitate improved utilization for downstream commonsense reasoning tasks.

Leveraging the existing framework, we build CKBP v2 by randomly sampling 2.5k instances from CKBP v1 and adding 2.5k adversarial instances, leading to a total of 5k instances as an evaluation set. These instances are then annotated by experts with substantial expertise in machine commonsense. Then, we present both in-

trinsic and extrinsic experiments based on CKBP v2. We study the performance of both supervised and semi-supervised task-specific models, together with powerful off-the-shelf language models, such as ChatGPT (OpenAI, 2022) and Vera (Liu et al., 2023), and show that the CKBP v2 evaluation set is still challenging even for advanced language models. Moreover, by employing a CSKB Population model that demonstrates satisfactory performance on CKBP v2, we can enrich existing CSKBs with diverse and novel knowledge that significantly benefits downstream reasoning. We present methodologies and experiments on generative commonsense inference (Bosselut et al., 2019) and zero-shot commonsense question answering (Ma et al., 2021), and show that the acquired commonsense knowledge can be valuable augmented data on the original CSKB and lead to improved downstream performance. In particular, CKBP v2-preferred population model exhibits better alignment than CKBP v1 with advancements in generative commonsense inference.

In summary, our contributions are three-fold: First, We introduce a new evaluation benchmark CKBP v2 for the CSKB Population task, which addresses the quality issues of its predecessor CKBP v1. Second, We launch a pioneer study to use populated commonsense knowledge as additional supervision signals to help downstream commonsense reasoning. Third, We conduct extensive experiments and evaluations with different models on both CKBP v2 itself as well as downstream generative commonsense inference and zero-shot question answering. The results show that CKBP v2 is still a hard task for language models, and the acquired populated knowledge can improve language models’ (zero-shot) commonsense reasoning ability on two downstream tasks across six datasets.

## 2 Related Work

In this section, we discuss 1) CSKBs and their role in the era of LLMs and 2) methods and benchmarks for completing and populating knowledge bases in general.

**Commonsense Knowledge Bases.** There are many commonsense knowledge bases<sup>1</sup> introduced in the past few years, such as ATOMIC2020 (Hwang et al., 2021), ComFact (Gao et al., 2022), CICERO (Ghosal et al.,

<sup>1</sup>Here, despite the subtle differences between datasets and knowledge bases, we refer to both as knowledge bases

2022), PIQA (Bisk et al., 2020a), Numersense (Lin et al., 2020). Unlike the decades-old knowledge base ConceptNet (Liu and Singh, 2004) that only focuses on taxonomic commonsense, these knowledge bases study a broad range of commonsense, including human-event-centric, contextualized, physical, numerical commonsense.

Along with pure-symbolic CSKBs whose knowledge is obtained from corpora and stored in textual format, there is a stream of research that works on developing neural(-symbolic) CSKBs, which are either knowledge models such as COMET (Bosselut et al., 2019) or symbolic CSKBs built by prompting knowledge from language models, such as ATOMIC<sup>10X</sup> (West et al., 2022a), SODA (Kim et al., 2022a). Although the approach seems highly scalable and seems promising to build more and larger CSKBs, knowledge from neural(-symbolic) CSKBs remains unreliable (Kim et al., 2022a; He et al., 2023a; Peng et al., 2023) thus often needs to have a robust critic model to filter for good/correct knowledge.

### Completing and Populating Knowledge Bases.

Regarding conventional knowledge bases like Wordnet (Miller, 1995) and Freebases (Bollacker et al., 2008), tasks involving completion and population have been well-studied as transductive and inductive link prediction problems in the field of graph neural network (Bordes et al., 2013; Yang et al., 2015; Sun et al., 2019; Shang et al., 2019; Fang et al., 2021b). Methods powered by pre-trained language models have also been studied in these tasks thanks to the models’ representation power (Yao et al., 2019). In that setting, knowledge instances of the knowledge bases are serialized to a text sequence, which serves as input to LMs such as BERT or RoBERTa.

Specific to CSKB Population task on CKBP v1, Fang et al. (2021a) proposed KGBertSAGE, a combination of KG-BERT (Yao et al., 2019) and GraphSAGE (Hamilton et al., 2017). The model showed higher performance over baselines yet still suffered from the out-of-domain problem. The follow-up work PseudoReasoner (Fang et al., 2022) employs the pseudo-labeling technique to solve that problem. Despite the significant gain in performance, PseudoReasoner is still far from human performance, suggesting that CKBP remains a challenging task in commonsense reasoning.

## 3 Dataset Construction

In this section, we introduce the task definition, the preparation of the candidate evaluation set, annotation guidelines, and data analysis.

### 3.1 Task Definition

The task of CKBP (Fang et al., 2021a) is defined as follows. Given  $G^C = \{(h, r, t) | h \in H, r \in R, t \in T\}$  (where  $H, R, T$  is the set of head events, relations, and tail events), the graph-like knowledge base formed by aligning a union of commonsense knowledge bases  $C$  and a much larger discourse knowledge graph  $G$  into the same format; the goal of CSKB population task is to learn a scoring function that gives a candidate knowledge triple  $(h, r, t)$  higher score if the triple is plausible commonsense. The training process is formulated as triple classification, with ground-truth positive triples from the CSKB  $C$  and negative triples randomly sampled from  $G^C - C$  with an equal amount. The model is then evaluated on a human-annotated evaluation set  $E$ . Here, CKBP v2 serves as the evaluation set.

### 3.2 Dataset Preparation

We randomly sampled 2.5k instances from CKBP v1 and 2.5k adversarial instances to form CKBP v2. Instances from CKBP v1 are sampled so that the ratio of the number of triples between relations remains unchanged. Meanwhile, the adversarial instances are ones from the candidate knowledge base ASER that the finetuned baseline KG-BERT (Yao et al., 2019) model confidently believes they are plausible, i.e., receives plausibility score  $\geq 0.9$ . To ensure the diversity of adversarial instances and hence the evaluation set, we adopt an additional diversity filter using self-BLEU following West et al. (2022a). The triples annotated as negative are considered *hard negatives* as they are what a standard CSKB Population model would favor. Note that we only consider instances of 15 relations other than general Want/React/Effect, because most of the triples on the three relations are broken sentences in CKBP v1. We also remove samples of these relations in the training set.

### 3.3 Annotation Process

**Setup** We recruited four human experts for the annotation work. The experts are graduate NLP researchers with at least one year of experience working on CSKBs. We randomly divide 5k samples into 4 parts, then for  $i$  from 0 to 3, assign the  $i^{th}$

	# Triples	% Plau.	% Unseen
<b>split</b>			
Dev	958	20.46	56.79
Test	4,048	22.06	60.43
<b>instance type</b>			
In-Domain	845	34.56	43.79
Out-of-Domain	1,653	11.92	63.37
<i>Adv.</i>	2,508	23.92	61.12
<b>relation</b>			
xWant	611	22.75	54.01
oWant	239	25.94	58.18
xEffect	603	29.68	55.23
oEffect	172	21.51	58.91
xReact	533	20.64	51.18
oReact	183	13.66	50.70
xAttr	605	23.47	52.91
xIntent	239	16.32	58.40
xNeed	378	25.66	55.37
Causes	236	21.61	55.41
xReason	5	40.0	30.0
isBefore	157	28.03	54.80
isAfter	182	24.73	55.40
HinderedBy	777	12.1	63.17
HasSubEvent	86	26.74	61.04

Table 1: Statistics of CKBP v2. # Triples, % Plausible, and % Unseen, respectively, indicate the number of triples in the subset, the proportion of plausible triples after label finalization, and the proportion of nodes that do not appear in the training set.

and  $(i + 1 \bmod 4)^{th}$  parts to the  $i^{th}$  expert. In this way, two different annotators annotate each triple, and we can fully compare the pairwise agreement between all four annotators. Experts are provided with knowledge triples in the format of  $(h, r, t)$ , referencing the definition and examples of all relations in Hwang et al. (2021). We ask annotators to judge the plausibility of triples in a three-point Likert scale with corresponding scores: Always/Often (1), Sometimes (0.5), Rarely/Never/Ambiguous/Invalid (0). The final label of an instance is determined as *plausible* if and only if it receives at least one score of 1 and the other score is at least 0.5. For remaining cases, the final label is *implausible*. After finalizing the annotation, we split the evaluation set into development and test sets with a ratio of 1:4 with the preservation of distribution w.r.t labels, relations, and instance types. To estimate human performance, we treat expert annotations as two sets of predictions and compare them to the final labels.

Similar to CKBP v1, we categorize the evaluation set into three groups based on their origin, which are 1) ID: in-domain, whose head and tail events are all from CSKBs, 2) OOD: out-of-

domain, which has at least one event outside of CSKBs (equivalent to “CSKB head + ASER tail” and “ASER Edges” in CKBP v1), and 3) *Adv.*: adversarial examples newly introduced in CKBP v2.

**Quality Control** Although annotators are experts with a clear understanding of the CSKB Population, we acknowledge the ambiguity of CSKB relations and the difficulty in discriminating between them. To control the quality, we provide guidance as a list of scoring criteria. We also carried out a dry run, which asked them to annotate 60 instances covering all relations in order to establish a unified understanding of the problem among participants.

After that, we carry out the main round, where the annotators perform their jobs individually and independently. Throughout the process, we regularly conduct random checks on the samples and engage in discussions with annotators to address any disagreements. We then use the insights gained from these discussions to update and refine our guidance iteratively. After the individual annotation, we facilitated a conflict resolution session to address instances with contrasting scores of 1 and 0. After resolving conflicts, we have the average inter-annotator agreement score IAA as 90.55%.

### 3.4 Data Analysis

The overall statistics of CKBP v2 are shown in Table 1. It can be easily observed that the new evaluation set has data imbalance issues. However, we do not down-sample the evaluation set to achieve the data balance since the imbalance better reflects the true distribution of plausible and implausible commonsense knowledge in ASER. Given this imbalance, we notice that the AUC scores of examined population models will naturally be high. Also, in the real application of population models, we focus on the precision and recall of the detection for plausible commonsense instances. Thus, in Section 4, along with AUC, we also report the binary F1 scores for each experimented model.

## 4 Intrinsic Evaluation

### 4.1 Setup

We examine several models which were previously evaluated on CKBP v1, including zero-shot GPT models (Radford et al., 2019), supervised-learning baselines KG-BERT (Yao et al., 2019) and COMET (Bosselut et al., 2019), and semi-supervised-learning models PseudoReasoner (Fang



Category	Model	AUC				F1			
		all	ID	OOD	<i>Adv.</i>	all	ID	OOD	<i>Adv.</i>
Zero-shot	GPT2-large	56.47	56.60	58.31	54.22	35.37	47.40	24.06	36.84
	GPT2-XL	56.79	54.47	56.70	54.63	35.22	47.62	23.49	36.65
	GPT3 <i>text-davinci-003</i>	61.63	65.93	59.17	59.98	39.44	51.09	28.57	38.20
	ChatGPT <i>gpt-3.5-turbo</i>	65.77	70.37	62.56	62.27	45.93	62.59	44.79	26.86
Supervised Learning	KG-BERT (BERT-base)	71.33	84.60	64.47	62.9	45.03	69.27	26.53	41.97
	KG-BERT (RoBERTa-L)	<u>73.70</u>	<u>85.53</u>	67.70	65.60	<u>46.70</u>	<u>69.73</u>	30.73	<u>43.27</u>
	COMET (GPT2-L)	70.00	79.02	66.43	62.62	45.55	61.90	<u>32.14</u>	42.15
	COMET (GPT2-XL)	70.32	79.66	66.53	63.22	45.32	63.34	31.18	40.83
	Vera (T5-xxlarge)	72.45	78.84	<u>68.40</u>	<b>68.16</b>	<b>52.13</b>	<b>71.73</b>	<b>36.74</b>	<b>50.02</b>
Semi-Supervised	PseudoReasoner BERT-base	71.93	84.23	66.67	63.43	45.47	68.67	30.17	41.77
	PseudoReasoner RoBERTa-L	<b>74.33</b>	<b>85.57</b>	<b>69.33</b>	<u>66.37</u>	46.63	69.70	30.87	43.13
Human		94.1	94.9	91.4	94.5	91.5	94.3	86.9	91.5

Table 2: Main experimental results on CKBP v2. Both AUC and F1 are used as evaluation metrics. The “all” column indicates the overall performance, and ID, OOD, *Adv.* indicate the performance of the In-domain, Out-of-domain, and Adversarial subset. The best results are **boldfaced**, and the second-best ones are underlined.

et al., 2022) with two backbone encoders, BERT-base-uncased (Devlin et al., 2019) and RoBERTa-large (Liu et al., 2019). We use Huggingface<sup>2</sup> Transformers (Wolf et al., 2020) to build our code base. For discriminative models, we set the learning rate as 1e-5, batch size 64/32 for base/large variants, respectively, and the number of training epochs as 1. For generative models (COMET), we use learning rate 1e-5 and batch size 32 to train in 3 epochs. Negative perplexity scores are used as the final prediction scores. For PseudoReasoner, we adopt the best settings in Fang et al. (2022), where we first finetune the KG-BERT model on pseudo-labeling data for one epoch, then from the best checkpoint, we resume the finetuning process on the original training data. Note that the training data and unlabeled data are taken from Fang et al. (2022). We run each baseline three times with different random seeds, then average the result and report in Table 2. For GPT3 (Brown et al., 2020a) and ChatGPT experiments, we use simple prompts asking them to decide whether an assertion is plausible or not.

## 4.2 Result and Analysis

The results are shown in Table 2. We provide the AUC score and F1 score of all the baselines on the test set in terms of overall performance (all), performance on the subset of ID, OOD, and *Adv.* samples. When calculating F1, for discriminative models, we set the decision threshold as 0.5 (as default), while for generative models, as perplexity

<sup>2</sup><https://huggingface.co/>

serves as the final prediction score, we tune the threshold to obtain the highest F1 score on the development set for each run.

In the zero-shot setting, the scores increase by the version of GPT. GPT3 *text-davinci-003* gives a significant improvement over GPT2 models, and ChatGPT surpasses its sibling *text-davinci-003* with a similar margin of improvement. Nonetheless, despite the performance improvement from ChatGPT, there is still a clear gap between the zero-shot and (semi-)supervised settings.

In terms of supervised and semi-supervised learning, we observe different scenarios between KG-BERT’s performance and COMET’s performance, comparing to the result on CKBP v1 reported in Fang et al. (2022). Here, on CKBP v2, KG-BERT outperforms COMET with a significant gap of 3 AUC overall and also outperforms in all subsets of the test set. This shows the importance of including negative (implausible) examples in the training for discriminating commonsense. This also explains why there is no significant improvement of PseudoReasoner over the baseline KG-BERT on this new evaluation set.

## 4.3 Artifacts Analysis

There is an uprising acknowledgment of “artifacts” (Gururangan et al., 2018; Poliak et al., 2018; Gardner et al., 2021) in a dataset, in other words, spurious correlations or confounding factors between the surface properties of textual instances and their labels, that may incidentally appear in the

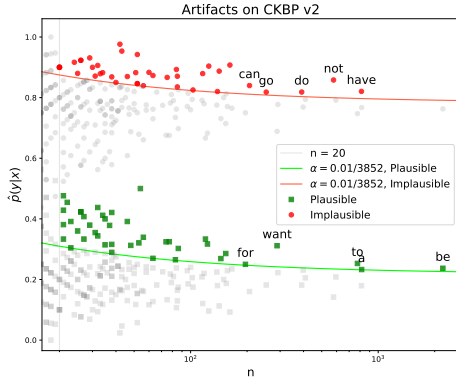


Figure 2: Artifacts statistics of CKBP v2. Colored dots (either square or circle) represent artifacts in the new evaluation set.

annotation process. “Artifacts” may undermine the designated evaluation purpose of the dataset. Thus, it is necessary for us to check if “artifacts” exist in CKBP v2.

We identify artifacts in CKBP v2 by following the previous work Gardner et al. (2021). Particularly, for each word  $x$  in the vocab list<sup>3</sup>, we compute all quantities appearing in the  $z$ -statistic formula

$$z = \frac{\hat{p}(y|x) - p_0}{\sqrt{p_0(1 - p_0)/n}}$$

These include word count  $n$ , estimated probability  $\hat{p}(y|x)$  as the fraction of the number of target label  $y$  in the corresponding  $n$  samples over  $n$ . After that, we compute the  $z$ -statistic and reject or not reject the null hypothesis  $\hat{p}(y|x) = p_0$  with a significance level  $\alpha = 0.01$  and a conservative Bonferroni correction (Bonferroni, 1936) for all 3852 vocabulary items. Note that the “true” probability  $p_0 = p(y|x)$  is taken to be the proportion of samples with label  $y$  in the whole evaluation set. Also, we do not consider artifacts with a word count less than 20, as they are not statistically significant.

Figure 2 shows the plot of word count against the estimated probability  $\hat{p}(y|x)$  for CKBP v2. The additional green and red curves correspond to the largest value of  $\hat{p}(y|x)$  w.r.t  $n$  to keep the null hypothesis from being rejected, where  $y$  takes value “Plausible” and “Implausible” respectively. This means that any dot above the corresponding curve with a frequency of at least 20 is marked as an artifact. The artifacts with the largest word count are labeled in the plot. Overall, CKBP v2 contains relatively few artifacts (83 artifacts out of 3852 vocabulary items), and the artifacts do not significantly

<sup>3</sup>We exclude all relation tokens, as well as special pronoun tokens, namely PersonX, PersonY, PersonZ, PeopleX

affect the evaluation set quality as their frequencies are not high.

## 5 Extrinsic Evaluation

In this section, we study two downstream applications of CKBP. After acquiring a population model, it act as a scoring function to determine whether a triple from the candidate knowledge base  $G$  is plausible or not, thus serving as a source of commonsense knowledge acquisition (Fang et al., 2021b). We leverage the populated knowledge as additional training data for both generative commonsense inference (COMET; Bosselut et al., 2019) and zero-shot commonsense question answering (Ma et al., 2021).

### 5.1 Generative Commonsense Inference (COMET)

**Setup** We follow the basic settings as in the original ATOMIC<sub>20</sub> paper (Hwang et al., 2021) to generate commonsense tails  $t$  given head  $h$  and relation  $r$  as input. The evaluation dataset is the annotated 5,000 test examples provided by Hwang et al. (2021). We use BLEU (Papineni et al., 2002), ROUGE-L (Lin, 2004), METEOR (Lavie and Agarwal, 2007), and CIDEr (Vedantam et al., 2015) as the automatic evaluation metrics.

Specifically, we compare the performance of the following training paradigms: 1) Training the model using the official training set of ATOMIC<sub>20</sub>. 2) Pre-training the model using a comparable amount of CKBP-acquired data, and subsequently fine-tune on ATOMIC<sub>20</sub> training set. 3) Training on a mixture of CKBP-acquired data and ATOMIC<sub>20</sub> training data.

We filter the CKBP-acquired data using two filters. First, we employ two typical population models, RoBERTa-L (Liu et al., 2019) fine-tuned on CKBP training set and Vera (Liu et al., 2023) to provide a plausibility score for each triple. We set an empirical threshold of 0.8 and selecting triples with plausibility score higher than that as populated commonsense knowledge. For the RoBERTa-L model, we select the best-performed checkpoints based on both CKBP v1 and CKBP v2 to evaluate which evaluation set is better aligned with downstream performance. Second, we utilize a diversity filter defined in G-DAUG (Yang et al., 2020), which is a heuristic favoring diverse n-grams. The diversity filter is applied such that we select the same amount of CKBP-acquired data as the training set

Training Data	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	ROUGE-L	CIDEr
ATOMIC	41.8	26.6	19.2	14.5	50.0	21.2	66.1
ATOMIC + CKBP <sub>RoBERTa-L (V1)</sub>	41.9	26.6	18.8	13.8	49.7	21.2	66.2
ATOMIC + CKBP <sub>RoBERTa-L (V2)</sub>	42.5	26.7	18.8	13.8	50.2	21.4	67.1
ATOMIC + CKBP <sub>Vera</sub>	42.9	27.2	19.4	14.4	50.2	21.4	<b>67.5</b>
ATOMIC + CKBP <sub>Vera (mix)</sub>	<b>43.3</b>	<b>27.6</b>	<b>19.7</b>	<b>14.7</b>	<b>50.3</b>	<b>21.5</b>	67.4

Table 3: Performance (%) of GPT2-Large on generative commonsense inference modeling (COMET). ATOMIC stands for ATOMIC<sub>20</sub> training set, and CKBP stands for our CKBP data. Subscripts under CKBP indicating the population model to select populated commonsense knowledge. The best performances are **bold-faced**.

of ATOMIC<sub>20</sub>.

We choose GPT2-Large as our backbone language model. We didn’t use GPT2-XL as in Hwang et al. (2021) because the XL version performs relatively poorer than the Large version in terms of most automatic evaluation metrics on the evaluation set of ATOMIC<sub>20</sub> despite twice the model size. The learning rate is set as 1e-5, and we train the model for three epochs on both CKBP-acquired data and ATOMIC<sub>20</sub> training data.

**Results and Analysis** The results of generative commonsense inference are presented in Table 3. First, adding CKBP-acquired commonsense knowledge for either pre-training or co-training can yield a general performance improvement in generative commonsense inference. Specifically, the model trained on ATOMIC + CKBP<sub>Vera</sub> achieves the best performance and outperforms that only fine-tuned on ATOMIC<sub>20</sub> on all automatic evaluation metrics. This indicates that leveraging the abundant unlabeled discourse knowledge from ASER ( $G$ ), accompanied by appropriate plausibility filtering through the population model, can effectively serve as valuable augmented data to enhance commonsense reasoning. Among the population models, we observe that a better population model, as evaluated by our CKBP v2 evaluation set, corresponds to a higher performance gain in the generative commonsense inference task. This finding highlights the promising potential of developing improved population models, which subsequently contribute to enhanced downstream applications.

Second, the RoBERTa-L model selected by CKBP v2 demonstrates greater efficacy in enhancing generative commonsense inference compared to the model selected by CKBP v1. This finding suggests that CKBP v2 exhibits improved alignment with real-world downstream applications, surpassing its predecessor in terms of practical utility. It’s also noteworthy that COMET is an important task that inherently benefits a pile of further downstream tasks that requires commonsense reason-

ing, including zero-shot commonsense question answering with self-talk (Shwartz et al., 2020) and dynamic graph construction (Bosselut et al., 2021), narrative reasoning (Peng et al., 2022), and dialogue generation (Tu et al., 2022). In this regard, our work exhibits significant potential for generalization to tasks extending beyond the realm of commonsense reasoning.

## 5.2 Zero-shot Commonsense QA

**Setup** For the zero-shot commonsense question answering (QA) task, we adopt the task definition and evaluation pipeline proposed by Ma et al. (2021) to evaluate the benefit CKBP v2 brings to extrinsic QA. Several methods have been proposed to tackle this task, including those by Shwartz et al. (2020); Bosselut et al. (2021); Kim et al. (2022b) The most effective pipeline, as proposed by Ma et al. (2021), injects commonsense knowledge into pre-trained language models through fine-tuning on QA pairs synthesized from knowledge in CSKBs. To perform this fine-tuning, the head  $h$  and relation  $r$  of a  $(h, r, t)$  triple are transformed into a question using natural language prompts, while the tail  $t$  is used as the correct answer option. Distractors or negative examples are created by randomly sampling tails from triples that do not share common keywords with the head. This fine-tuning process enhances the model’s knowledge not only for QA benchmarks constructed from CSKBs, such as SocialQA (Sap et al., 2019b) derived from ATOMIC, but also improves its ability to answer previously unseen commonsense questions in a more generalized manner.

We adopt the original QA synthesis and model training pipeline by Ma et al. (2021) on the original ATOMIC and the one augmented with populated knowledge from CKBP v2 to ablatively study the sole benefit that knowledge in CKBP v2 brings. Similar with that in COMET experiments, we use the best-performed CKBP model, Vera, to score the whole population space in ASER and select

Model	CSKB	a-NLI	CSQA	PIQA	SIQA	WG	Avg.
<b>Zero-shot Baselines</b>							
Random	-	50.0	20.0	50.0	33.3	50.0	40.7
Majority	-	50.8	20.9	50.5	33.6	50.4	41.2
RoBERTa-L (Liu et al., 2019)	-	65.5	45.0	67.6	47.3	57.5	56.6
DeBERTa-v3-L (He et al., 2023b)	-	59.9	25.4	44.8	47.8	50.3	45.6
Self-talk (Shwartz et al., 2020)	-	-	32.4	70.2	46.2	54.7	-
COMET-DynGen (Bosselut et al., 2021)	ATOMIC	-	-	-	50.1	-	-
SMLM (Banerjee and Baral, 2020)	*	65.3	38.8	-	48.5	-	-
MICO (Su et al., 2022)	ATOMIC	-	44.2	-	56.0	-	-
STL-Adapter (Kim et al., 2022b)	ATOMIC	71.3	66.5	71.1	64.4	60.3	66.7
<b>Backbone: DeBERTa-v3-Large <sup>435M</sup></b>							
DeBERTa-v3-L (MR) (Ma et al., 2021)	ATM-10X	75.1	<u>71.6</u>	<b>79.0</b>	59.7	71.7	71.4
DeBERTa-v3-L (MR) (Ma et al., 2021)	ATOMIC	76.0	67.0	<u>78.0</u>	62.1	<u>76.0</u>	71.8
DeBERTa-v3-L (MR) (Ma et al., 2021)	CKBP (our)	<b>79.2</b>	69.6	<u>77.9</u>	64.3	<b>77.2</b>	<b>73.6</b>
<b>Large Language Models</b>							
GPT-3.5 (text-davinci-003)	-	61.8	68.9	67.8	<u>68.0</u>	60.7	65.4
ChatGPT (gpt-3.5-turbo)	-	69.3	<b>74.5</b>	75.1	<b>69.5</b>	62.8	70.2
<b>Supervised Learning &amp; Human Performance</b>							
RoBERTa-L (Supervised)	-	85.6	78.5	79.2	76.6	79.3	79.8
DeBERTa-v3-L (Supervised)	-	89.0	82.1	84.5	80.1	84.1	84.0
Human Performance	-	91.4	88.9	94.9	86.9	94.1	91.2

Table 4: Zero-shot evaluation results (%) on five commonsense question answering benchmarks. The best results are **bold-faced**, and the second-best ones are underlined. The performance of supervised learning and human are for reference only.

the populated knowledge with plausibility scores of over 0.8. Then the same diversity filter as in Section 5.1 is used to downsample the number of populated triples to be comparable with the size of the training set in ATOMIC<sub>20</sub><sup>20</sup>. For the QA model, DeBERTa-v3-Large (He et al., 2023b) is used as the backbone, and we train the model using a learning rate of 7e-6 for one epochs on both the CKBP-acquired data and ATOMIC-synthesized data as provided by Ma et al. (2021).

Once trained, we evaluate the model on the validation splits of five commonsense QA benchmarks: Abductive NLI (aNLI; Bhagavatula et al., 2020), CommonsenseQA (CSQA; Talmor et al., 2019), PhysicalIQA (PIQA; Bisk et al., 2020b), SocialIQA (SIQA; Sap et al., 2019b), and WinoGrande (WG; Sakaguchi et al., 2021). Accuracy is used as the evaluation metric. Furthermore, we compare our model not only against existing zero-shot knowledge injection methods (Shwartz et al., 2020; Bosselut et al., 2021; Banerjee and Baral, 2020; Su et al., 2022; Kim et al., 2022b; Ma et al., 2021) but also against large language models such as ChatGPT (OpenAI, 2022) and GPT-3.5 (Brown et al., 2020b).

**Results and Analysis** The zero-shot commonsense QA results are shown in Table 4. Among all the zero-shot methods, the model trained on CKBP v2 demonstrates the highest performance. It out-

performs models trained solely on ATOMIC (with an increase of 2.2%) and ATOMIC10X (West et al., 2022b) (with an increase of 1.8%). Importantly, our method surpasses large language models by an average of 3.4%. This performance gain highlights the significant advantage of our populated commonsense knowledge over both human annotations and distilled knowledge from large language models. Furthermore, we observe that the model trained on CKBP-acquired data shows the most improvement on the aNLI and WinoGrande benchmarks. One potential reason for this is that the populated knowledge in CKBP v1 encompasses a wider range of commonsense knowledge beyond only social commonsense, which benefits tasks involving abductive reasoning (based on narrative) and pronoun coreference resolution.

## 6 Conclusion

In this paper, we introduce a new CSKB Population benchmark CKBP v2 which addresses two problems of the predecessor CKBP v1. Besides, we conduct a broad range of experiments with different models, including GPT3.5 and ChatGPT, on the new evaluation set. The result shows that the CSKB Population task remains a hard task of commonsense reasoning even for state-of-the-art LLMs, which challenges the community for future research.



## 610 Limitations

611 We observe several limitations of this work. First,  
612 CKBP v2 still follows the lemmatized format of  
613 events, which may hinder the usage of the resulting  
614 population model on knowledge bases other than  
615 ASER. Second, the paradigm of CSKB is context-  
616 free, which may have difficulty in directly applying  
617 to actual downstream tasks. Third, As this paper  
618 focuses on proposing a new evaluation set of the  
619 CSKB Population, we do not present novel tailored  
620 methods for solving this task, leaving it to future  
621 research.

## 622 Ethical Statements

623 This work presents CKBP v2, an open-source  
624 benchmark for the research community to study  
625 the CSKB population problem. The training set is  
626 directly adapted from CKBP v1 and ATOMIC<sub>(20)</sub>,  
627 GLUCOSE, and ConceptNet, which would have  
628 the same ethical issues as in those previous works.  
629 Instances in the evaluation set are retrieved from  
630 CKBP v1 and ASER, both being open-source with  
631 an MIT license. Events in all data instances are  
632 anonymized. Thus, the benchmark does not pose  
633 any privacy problems about any specific entities  
634 (e.g., a person or company). We carried out human  
635 expert annotation, where annotators are fairly paid  
636 according to the minimum wage requirement of the  
637 local government.

## 638 References

639 Pratyay Banerjee and Chitta Baral. 2020. [Self-](#)  
640 [supervised knowledge triplet learning for zero-shot](#)  
641 [question answering](#). In *Proceedings of the 2020 Con-*  
642 *ference on Empirical Methods in Natural Language*  
643 *Processing, EMNLP 2020, Online, November 16-20,*  
644 *2020*, pages 151–162. Association for Computational  
645 Linguistics.

646 Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wen-  
647 liang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei  
648 Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu,  
649 and Pascale Fung. 2023. [A multitask, multilingual,](#)  
650 [multimodal evaluation of chatgpt on reasoning, hal-](#)  
651 [lucination, and interactivity](#). *CoRR*, abs/2302.04023.

652 Chandra Bhagavatula, Ronan Le Bras, Chaitanya  
653 Malaviya, Keisuke Sakaguchi, Ari Holtzman, Han-  
654 nah Rashkin, Doug Downey, Wen-tau Yih, and Yejin  
655 Choi. 2020. [Abductive commonsense reasoning](#). In  
656 *8th International Conference on Learning Representa-*  
657 *tions, ICLR 2020, Addis Ababa, Ethiopia, April*  
658 *26-30, 2020*. OpenReview.net.

Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng  
Gao, and Yejin Choi. 2020a. [PIQA: reasoning about](#)  
[physical commonsense in natural language](#). In *The*  
*Thirty-Fourth AAAI Conference on Artificial Intelli-*  
*gence, AAAI 2020, The Thirty-Second Innovative Ap-*  
*plications of Artificial Intelligence Conference, IAAI*  
*2020, The Tenth AAAI Symposium on Educational*  
*Advances in Artificial Intelligence, EAAI 2020, New*  
*York, NY, USA, February 7-12, 2020*, pages 7432–  
7439. AAAI Press.

Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng  
Gao, and Yejin Choi. 2020b. [PIQA: reasoning about](#)  
[physical commonsense in natural language](#). In *The*  
*Thirty-Fourth AAAI Conference on Artificial Intelli-*  
*gence, AAAI 2020, The Thirty-Second Innovative Ap-*  
*plications of Artificial Intelligence Conference, IAAI*  
*2020, The Tenth AAAI Symposium on Educational*  
*Advances in Artificial Intelligence, EAAI 2020, New*  
*York, NY, USA, February 7-12, 2020*, pages 7432–  
7439. AAAI Press.

Kurt D. Bollacker, Colin Evans, Praveen K. Paritosh,  
Tim Sturge, and Jamie Taylor. 2008. [Freebase: a](#)  
[collaboratively created graph database for structuring](#)  
[human knowledge](#). In *Proceedings of the ACM SIG-*  
*MOD International Conference on Management of*  
*Data, SIGMOD 2008, Vancouver, BC, Canada, June*  
*10-12, 2008*, pages 1247–1250. ACM.

C. Bonferroni. 1936. [Teoria statistica delle classi e](#)  
[calcolo delle probabilita](#). *Pubblicazioni del R Istituto*  
*Superiore di Scienze Economiche e Commerciali di*  
*Firenze*, 8:3–62.

Antoine Bordes, Nicolas Usunier, Alberto García-  
Durán, Jason Weston, and Oksana Yakhnenko.  
2013. [Translating embeddings for modeling multi-](#)  
[relational data](#). In *Advances in Neural Information*  
*Processing Systems 26: 27th Annual Conference on*  
*Neural Information Processing Systems 2013. Pro-*  
*ceedings of a meeting held December 5-8, 2013, Lake*  
*Tahoe, Nevada, United States*, pages 2787–2795.

Antoine Bosselut, Ronan Le Bras, and Yejin Choi. 2021. [Dynamic neuro-symbolic knowledge graph construc-](#)  
[tion for zero-shot commonsense question answering](#).  
In *Thirty-Fifth AAAI Conference on Artificial Intel-*  
*ligence, AAAI 2021, Thirty-Third Conference on In-*  
*novative Applications of Artificial Intelligence, IAAI*  
*2021, The Eleventh Symposium on Educational Ad-*  
*vances in Artificial Intelligence, EAAI 2021, Virtual*  
*Event, February 2-9, 2021*, pages 4923–4931. AAAI  
Press.

Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chai-  
tanya Malaviya, Asli Celikyilmaz, and Yejin Choi.  
2019. [COMET: commonsense transformers for auto-](#)  
[matic knowledge graph construction](#). In *Proceedings*  
*of the 57th Conference of the Association for Compu-*  
*tational Linguistics, ACL 2019, Florence, Italy, July*  
*28- August 2, 2019, Volume 1: Long Papers*, pages  
4762–4779. Association for Computational Linguis-  
tics.

717	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	778
718	Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind	Kristina Toutanova. 2019. <a href="#">BERT: pre-training of</a>	779
719	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	<a href="#">deep bidirectional transformers for language under-</a>	780
720	Askeff, Sandhini Agarwal, Ariel Herbert-Voss,	<a href="#">standing</a> . In <i>Proceedings of the 2019 Conference of</i>	781
721	Gretchen Krueger, Tom Henighan, Rewon Child,	<i>the North American Chapter of the Association for</i>	782
722	Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,	<i>Computational Linguistics: Human Language Tech-</i>	783
723	Clemens Winter, Christopher Hesse, Mark Chen, Eric	<i>nologies, NAACL-HLT 2019, Minneapolis, MN, USA,</i>	784
724	Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,	<i>June 2-7, 2019, Volume 1 (Long and Short Papers),</i>	785
725	Jack Clark, Christopher Berner, Sam McCandlish,	pages 4171–4186. Association for Computational	786
726	Alec Radford, Ilya Sutskever, and Dario Amodei.	Linguistics.	787
727	2020a. <a href="#">Language models are few-shot learners</a> . In	Tianqing Fang, Quyet V. Do, Hongming Zhang,	788
728	<i>Advances in Neural Information Processing Systems</i>	Yangqiu Song, Ginny Y. Wong, and Simon See. 2022.	789
729	<i>33: Annual Conference on Neural Information Pro-</i>	<a href="#">Pseudoreasoner: Leveraging pseudo labels for com-</a>	790
730	<i>cessing Systems 2020, NeurIPS 2020, December 6-</i>	<a href="#">monsense knowledge base population</a> . In <i>Findings</i>	791
731	<i>12, 2020, virtual</i> .	<i>of the Association for Computational Linguistics:</i>	792
732	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie	<i>EMNLP 2022, Abu Dhabi, United Arab Emirates,</i>	793
733	Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind	<i>December 7-11, 2022</i> , pages 3379–3394. Associa-	794
734	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	tion for Computational Linguistics.	795
735	Askeff, Sandhini Agarwal, Ariel Herbert-Voss,	Tianqing Fang, Weiqi Wang, Sehyun Choi, Shibo Hao,	796
736	Gretchen Krueger, Tom Henighan, Rewon Child,	Hongming Zhang, Yangqiu Song, and Bin He. 2021a.	797
737	Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,	<a href="#">Benchmarking commonsense knowledge base pop-</a>	798
738	Clemens Winter, Christopher Hesse, Mark Chen, Eric	<a href="#">ulation with an effective evaluation dataset</a> . In <i>Pro-</i>	799
739	Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,	<i>ceedings of the 2021 Conference on Empirical Meth-</i>	800
740	Jack Clark, Christopher Berner, Sam McCandlish,	<i>ods in Natural Language Processing, EMNLP 2021,</i>	801
741	Alec Radford, Ilya Sutskever, and Dario Amodei.	<i>Virtual Event / Punta Cana, Dominican Republic, 7-</i>	802
742	2020b. <a href="#">Language models are few-shot learners</a> . In	<i>11 November, 2021</i> , pages 8949–8964. Association	803
743	<i>Advances in Neural Information Processing Systems</i>	for Computational Linguistics.	804
744	<i>33: Annual Conference on Neural Information Pro-</i>	Tianqing Fang, Hongming Zhang, Weiqi Wang,	805
745	<i>cessing Systems 2020, NeurIPS 2020, December 6-</i>	Yangqiu Song, and Bin He. 2021b. <a href="#">DISCOS: bridg-</a>	806
746	<i>12, 2020, virtual</i> .	<a href="#">ing the gap between discourse knowledge and com-</a>	807
747	Chunkit Chan, Jiayang Cheng, Weiqi Wang, Yuxin	<a href="#">monsense knowledge</a> . In <i>WWW '21: The Web Con-</i>	808
748	Jiang, Tianqing Fang, Xin Liu, and Yangqiu Song.	<i>ference 2021, Virtual Event / Ljubljana, Slovenia,</i>	809
749	2023. <a href="#">Chatgpt evaluation on sentence level relations:</a>	<i>April 19-23, 2021</i> , pages 2648–2659. ACM / IW3C2.	810
750	<a href="#">A focus on temporal, causal, and discourse relations.</a>	Silin Gao, Jena D. Hwang, Saya Kanno, Hiromi Wakaki,	811
751	<i>CoRR</i> , abs/2304.14827.	Yuki Mitsufuji, and Antoine Bosselut. 2022. <a href="#">Com-</a>	812
752	Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,	<a href="#">fact: A benchmark for linking contextual common-</a>	813
753	Maarten Bosma, Gaurav Mishra, Adam Roberts,	<a href="#">sense knowledge</a> . In <i>Findings of the Association</i>	814
754	Paul Barham, Hyung Won Chung, Charles Sutton,	<i>for Computational Linguistics: EMNLP 2022, Abu</i>	815
755	Sebastian Gehrmann, Parker Schuh, Kensen Shi,	<i>Dhabi, United Arab Emirates, December 7-11, 2022,</i>	816
756	Sasha Tsvyashchenko, Joshua Maynez, Abhishek	pages 1656–1675. Association for Computational	817
757	Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vin-	Linguistics.	818
758	odkumar Prabhakaran, Emily Reif, Nan Du, Ben	Matt Gardner, William Merrill, Jesse Dodge, Matthew E.	819
759	Hutchinson, Reiner Pope, James Bradbury, Jacob	Peters, Alexis Ross, Sameer Singh, and Noah A.	820
760	Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin,	Smith. 2021. <a href="#">Competency problems: On finding</a>	821
761	Toju Duke, Anselm Levskaya, Sanjay Ghemawat,	<a href="#">and removing artifacts in language data</a> . In <i>Proceed-</i>	822
762	Sunipa Dev, Henryk Michalewski, Xavier Garcia,	<i>ings of the 2021 Conference on Empirical Methods</i>	823
763	Vedant Misra, Kevin Robinson, Liam Fedus, Denny	<i>in Natural Language Processing, EMNLP 2021, Vir-</i>	824
764	Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim,	<i>tual Event / Punta Cana, Dominican Republic, 7-11</i>	825
765	Barret Zoph, Alexander Spiridonov, Ryan Sepassi,	<i>November, 2021</i> , pages 1801–1813. Association for	826
766	David Dohan, Shivani Agrawal, Mark Omernick, An-	Computational Linguistics.	827
767	drew M. Dai, Thanumalayan Sankaranarayanan Pil-	Deepanway Ghosal, Siqi Shen, Navonil Majumder,	828
768	lai, Marie Pellat, Aitor Lewkowycz, Erica Moreira,	Rada Mihalcea, and Soujanya Poria. 2022. <a href="#">CICERO:</a>	829
769	Rewon Child, Oleksandr Polozov, Katherine Lee,	<a href="#">A dataset for contextualized commonsense inference</a>	830
770	Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark	<a href="#">in dialogues</a> . In <i>Proceedings of the 60th Annual Meet-</i>	831
771	Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy	<i>ing of the Association for Computational Linguistics</i>	832
772	Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov,	<i>(Volume 1: Long Papers), ACL 2022, Dublin, Ireland,</i>	833
773	and Noah Fiedel. 2022. <a href="#">Palm: Scaling language mod-</a>	<i>May 22-27, 2022</i> , pages 5010–5028. Association for	834
774	<a href="#">eling with pathways</a> . <i>CoRR</i> , abs/2204.02311.	Computational Linguistics.	835
775	Ernest Davis. 2023. <a href="#">Benchmarks for automated</a>		
776	<a href="#">commonsense reasoning: A survey</a> . <i>CoRR</i> ,		
777	abs/2302.04752.		

836	Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel R. Bowman, and Noah A. Smith. 2018. <a href="#">Annotation artifacts in natural language inference data</a> . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 2 (Short Papers)</i> , pages 107–112. Association for Computational Linguistics.	894
837		895
838		
839		896
840		897
841		898
842		899
843		900
844		901
845		902
846	William L. Hamilton, Zhitao Ying, and Jure Leskovec. 2017. <a href="#">Inductive representation learning on large graphs</a> . In <i>Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA</i> , pages 1024–1034.	903
847		904
848		905
849		906
850		907
851		908
852	Hangfeng He, Hongming Zhang, and Dan Roth. 2023a. <a href="#">Rethinking with retrieval: Faithful large language model inference</a> . <i>CoRR</i> , abs/2301.00303.	909
853		910
854		911
855	Pengcheng He, Jianfeng Gao, and Weizhu Chen. 2023b. <a href="#">DeBERTav3: Improving deBERTa using ELECTRA-style pre-training with gradient-disentangled embedding sharing</a> . In <i>The Eleventh International Conference on Learning Representations</i> .	912
856		913
857		914
858		915
859		916
860	Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. 2022. <a href="#">Training compute-optimal large language models</a> . <i>CoRR</i> , abs/2203.15556.	917
861		918
862		919
863		920
864		
865		921
866		922
867		923
868		924
869	Jena D. Hwang, Chandra Bhagavatula, Ronan Le Bras, Jeff Da, Keisuke Sakaguchi, Antoine Bosselut, and Yejin Choi. 2021. <a href="#">(comet-) atomic 2020: On symbolic and neural commonsense knowledge graphs</a> . In <i>Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021</i> , pages 6384–6392. AAAI Press.	925
870		926
871		927
872		928
873		
874		929
875		930
876		931
877		932
878		933
879		
880	Hyunwoo Kim, Jack Hessel, Liwei Jiang, Ximing Lu, Youngjae Yu, Pei Zhou, Ronan Le Bras, Malihe Alikhani, Gunhee Kim, Maarten Sap, and Yejin Choi. 2022a. <a href="#">SODA: million-scale dialogue distillation with social commonsense contextualization</a> . <i>CoRR</i> , abs/2212.10465.	934
881		935
882		936
883		937
884		938
885		939
886	Yu Jin Kim, Beong-woo Kwak, Youngwook Kim, Reinald Kim Amplayo, Seung-won Hwang, and Jinyoung Yeo. 2022b. <a href="#">Modularized transfer learning with multiple knowledge graphs for zero-shot commonsense reasoning</a> . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2022, Seattle, WA, United States, July 10-15, 2022</i> , pages 2244–2257. Association for Computational Linguistics.	940
887		941
888		942
889		943
890		944
891		
892		945
893		946
		947
		948
		949
		950
		896
		897
		898
		899
		900
		901
		902
		903
		904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950



951	<i>Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020</i> , pages 2925–2933. AAAI Press.	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. <a href="#">Winogrande: an adversarial winograd schema challenge at scale</a> . <i>Commun. ACM</i> , 64(9):99–106.	1005 1006 1007 1008
955	George A. Miller. 1995. <a href="#">Wordnet: A lexical database for english</a> . volume 38, pages 39–41.	Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A. Smith, and Yejin Choi. 2019a. <a href="#">ATOMIC: an atlas of machine commonsense for if-then reasoning</a> . In <i>The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019</i> , pages 3027–3035. AAAI Press.	1009 1010 1011 1012 1013 1014 1015 1016 1017 1018 1019 1020
956	Nasrin Mostafazadeh, Aditya Kalyanpur, Lori Moon, David W. Buchanan, Lauren Berkowitz, Or Biran, and Jennifer Chu-Carroll. 2020. <a href="#">GLUCOSE: generalized and contextualized story explanations</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020</i> , pages 4569–4586. Association for Computational Linguistics.	Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019b. <a href="#">Social iqa: Commonsense reasoning about social interactions</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019</i> , pages 4462–4472. Association for Computational Linguistics.	1021 1022 1023 1024 1025 1026 1027 1028 1029
965	OpenAI. 2022. <a href="#">Chatgpt: Optimizing language models for dialogue</a> . <i>OpenAI</i> .	Chao Shang, Yun Tang, Jing Huang, Jinbo Bi, Xiaodong He, and Bowen Zhou. 2019. <a href="#">End-to-end structure-aware convolutional networks for knowledge base completion</a> . pages 3060–3067.	1030 1031 1032 1033
966	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. <a href="#">Bleu: a method for automatic evaluation of machine translation</a> . In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, July 6-12, 2002, Philadelphia, PA, USA</i> , pages 311–318. ACL.	Vered Shwartz, Peter West, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2020. <a href="#">Unsupervised commonsense question answering with self-talk</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020</i> , pages 4615–4629. Association for Computational Linguistics.	1034 1035 1036 1037 1038 1039 1040
967	Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou Yu, Weizhu Chen, and Jianfeng Gao. 2023. <a href="#">Check your facts and try again: Improving large language models with external knowledge and automated feedback</a> . <i>CoRR</i> , abs/2302.12813.	Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. <a href="#">Conceptnet 5.5: An open multilingual graph of general knowledge</a> . In <i>Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA</i> , pages 4444–4451. AAAI Press.	1041 1042 1043 1044 1045 1046
973	Xiangyu Peng, Siyan Li, Sarah Wiegrefe, and Mark Riedl. 2022. <a href="#">Inferring the reader: Guiding automated story generation with commonsense reasoning</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 7008–7029, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	Ying Su, Zihao Wang, Tianqing Fang, Hongming Zhang, Yangqiu Song, and Tong Zhang. 2022. <a href="#">MICO: A multi-alternative contrastive learning framework for commonsense knowledge representation</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022</i> , pages 1339–1351. Association for Computational Linguistics.	1047 1048 1049 1050 1051 1052 1053 1054
974	Adam Poliak, Jason Naradowsky, Aparajita Haldar, Rachel Rudinger, and Benjamin Van Durme. 2018. <a href="#">Hypothesis only baselines in natural language inference</a> . In <i>Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics, *SEM@NAACL-HLT 2018, New Orleans, Louisiana, USA, June 5-6, 2018</i> , pages 180–191. Association for Computational Linguistics.	Zhiqing Sun, Zhi-Hong Deng, Jian-Yun Nie, and Jian Tang. 2019. <a href="#">Rotate: Knowledge graph embedding by relational rotation in complex space</a> . In <i>7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019</i> . OpenReview.net.	1055 1056 1057 1058 1059 1060
975	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. <a href="#">Commonsenseqa: A question</a>	1061 1062
976	Itsumi Saito, Kyosuke Nishida, Hisako Asano, and Junji Tomita. 2018. <a href="#">Commonsense knowledge base completion and generation</a> . In <i>Proceedings of the 22nd Conference on Computational Natural Language Learning, CoNLL 2018, Brussels, Belgium, October 31 - November 1, 2018</i> , pages 141–150. Association for Computational Linguistics.		
977			
978			
979			
980			
981			
982			
983			
984			
985			
986			
987			
988			
989			
990			
991			
992			
993			
994			
995			
996			
997			
998			
999			
1000			
1001			
1002			
1003			
1004			



1063	<a href="#">answering challenge targeting commonsense knowledge</a> . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)</i> , pages 4149–4158. Association for Computational Linguistics.	1121
1064		1122
1065		1123
1066		
1067		
1068		
1069		
1070		
1071	Quan Tu, Yanran Li, Jianwei Cui, Bin Wang, Ji-Rong Wen, and Rui Yan. 2022. <a href="#">MISC: A mixed strategy-aware model integrating COMET for emotional support conversation</a> . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 308–319, Dublin, Ireland. Association for Computational Linguistics.	
1072		
1073		
1074		
1075		
1076		
1077		
1078		
1079	Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. 2015. <a href="#">Cider: Consensus-based image description evaluation</a> . In <i>IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015</i> , pages 4566–4575. IEEE Computer Society.	
1080		
1081		
1082		
1083		
1084		
1085	Peter West, Chandra Bhagavatula, Jack Hessel, Jena D. Hwang, Liwei Jiang, Ronan Le Bras, Ximing Lu, Sean Welleck, and Yejin Choi. 2022a. <a href="#">Symbolic knowledge distillation: from general language models to commonsense models</a> . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2022, Seattle, WA, United States, July 10-15, 2022</i> , pages 4602–4625. Association for Computational Linguistics.	
1086		
1087		
1088		
1089		
1090		
1091		
1092		
1093		
1094		
1095	Peter West, Chandra Bhagavatula, Jack Hessel, Jena D. Hwang, Liwei Jiang, Ronan Le Bras, Ximing Lu, Sean Welleck, and Yejin Choi. 2022b. <a href="#">Symbolic knowledge distillation: from general language models to commonsense models</a> . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL 2022, Seattle, WA, United States, July 10-15, 2022</i> , pages 4602–4625. Association for Computational Linguistics.	
1096		
1097		
1098		
1099		
1100		
1101		
1102		
1103		
1104		
1105	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. <a href="#">Transformers: State-of-the-art natural language processing</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, EMNLP 2020 - Demos, Online, November 16-20, 2020</i> , pages 38–45. Association for Computational Linguistics.	
1106		
1107		
1108		
1109		
1110		
1111		
1112		
1113		
1114		
1115		
1116		
1117		
1118	Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. 2015. <a href="#">Embedding entities and relations for learning and inference in knowledge bases</a> . In <i>3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings</i> .	1124
1119		1125
1120		1126
		1127
		1128
		1129
		1130
		1131
		1132
		1133
		1134
		1135
		1136
		1137
		1138
		1139
		1140
		1141
		1142
		1143
		1144
		1145