

End-to-end Case-Based Reasoning for Commonsense Knowledge Base Completion

Zonglin Yang^{♣*} Xinya Du^{♣*} Erik Cambria[♣] Claire Cardie[♦]

[♣] Nanyang Technological University

[♣] University of Texas at Dallas

[♦] Cornell University

{zonglin.yang, cambria}@ntu.edu.sg

xinya.du@utdallas.edu

cardie@cs.cornell.edu

Abstract

Pretrained language models have been shown to store knowledge in their parameters and have achieved reasonable performance in commonsense knowledge base completion (CKBC) tasks. However, CKBC is knowledge-intensive and it is reported that pretrained language models’ performance in knowledge-intensive tasks are limited because of their incapability of accessing and manipulating knowledge. As a result, we hypothesize that providing retrieved passages that contain relevant knowledge as additional input to the CKBC task will improve performance. In particular, we draw insights from Case-Based Reasoning (CBR) – which aims to solve a new problem by reasoning with retrieved relevant cases, and investigate the direct application of it to CKBC. On two benchmark datasets, we demonstrate through automatic and human evaluations that our **End-to-end Case-Based Reasoning Framework (ECBRF)** generates more valid knowledge than the state-of-the-art COMET model for CKBC in both the fully supervised and few-shot settings. From the perspective of CBR, our framework addresses a fundamental question on whether CBR methodology can be utilized to improve deep learning models.

1 Introduction

Commonsense knowledge helps humans navigate everyday situations seamlessly (Apperly, 2010) and is required for many intelligent scenarios (Davis and Marcus, 2015). To automatically enlarge the scale of commonsense knowledge base for the benefit of reducing labeling labor and expense, Knowledge Graph Completion (KGC) has become a hot research topic (Ji et al., 2022). The general KGC task is to expand existing knowledge graphs by using well-trained classifiers—they are trained with existing annotated samples and predict whether or not there is a relationship between two *existing* entities in a knowledge graph (Wang et al., 2017).

*Contribution starts from their stay at Cornell.

Subject	Relation	Object
hardware shop	at location	mall
world map	has property	draw with grid-lines
PersonX receives its reward	wants to	keep the prize
PersonX wins the big Jackpot	wants to	get its money

Table 1: Commonsense knowledge base tuples. Examples are from ConceptNet and Atomic.

Although KGC methods can automatically find unlabeled relationships, they are always classification or ranking tasks and are limited to *existing* entities in a knowledge graph and can’t extend to new entities (Ji et al., 2022). To extend to new entities, COMET (Bosselut et al., 2019) proposes to use *text generation* for exploring and discovering new entities, which is called commonsense knowledge base completion (CKBC) task, utilizing the knowledge within the pretrained language models (PLM), which has been with process in recent years (Devlin et al., 2019; Radford et al., 2019). Specifically, COMET uses subject and relation as direct input to PLM and aims to *generate* objects, most of which are novel and unseen entities.

However, CKBC is knowledge-intensive, requiring wide-ranging and detailed knowledge; and it is reported that the ability of pretrained language models to access and precisely manipulate knowledge is limited (Lewis et al., 2020b). One potential solution to this is to provide “non-parametric knowledge” through additional input. Lewis et al. (2020b) and Guu et al. (2020), for example, have shown that by retrieving passages that contain knowledge relevant to the current task, performance can be improved. For CKBC, unfortunately, it might be especially difficult to find useful passages that contain relevant commonsense knowledge from the web due to a *reporting bias* (Gordon and Van Durme, 2013) in which people rarely express the obvious (i.e., commonsense knowledge).

An example of reporting bias from Table 1 is

that people rarely say “when a person wins a big Jackpot, he/she will want to get its money” because it’s too obvious and meaningless to say. Therefore, instead of retrieving passages from the web, we propose that benefits can still be gained by retrieving relevant knowledge from a “case base” of existing commonsense knowledge tuples¹ and using the retrieved knowledge as non-parametric knowledge (i.e., beyond that represented in the model parameters) to augment the current CKBC input example. In addition, to prevent ECBRF from overfitting to some commonly retrieved cases, we propose *random mask* as a training strategy that randomly masks the retrieved cases during training, which functions similar to dropout (Srivastava et al., 2014) and further improves the performance of the framework. We also analyze several variations to better understand the process.

Although past attempts suggest that similar retrieval-based methods cannot improve the performance of CKBC (Wang et al., 2021), on two benchmark datasets, we demonstrate through automatic and human evaluations that our End-to-end Case-Based Reasoning Framework (ECBRF)² generates more valid and informative knowledge than (1) the state-of-the-art COMET model (Bosse-lut et al., 2019) for CKBC which employs no case retrieval, and (2) a baseline model that employs random case retrieval – on both fully supervised and few-shot settings. We also provide an analysis on why different conclusions are reached.

In addition, our framework draws insight from Case-Based Reasoning (CBR), and also has contributions to the CBR research. CBR is a subject in classical AI that solves a new problem by reusing the solutions of retrieved seen similar problems stored in the case base (Aamodt and Plaza, 1994). CBR’s methodology has four steps — case retrieval, reuse, revise and retain. Past years of accomplishment in deep learning (DL) have led to enthusiasm in the CBR community to apply DL in the service of CBR. However, based on the observation that many challenges remain in DL where CBR has advantages (e.g. few-shot learning), some CBR researchers (Leake and Crandall, 2020) advocate using CBR to complement DL. However, past works on using CBR to complement DL only limit to shallow Neural Networks (NN) (Liao et al., 2018; Leake et al., 2021; Ye et al., 2021, 2022).

¹Initialized with tuples from training set or external data.

²Code available at https://github.com/ZonglinY/ECBRF_Case-Based_Reasoning_with_PLM.git

The latest work even suggests that in many tasks NN itself outperforms CBR-complemented NN (Ye et al., 2022), which raises fundamental questions on whether CBR methodology is useful for DL.

Our work addresses this doubt by being the first to show a concrete implementation of the integration of the full methodology of CBR to PLM (as one typical model in DL) (we show in §6 that we simulate the third step in the methodology of CBR instead of actually implement it, since it requires huge human efforts) and show that the integration method can benefit from multiple steps in the methodology of CBR, and can lead to better performance over PLM itself in both fully supervised settings and few-shot settings on CKBC. Notably our proposed framework has a larger advantage in few-shot settings, where CBR methods typically have advantage. We also find that the generation of our framework is largely related to the retrieved case especially when they are similar, which exhibits strong case-based reasoning patterns. In addition, a detailed analysis of our framework from a CBR perspective is provided in §6.

Our contributions can be summarized as follows: (1) Drawing insights from CBR, we introduce a new end-to-end framework for CKBC task. We also propose training strategies that can better utilize the retrieved knowledge. (2) We conduct extensive experiments on the CKBC task in various settings (e.g. fully supervised and few-shot), and the results consistently demonstrate that our proposed framework achieves improvements over the state-of-the-art baseline methods. (3) From the perspective of the CBR community, whether CBR methodology can be used to improve DL models remains a fundamental research question. We address this doubt by being the first to show a concrete implementation of the integration of the full methodology of CBR to PLM, and showing that such integration can achieve better performance than single PLM. A thorough analysis of the integration from CBR perspective is also provided.

2 Related Work

Case-Based Reasoning CBR is a subject in classical AI which consists of 4 sub-processes in its methodology: *retrieve*, *reuse*, *revise* and *retain* (Aamodt and Plaza, 1994). Leake and Crandall (2020) advocate using CBR to complement the challenges in deep learning (e.g., few-shot learning). §A.4 provides more detailed related works

relevant to this line. Specifically, our framework is inspired by [Watson \(1999\)](#)’s proposal that compared to CBR being described as an artificial intelligence technology, it is better to describe CBR as a methodology for problem solving, that may use any appropriate technology. Here we treat CBR as a methodology and deep learning as technology that uses CBR as the general high-level process and deep learning as components of the process.

Reasoning in NLP CBR could be seen as a type of analogical reasoning ([Kolodner, 1997](#)), and analogical reasoning belongs to inductive reasoning ([Salmon, 1989](#)). Inductive reasoning ([Yang et al., 2022](#)) is different from deductive reasoning ([Clark et al., 2020](#)) (both belong to logical reasoning) that the premise in inductive reasoning can not provide conclusive support to its conclusion.

Commonsense Knowledge Base Completion

Here we mainly describe works that use text generation models for this task. [Li et al. \(2016\)](#) propose models to evaluate the full knowledge tuple rather than generate new knowledge. [Saito et al. \(2018\)](#) make an extension by proposing a joint model for the completion and generation of commonsense tuples. However, their work focuses on augmenting knowledge base completion model, rather than to increase coverage in commonsense knowledge base construction. [Yao et al. \(2019\)](#) and [Malaviya et al. \(2020\)](#) focus on link prediction and ranking of knowledge, which is a different task with our generative CKBC task. [Sap et al. \(2019\)](#) use LSTM ([Hochreiter and Schmidhuber, 1997](#)) to generate commonsense knowledge and [Bosselut et al. \(2019\)](#) further leverage pre-trained language models to generate commonsense knowledge. [Gabriel et al. \(2021\)](#) present the task of discourse-aware commonsense inference and proposes a memory-based model to generate commonsense knowledge that is more coherent with context. [Wang et al. \(2021\)](#) give an analysis on knowledge capacity, transferability, and induction of pre-trained language models to perform generalizable commonsense inference. [Da et al. \(2021\)](#) analyze the few-shot learning ability of pretrained language models for CKBC task. Unlike these works, we propose a model that can improve the performance of generative CKBC tasks in both fully supervised settings and few-shot settings.

Language Model Prompting First developed by the GPT series ([Brown et al., 2020](#)), retrieved data

are used as augmented input to improve few-shot performance of remarkable large models. However, past research suggest that such in-context learning cannot improve the CKBC task ([Wang et al., 2021](#)), and we are the first to show how in-context learning is useful for CKBC. In addition, such large models are hard to obtain and [Brown et al. \(2020\)](#) do not explore the finetuning performance, neither do they explore the full CBR methodology’s effect on PLM. [Gao et al. \(2021\)](#) use prompting and also incorporate demonstrations into context to improve few-shot performance. Their work, however, only focuses on classification tasks and regression tasks, which is different from the CKBC. Similar to our work, [Das et al. \(2021\)](#) use retrieved cases as prompt to improve the performance of PLM. However, they only focus on question answering task and do not integrate the full methodology of CBR, missing important steps such as retain.

3 Task Definition

In the generative CKBC task, a knowledge data instance is represented as a tuple of subject, relation, and object: (sub, rel, obj) . All sub and obj are in natural language phrases (Figure 1). rel can be used as either a special token or the corresponding natural language phrases ([Bosselut et al., 2019](#)). Here we use rel as natural language phrases. The task is that given a pair of sub and rel , the goal is to *generate* the corresponding obj .

4 Methodology

We start by formalizing our framework as a retrieve-then-predict generative process. Then in §4.2, we describe our ECBRF’s modules for the generative process in detail. Finally, we present a hybrid training strategy for better regularization.

Figure 1 describes our method. In the figure, “query” stands for a sub and rel pair which is used as input to ECBRF to generate obj . “Case Base” is initialized with knowledge triples from the training set. “Cases” means the retrieved knowledge triples from the “Case Base”. “In-context demonstrations” stand for the retrieved cases that are used for input augmentation (concatenate with the query). The subject, relation, and object of the retrieved cases are sub-scripted with “r” (e.g., sub_r).

4.1 ECBRF’s Generative Process

ECBRF takes x as input and learns a distribution $p(y|x)$ over possible outputs y . Here x consists of

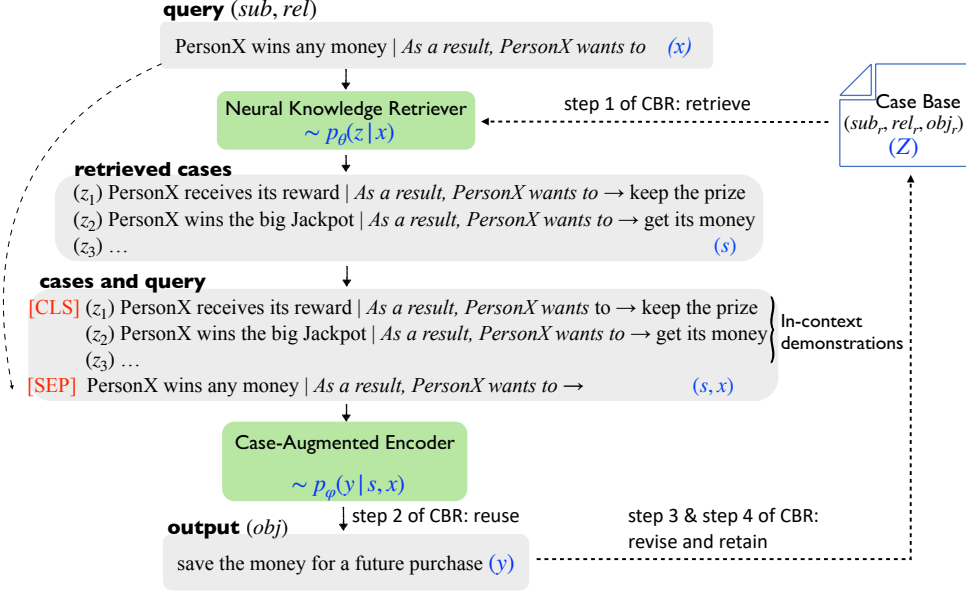


Figure 1: Our end-to-end case-based reasoning framework (ECBRF) for commonsense knowledge base completion. It involves all four steps of the CBR methodology (*retrieve, reuse, revise and retain*).

sub and rel , and y consists of obj . More specifically, ECBRF decomposes $p(y|x)$ into two steps: *retrieve* and *predict*. Given an input x , we first retrieve similar cases z_1, z_2, \dots (each case z_i consists of sub_r, rel_r and obj_r) from case base Z , while $(x, y) \notin z_i$. We model this as a sample from the distribution $p(z|x)$.

Then we use z_i in a number of m to compose a supporting set s (for each query, one supporting set is used for input augmentation). Specifically,

$$\hat{p}(s|x) = \sum_{z_i \in \text{top-}m(p(\cdot|x))} p(z_i|x) \quad ((x, y) \notin z) \quad (1)$$

$$p(s_i|x) = \frac{\exp \hat{p}(s_i|x)}{\sum_j \exp \hat{p}(s_j|x)} \quad (2)$$

Then we condition on both the supporting set s and the query x to generate the output y , modeled as $p(y|s, x)$. To obtain the overall likelihood of generating y , we treat s as a latent variable and marginalize s via a top- k approximation, yielding:

$$p(y|x) = \sum_{s \in \text{top-}k(p_\theta(\cdot|x))} p_\theta(s|x) p_\phi(y|s, x) \quad (3)$$

4.2 Model Architecture

We now detail two key components – the *neural knowledge retriever*, which models $p_\theta(z|x)$; and *case-augmented encoder*, which models $p_\phi(y|s, x)$.

Neural Knowledge Retriever The retriever uses max inner product search (MIPS) to retrieve z . Specifically, the retriever is defined using a dense inner product model:

$$p_\theta(z|x) = \frac{\exp f(x, z)}{\sum_{z'} \exp f(x, z')} \quad (4)$$

$$f(x, z) = \text{Embed}_{query}(x)^T \text{Embed}_{case}(z) \quad (5)$$

where Embed_{query} is an embedding function that maps sub and rel in the query input to a d -dimensional vector, and Embed_{case} is an embedding function that maps sub_r, rel_r and obj_r in the knowledge tuples in memory store to a d -dimensional vector. The relevance score $f(x, z)$ between x and z is defined as the inner product of the vector embeddings. The retrieval distribution is the softmax over relevance scores between top- k retrieved cases and current query input.

We implement the embedding functions Embed_{query} and Embed_{case} using two DPR-based models (Karpukhin et al., 2020). The input format for query x is the concatenation of subject and relation: [CLS] sub [SEP] rel [SEP]; And the input format for case z is the concatenation of the subject, relation, and object: [CLS] sub_r [SEP] rel_r obj_r [SEP].

Case-Augmented Encoder Given an input x and a supporting set s , the case-augmented encoder

defines:

$$p_\varphi(y|s, x) = \prod_i^N p_\varphi(y_i|x, s, y_{1:i-1}) \quad (6)$$

We use BART (Lewis et al., 2020a) and GPT-2 (Radford et al., 2019) as the base model for case-augmented encoder.

We also add prompts which we find is helpful. The input format (with prompts and in-context demonstrations) for case-augmented encoder is:

*Here are some similar cases to infer from: $z_0 z_1$
... z_{m-1} From the similar cases we can infer that:
[SEP] sub rel*

Zhao et al. (2021) show that pre-trained language model has “Recency Bias”, which is the tendency to repeat answers that appear in the last in-context demonstration in classification tasks. We analyse this strategy for the generative CKBC task (we call it “reverse demonstration”) that the most similar case from the retriever is placed as the last demonstration, the second most similar case in the second last demonstration, and so on.

4.3 Training Method

Since the purpose of in-context demonstration is only to provide ancillary information, the model should be able to predict *obj* w/ or w/o it. Therefore here we design a specific training strategy for ECBRF – during training, we *randomly mask* out in-context demonstrations and only keep the (*sub, rel*) query for some training examples with probability p_{mask} . It is designed to function similarly to dropout to prevent overly relying on retrieved cases.

5 Experiments & Analysis

In this section, we introduce the experiment datasets and evaluation details, as well as experiment setups and the experiment results, measured with automatic and human evaluations.

5.1 Datasets and statistics

We evaluate ECBRF using two automatic common-sense knowledge base completion benchmarks — ConceptNet (Speer et al., 2017) and ATOMIC (Sap et al., 2019). In total, ConceptNet contains 101,800 tuples and ATOMIC contains 877,077 tuples. We use the same data split as COMET (Bosselut et al., 2019) did. In ATOMIC, around 17% of the labeled knowledge tuples use “None” as the object. As a result, models can easily get high performance by

always generating “None”. To better evaluate, we don’t use knowledge tuples with the object being “None” for both training and evaluation. Apart from using the entire train set for training, we also conduct experiments in the few-shot settings — where the model is only trained with 5 to 320 knowledge tuples³. More details on data pre-processing is shown in §A.7.

5.2 Evaluation Details

For automatic evaluation metrics, following Bosselut et al. (2019) we use BLEU-2, perplexity, and *novelty* metrics (including %N/T-sro, %N/T-o, and %N/U-o). Specifically for *novelty* metrics, we report the proportion of all generated tuples that are novel tuple (%N/T-sro) (here novel means unseen in train set), have a novel *obj* (%N/T-o), and the proportion of the set of unique *obj* in all generated objects (%N/U-o).

In addition to automatic evaluation, we also perform human evaluation, including *validness*, *informativeness*, and *preference score*. For *validness* and *informativeness*, following Gabriel et al. (2021), the score is based on a 5-point Likert scale (with 5 points the highest score). For *validness*, following Gabriel et al. (2021), we judge the validness of the generated new knowledge by the likelihood of inferences based on a 5-point Likert scale (with 5 points the highest score). Specifically, obviously true (5), generally true (4), plausible (3), neutral or unclear or basically a repetition (sub-sentence) of the query (2), and doesn’t make sense (1). For *informativeness*, the rating standard is also based on a 5-point Likert scale. Specifically, rich in relevant details (5), has relevant details (4), it seems some details are provided (3), basically a repetition (sub-sentence) of the query (2), unfinished generation (1). For *preference score*, We ask the human raters to *compare* the generations between ECBRF and COMET. Specifically, a valid generation with more information provided will be assigned 1.0 point, and a generation that is not valid or with less information will be assigned 0.0 instead. However, if the two generations perform comparably, both generations will be assigned 0.5 points.

Following Bosselut et al. (2019), for each experiment and for each model, we sample 100 generations for human evaluation. Each generation is rated by three graduate students. During the

³Note that for the few-shot settings, our ECBRF’s case base is also initialized with 5 to 320 tuples

ConceptNet	5-shot	20-shot	40-shot	160-shot	320-shot	Full (100%)
COMET (GPT2)*	374.32 / 0.33	339.19 / 0.58	282.76 / 0.38	58.59 / 1.44	41.23 / 2.29	13.04 / 3.37
ECBRF (GPT2)	284.67 / 0.51	220.07 / 0.59	102.46 / 1.63	52.10 / 1.61	38.63 / 2.59	13.60 / 2.83
COMET (BART)*	14.26 / 1.15	11.31 / 1.64	9.48 / 3.70	6.60 / 6.70	5.44 / 9.57	2.90 / 20.19
ECBRF (BART)	12.85 / 1.20	9.53 / 2.11	8.70 / 3.17	6.19 / 6.20	5.22 / 9.31	2.93 / 18.55
w/ <i>random mask</i>	13.68 / 1.40	9.82 / 2.02	8.98 / 2.96	6.14 / 6.88	5.05 / 10.71	2.86 / 19.97
w/o <i>reverse demonstration</i>	13.02 / 1.20	9.60 / 1.84	8.74 / 2.95	6.21 / 6.13	5.06 / 10.12	2.92 / 19.80
w/ rand retrieval	13.23 / 1.00	10.15 / 2.61	9.24 / 3.45	6.42 / 6.21	5.40 / 8.93	2.91 / 20.29

ATOMIC	5-shot	20-shot	40-shot	160-shot	320-shot	Full (100%)
COMET (GPT2)*	753.93 / 2.11	512.11 / 3.44	409.30 / 2.32	209.78 / 2.73	165.28 / 2.68	67.95 / 4.00
ECBRF (GPT2)	653.90 / 2.30	416.16 / 2.89	319.12 / 2.26	182.43 / 2.92	163.56 / 2.86	67.35 / 4.05
COMET (BART)*	19.72 / 5.76	16.83 / 5.38	13.58 / 9.20	14.30 / 11.56	14.45 / 12.67	6.98 / 19.34
ECBRF (BART)	18.17 / 5.16	14.73 / 3.85	13.05 / 10.13	15.19 / 10.13	14.61 / 11.66	6.95 / 19.06
w/ <i>random mask</i>	18.50 / 5.29	14.93 / 4.04	13.01 / 7.47	14.52 / 12.42	14.53 / 12.64	6.96 / 19.22
w/o <i>reverse demonstration</i>	18.13 / 5.58	14.70 / 3.89	12.99 / 5.14	15.15 / 10.18	15.65 / 10.95	6.95 / 19.24
w/ rand retrieval	18.16 / 4.33	14.66 / 3.64	12.96 / 4.70	14.60 / 11.76	14.80 / 12.91	6.96 / 19.02

Table 2: Perplexity (\downarrow) / BLEU (\uparrow) scores on ConceptNet (upper) and ATOMIC (down). The best scores for each setting are boldfaced. *: baseline models (our own implementation).

	20-shot (BART)	160-shot (BART)	Full (BART)
ConceptNet			
COMET	0.37 / 1.76 / 1.74	0.42 / 2.86 / 2.69	0.47 / 3.87 / 3.49
ECBRF	0.63 / 2.47 / 2.38	0.58 / 3.26 / 3.13	0.53 / 3.95 / 3.58
ATOMIC			
COMET	0.43 / 2.21 / 2.27	0.44 / 3.05 / 3.05	0.47 / 3.59 / 3.36
ECBRF	0.57 / 2.44 / 2.58	0.56 / 3.22 / 3.17	0.53 / 3.64 / 3.43

Table 3: Human evaluation results using *preference score*, *validness*, and *informativeness*.

<i>sub</i> : PersonX spends ___ working;
<i>rel</i> : As a result, others feel
Ground truth: ['happy', 'happy to have x in their life']
COMET's generation: happy (BLEU: 31.62)
ECBRF's generation: satisfied with personx's work (BLEU: 0.00)

Table 4: An example to show that BLEU is not a perfect metric for CKBC.

evaluation the order of the two generations to be compared are randomized for each selection, therefore human raters have no clue on which choice is associated with which model. More details about human evaluation can be found at §A.6.

5.3 Experimental Setup

Baselines We use COMET (Bosselut et al., 2019) as our baseline. COMET is originally implemented with GPT (Radford et al.), a pretrained language model as the base model and uses subject and relation as direct input and uses the generation result as object. Here we compare two versions of COMET,

	20-shot (BART)	160-shot (BART)	Full (BART)
ConceptNet			
COMET	98.00 / 13.46 / 58.99	92.22 / 14.83 / 61.83	57.83 / 5.57 / 71.29
ECBRF	96.64 / 21.24 / 51.10	93.27 / 16.09 / 65.83	59.62 / 4.84 / 73.08
ATOMIC			
COMET	100.0 / 58.18 / 15.83	100.0 / 30.92 / 29.11	100.0 / 9.21 / 17.93
ECBRF	100.0 / 82.16 / 22.79	100.0 / 34.96 / 29.70	100.0 / 6.49 / 15.37

Table 5: Novelty evaluation results using %N/T-sro, %N/T-o, and %N/U-o.

one is GPT-2 (Radford et al., 2019) based and another is BART (Lewis et al., 2020a) based. Both GPT-2 and BART are more powerful pretrained language models than GPT. We leave the details of hyperparameters in §A.1.

5.4 Main Results

Automatic evaluation results on BLEU-2 and perplexity are shown in Table 2. We present results with human evaluations in Table 3. Automatic evaluation of *novelty* is shown in Table 5. Table 2 shows that *random mask* is constantly helpful for ECBRF when the train set is equal or larger than 160. Therefore we adopt *random mask* for ECBRF when the train set is equal to or larger than 160 in Table 3 and Table 5.

Table 2 shows that regardless of the selection of base models, ECBRF consistently outperforms the COMET baseline in almost all perplexity measures and most BLEU measures. We argue that BLEU is not a perfect metric for CKBC (Sai et al.,

<i>sub</i> : PersonX wants to play with PersonY
<i>rel</i> : Before, this person needed
One retrieved case by ECBRF: PersonX plays tennis with PersonY’s friend, Before, this person needed, get a tennis racket
COMET’s generation: to have a game ECBRF’s generation: to find a tennis court

Table 6: An example to show how the retrieved cases influence the generated *obj*.

2023), since each *sub* and *rel* pair can lead to more than one feasible *obj*, and BLEU can only refer to a limited set of ground truth *obj*. Even if a model generates a reasonable *obj*, it may yield a low BLEU, because the generated *obj* is not in the ground truth set.

Table 4 shows one typical example from ATOMIC that shows although ECBRF’s generations are reasonable, they only receive low BLEU scores (More examples in §A.2). Table 3 shows that ECBRF consistently outperforms COMET in *preference score*, *validness*, and *informativeness* in human evaluation, especially in few-shot setting. In practice, we observe that in few-shot setting, BART without retrieval tends to repeat the query during generation, while retrieved cases seem to be able to provide knowledge and guidance to generate more proper *obj*.

Table 5 shows that the generated *obj* of ECBRF are generally *novel* especially in few-shot settings. We empirically attribute the lower novelty of ECBRF in full train set to that ECBRF sometimes tends to copy proper retrieved *obj* as generation. The reason %N/T-sro score is always 100.0 in ATOMIC is that the (*sub*, *rel*) pairs in ATOMIC’s train set and test set do not overlap.

We attribute the different conclusions reached on whether in-context demonstrations (ICD) with finetuning can be beneficial to CKBC (Wang et al., 2021) to that (1) ICD is more useful in few-shot settings, so that investigation on full train set setting might not discover this advantage; (2) human evaluation is the most precise metric for the task while BLEU is not so only evaluating with automatic metrics could not be precise; (3) ECBRF uses *random mask*, which is empirically found to be helpful in performance when using a large train set.

5.5 Ablation Study of ECBRF

In Table 2, we show some ablation studies of ECBRF. “w/ rand mask” stands for the ECBRF

model using random mask ($p_{mask} = 30\%$); “w/o reverse demonstrations” stands for the ECBRF model without using reverse demonstrations; and “w/ rand retrieval” represents an ECBRF model that uses randomly searched cases instead of MIPS search.

Both tables for automatic evaluation show that using *random mask* can generally lead to better performance for ECBRF in both perplexity and BLEU when the training set is larger, while lead to worse performance when the number of training set is less than 160. Our interpretation is that, when the train set is very small, the model can benefit from over-relying to the retrieved cases; while when the train set is large, PLM can still benefit from the retrieved cases but the over-relying could be harmful.

The tables also show that “reverse demonstration” only leads to comparable performance, which might indicate that the order of retrieved cases does not make a difference in a generative task (as CKBC).

From the tables, we also observe that ECBRF with MIPS retrieval consistently leads to better performance than ECBRF with random retrieval in terms of perplexity in ConceptNet experiments, while performs comparably with ECBRF in ATOMIC experiments. Notice that the required generation in ConceptNet is usually shorter and more similar compared to ATOMIC. Therefore our interpretation is that, only when the retrieved cases are enough similar to the input query and its designed golden generation, can the retrieved cases significantly benefit the generation process (towards golden generation).

5.6 Qualitative Analysis on How Retrieved Cases Influence *obj* Generation

Table 6 shows one example of the generation of ECBRF and COMET (more examples are shown in §A.3). It shows that ECBRF’s generation is related to the retrieved case, exhibiting the case-based reasoning ability of *reusing the retrieved* old experience to solve new problems.

6 Further Analysis from Perspective of CBR

CBR methodology contains 4 sub-processes, which are *retrieve*, *reuse*, *revise* and *retain*. More specifically, when given a new problem, the method first *retrieves* the most similar cases, then *reuses* the information in that case to solve the new problem by

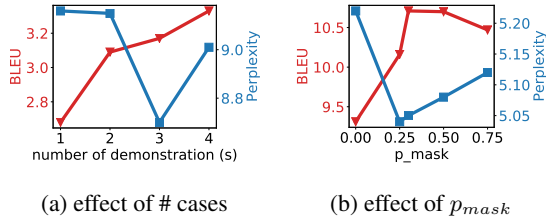


Figure 2: ECBRF’s performance (perplexity (\downarrow) / BLEU (\uparrow)) with regard to different numbers of demonstrations and different p_{mask} . Experiments on the left figure use ConceptNet 40-shot train set, and the right figure use ConceptNet 320-shot train set (since *random mask* is only helpful on large train set).

proposing a new solution, then *revises* the proposed solution according to the feedback of adopting it in real application scenarios (*revise* step usually involves human’s effort), and finally select high quality revised solutions together with their problems as new cases to *retain* to case base.

We provide an analysis of how the high-level methodology of CBR (*retrieve*, *reuse*, *revise* and *retain*) shapes the design and how the selection details of CBR-related components improve the performance of our end-to-end DL framework.

Step 1: Retrieve *Retrieve* is an important step since the effectiveness of a CBR system largely relies on its ability to retrieve useful previous cases (Montazemi and Gupta, 1997). Here we use *neural knowledge retriever* (DPR) for retrieving the most similar cases. Table 2 shows the results of ECBRF using MIPS retrieval and random retrieval. As illustrated in §5.5, from the experimental results we hypothesize that ECBRF tends to make generations that are similar to the retrieved cases. This hypothesis is consistent with insights from CBR that the *retrieve* step is essential for guiding the *reuse* step. However, the difference lies in that CBR insights rely on *retrieve* step more (with irrelevant retrieval it would be particularly hard for *reuse*), while PLM seems to be able to benefit from even random retrieval.

Step 2: Reuse Here we use *case-augmented encoder* to automatically *reuse* the retrieved cases.

Figure 2a shows the effects of number of retrieved cases. We observe that when *case-augmented encoder* uses 3 cases, it reaches the best perplexity, and nearly the best BLEU performance.

Figure 2b shows the effects of p_{mask} . Only when p_{mask} is 1.0, in-context demonstrations are not used at test time, which makes the model the same

	Perplexity	BLEU
ECBRF	8.70	3.17
w/ only obj_r	8.90	3.26
w/o prompt	9.01	3.44
w/ larger case base	7.90	3.69

Table 7: Ablation Study: effect of sub_r , prompt, and the retain step (perplexity (\downarrow) / BLEU (\uparrow)). Results of this table use ConceptNet 40-shot train set.

as COMET. As we gradually increase p_{mask} , perplexity keeps improving and BLEU-2 reaches the global maximum when p_{mask} is 0.3. It is also interesting to see empirically how the *case-augmented encoder* gradually learns to *reuse* the retrieved cases to increase the performance of the deep learning model as we gradually decrease p_{mask} .

In CBR, the *reuse* of the retrieved case’s solution contains two steps: (a) find the difference between the past and the current queries and (b) adapt the retrieved solution to the current query (Aamodt and Plaza, 1994). So it is important to know the difference between the past queries and the current input query for better adaptation. Table 7 shows the comparison between the result of only using obj_r (the retrieved cases’ object phrases) as in-context demonstrations and the result of ECBRF (uses both sub_r , rel_r and obj_r), and we observe that ECBRF performs better in perplexity. but a little bit worse in BLEU. This result indicates that deep learning based *case-augmented encoder* is possible to automatically learn and reason from the difference between the past queries and the current input query for *reuse*. We leave further investigations on whether PLMs can learn to compare the difference between the past queries and the current input query as an open research question.

We use prompts to indicate the role of retrieved cases and current query in input. Table 7 shows that *case-augmented encoder* with the prompt performs better in perplexity while a bit worse on BLEU, indicating that usage of a prompt is possible to help the model better *reuse* the retrieved cases.

Step 3 & 4: Revise and Retain Since *revise* typically involves human efforts, here we simulate *revise* and *retain* and see their effect on our framework. The result of *revise* and *retain* is a larger case base with more high quality data, and the parameters of the model for *reuse* are not necessarily updated according to the new data. Here we simulate the effect of *revise* and *retain* by first training

ECBRF in a low-resource experiment (with a small case base), then at test time we expand the case base to the full train set. Table 7 shows that, at test time, ECBRF with access to a larger case base substantially outperforms ECBRF (with access to only a small case base), although the parameters have not been updated with the new data. This result demonstrates that our framework can benefit from CBR’s methodology as *revise* and *retain*.

7 Conclusion

Drawing insights from CBR, we propose an end-to-end framework for the CKBC task. We demonstrate through automatic and human evaluations that our framework generates more valid knowledge than the state-of-the-art COMET model in both the fully supervised and few-shot settings. From the perspective of CBR, our framework addresses a fundamental question on whether CBR methodology can be utilized to improve deep learning models.

8 Future Works and Challenges

In general, we hope this work could provide some insights to bridge the two research areas, classic AI (Case-Based Reasoning) and deep learning based NLP methods together, and therefore to advance the research of both fields from each other’s research developments.

From the aspect of NLP methods, for example, new prompting methods could be further developed based on insights from CBR research; The concept of *revise* and *retain* from CBR could be paid more attention to investigate their interaction with in-context demonstrations (prompting).

From the aspect of CBR, this work provides a tentative answer to the two long-remaining challenges — (1) whether CBR can be used to complement DL (Leake and Crandall, 2020), given that the latest work even suggests that in many tasks NN itself outperforms CBR-complemented NN (Ye et al., 2022); (2) the adaptation (*reuse* step) of previous cases to the current case is a very challenging problem, so that in many fields the CBR methodology is used only as a retriever (Choudhury and Begum, 2016). How to further answer these two questions could be a challenging research topic.

9 Limitations

From the perspective of CBR, we have shown through experiments that our framework can per-

form *retrieve* and *reuse* steps, and can benefit from *revise* and *retain* steps. But the *revise* step in CBR typically involves human efforts, and this paper does not focus on addressing this challenge. As a result, our framework might still need manual efforts to benefit from *revise* and *retain*.

However, human efforts could be more efficiently utilized for *revise* than writing new data from scratch. Since comparing with requesting the workers to write the knowledge from scratch, *revising* the existing generations of ECBRF could be much faster.

10 Acknowledgements

We thank Zefan Zhang and Jie Zheng for their participation and provision of human evaluation to this research work.

References

- A. Aamodt and E. Plaza. 1994. Case-based reasoning: foundational issues, methodological variations, and system approaches. *AI communications*.
- Ian Apperly. 2010. *Mindreaders: the cognitive basis of "theory of mind"*. Psychology Press.
- Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, and Yejin Choi. 2019. COMET: Commonsense transformers for automatic knowledge graph construction. In *ACL 2019*. Association for Computational Linguistics.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Erik Cambria, Qian Liu, Sergio Decherchi, Frank Xing, and Kenneth Kwok. 2022. [Senticnet 7: A commonsense-based neurosymbolic AI framework for explainable sentiment analysis](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference, LREC 2022, Marseille, France, 20-25 June 2022*, pages 3829–3839. European Language Resources Association.
- Nabanita Choudhury and Shahin Ara Begum. 2016. A survey on case-based reasoning in medicine. *International Journal of Advanced Computer Science and Applications*, 7(8).
- Peter Clark, Oyvind Tafjord, and Kyle Richardson. 2020. [Transformers as soft reasoners over language](#). In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020*, pages 3882–3890. ijcai.org.

- Jeff Da, Ronan Le Bras, Ximing Lu, Yejin Choi, and Antoine Bosselut. 2021. Understanding few-shot commonsense knowledge models. *AKBC*.
- Rajarshi Das, Ameya Godbole, Shehzaad Dhuliawala, Manzil Zaheer, and Andrew McCallum. 2020a. A simple approach to case-based reasoning in knowledge bases. In *Automated Knowledge Base Construction*.
- Rajarshi Das, Ameya Godbole, Nicholas Monath, Manzil Zaheer, and Andrew McCallum. 2020b. Probabilistic case-based reasoning for open-world knowledge graph completion. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4752–4765.
- Rajarshi Das, Manzil Zaheer, Dung Thai, Ameya Godbole, Ethan Perez, Jay Yoon Lee, Lizhen Tan, Lazaros Polymenakos, and Andrew McCallum. 2021. Case-based reasoning for natural language queries over knowledge bases. In *EMNLP 2021*, pages 9594–9611.
- Ernest Davis and Gary Marcus. 2015. Commonsense reasoning and commonsense knowledge in artificial intelligence. *Communications of the ACM*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Saadia Gabriel, Chandra Bhagavatula, Vered Shwartz, Ronan Le Bras, Maxwell Forbes, and Yejin Choi. 2021. Paragraph-level commonsense transformers with recurrent memory. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 12857–12865.
- Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. Making pre-trained language models better few-shot learners. In *ACL 2021*, pages 3816–3830, Online. Association for Computational Linguistics.
- Jonathan Gordon and Benjamin Van Durme. 2013. Reporting bias and knowledge acquisition. In *Proceedings of the 2013 workshop on Automated knowledge base construction*, pages 25–30.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and M-W Chang. 2020. Realm: Retrieval-augmented language model pre-training. *ICML*.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*.
- Shaoxiong Ji, Shirui Pan, Erik Cambria, Pekka Martinen, and Philip S. Yu. 2022. A survey on knowledge graphs: Representation, acquisition, and applications. *IEEE Trans. Neural Networks Learn. Syst.*, 33(2):494–514.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. In *EMNLP 2020*, pages 6769–6781, Online. Association for Computational Linguistics.
- Mark T. Keane, Eoin M. Kenny, Mohammed Temraz, Derek Greene, and Barry Smyth. 2021. Twin systems for deepcbr: A menagerie of deep learning and case-based reasoning pairings for explanation and data augmentation. *CoRR*, abs/2104.14461.
- Janet L Kolodner. 1997. Educational implications of analogy: A view from case-based reasoning. *American psychologist*, 52(1):57.
- David Leake and David Crandall. 2020. On bringing case-based reasoning methodology to deep learning. In *International Conference on Case-Based Reasoning*. Springer.
- David Leake, Xiaomeng Ye, and David J. Crandall. 2021. Supporting case-based reasoning with neural networks: An illustration for case adaptation. In *Proceedings of the AAAI 2021 Spring Symposium on Combining Machine Learning and Knowledge Engineering (AAAI-MAKE 2021)*, Stanford University, Palo Alto, California, USA, March 22-24, 2021, volume 2846 of *CEUR Workshop Proceedings*. CEUR-WS.org.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020a. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *ACL 2020*, pages 7871–7880, Online. Association for Computational Linguistics.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020b. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Xiang Li, Aynaz Taheri, Lifu Tu, and Kevin Gimpel. 2016. Commonsense knowledge base completion. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1445–1455, Berlin, Germany. Association for Computational Linguistics.
- Chieh-Kang Liao, Alan Liu, and Yu-Sheng Chao. 2018. A machine learning approach to case adaptation. In *First IEEE International Conference on Artificial Intelligence and Knowledge Engineering, AIKE 2018, Laguna Hills, CA, USA, September 26-28, 2018*, pages 106–109. IEEE Computer Society.
- Chaitanya Malaviya, Chandra Bhagavatula, Antoine Bosselut, and Yejin Choi. 2020. Commonsense knowledge base completion with structural and semantic context. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The*

- Thirty-Second Innovative Applications 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 2925–2933. AAAI Press.
- Rui Mao, Chenghua Lin, and Frank Guerin. 2018. Word embedding and WordNet based metaphor identification and interpretation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, volume 1, pages 1222–1231.
- Ali Reza Montazemi and Kalyan Moy Gupta. 1997. A framework for retrieval in cbr systems. *Annals of operations research*.
- Feng Pan, Rutu Mulkar-Mehta, and Jerry R. Hobbs. 2011. [Annotating and learning event durations in text](#). *Comput. Linguistics*, 37(4):727–752.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI blog*.
- Ananya B. Sai, Akash Kumar Mohankumar, and Mitesh M. Khapra. 2023. [A survey of evaluation metrics used for NLG systems](#). *ACM Comput. Surv.*, 55(2):26:1–26:39.
- Itsumi Saito, Kyosuke Nishida, Hisako Asano, and Junji Tomita. 2018. [Commonsense knowledge base completion and generation](#). In *Proceedings of the 22nd Conference on Computational Natural Language Learning*, pages 141–150, Brussels, Belgium. Association for Computational Linguistics.
- Merrilee H Salmon. 1989. Introduction to logic and critical thinking.
- Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A Smith, and Yejin Choi. 2019. Atomic: An atlas of machine commonsense for if-then reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3027–3035.
- Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Thirty-first AAAI conference on artificial intelligence*.
- Nitish Srivastava, Geoffrey E. Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. [Dropout: a simple way to prevent neural networks from overfitting](#). *J. Mach. Learn. Res.*, 15(1):1929–1958.
- Peifeng Wang, Filip Ilievski, Muhao Chen, and Xiang Ren. 2021. [Do language models perform generalizable commonsense inference?](#) In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3681–3688, Online. Association for Computational Linguistics.
- Quan Wang, Zhendong Mao, Bin Wang, and Li Guo. 2017. Knowledge graph embedding: A survey of approaches and applications. *IEEE Transactions on Knowledge and Data Engineering*, 29(12):2724–2743.
- Ian Watson. 1999. Case-based reasoning is a methodology not a technology. In *Research and Development in Expert Systems XV*. Springer.
- Zonglin Yang, Li Dong, Xinya Du, Hao Cheng, Erik Cambria, Xiaodong Liu, Jianfeng Gao, and Furu Wei. 2022. [Language models as inductive reasoners](#). *CoRR*, abs/2212.10923.
- Zonglin Yang, Xinya Du, Alexander M. Rush, and Claire Cardie. 2020. [Improving event duration prediction via time-aware pre-training](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020*, volume EMNLP 2020 of *Findings of ACL*, pages 3370–3378. Association for Computational Linguistics.
- Liang Yao, Chengsheng Mao, and Yuan Luo. 2019. [KG-BERT: BERT for knowledge graph completion](#). *CoRR*, abs/1909.03193.
- Xiaomeng Ye, David Leake, and David J. Crandall. 2022. [Case adaptation with neural networks: Capabilities and limitations](#). In *Case-Based Reasoning Research and Development - 30th International Conference, ICCBR 2022, Nancy, France, September 12-15, 2022, Proceedings*, volume 13405 of *Lecture Notes in Computer Science*, pages 143–158. Springer.
- Xiaomeng Ye, David Leake, Vahid Jalali, and David J. Crandall. 2021. [Learning adaptations for case-based classification: A neural network approach](#). In *Case-Based Reasoning Research and Development - 29th International Conference, ICCBR 2021, Salamanca, Spain, September 13-16, 2021, Proceedings*, volume 12877 of *Lecture Notes in Computer Science*, pages 279–293. Springer.
- Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models. In *International Conference on Machine Learning*, pages 12697–12706. PMLR.

A Appendix

A.1 Hyperparameters

We use a common training hyperparameters setup for PLM, and the variance of hyperparameters such as batch size is mainly determined by the computational resources.

Specifically, batch size is 32 for all ATOMIC experiments; is 16 for all BART-based ConceptNet

experiments; is 8 for all GPT2-based ConceptNet experiments. Learning rate is $1e-5$ for all experiments. The remainder of our training hyperparameters is the same as COMET (Bosselut et al., 2019) for full train set experiments. For few-shot experiments, we adapt the warm up steps accordingly (more details can be found in our public release code).

For decoding hyperparameters, we use greedy decoding for all BART experiments (since Bosselut et al. (2019) suggests that greedy decoding can lead to the best human evaluation results); and use topk ($k=50$) decoding for all GPT2 experiments (since with greedy decoding many generations of GPT2 model is only the end of text token; with randomly sampling a concrete generation could be made via multiple attempts).

A.2 Examples that BLEU not a Perfect Metric

Table 8 shows three examples with three different *rel* from ATOMIC that shows although ECBRF’s generations are reasonable, but they only receive low BLEU scores.

This table shows that sometimes a good generation is assigned with a low BLEU score, especially when the generation is novel and unseen from the ground truth set. Table 5 shows that ECBRF produces more novel generations compared with the COMET baseline, which might make ECBRF suffer more from the imperfectness of the BLEU score.

A.3 Examples on How Retrieved Cases Influence ECBRF

Table 9 shows three examples of the generation of ECBRF and COMET.

This table shows that many ECBRF’s generations are related to the retrieved case, exhibiting the case-based reasoning ability of *reusing* the *retrieved* old experience to solve new problems.

A.4 More Related Works on CBR

Leake and Crandall (2020) advocate using CBR to complement the challenges in deep learning (e.g., few-shot learning). However, past works on using CBR to complement DL only limit to shallow Neural Networks (NN) (Liao et al., 2018; Leake et al., 2021; Ye et al., 2021, 2022). The latest work even suggests that in many tasks NN itself outperforms CBR-complemented NN (Ye et al., 2022), which raises fundamental questions on whether CBR methodology is useful for DL. Other relevant

works include only use CBR to improve the explainability of deep learning models (Keane et al., 2021), or use CBR to improve the performance of symbolic reasoning (Das et al., 2020a,b).

A.5 Future Applications of ECBRF

Although experiments in this paper are only conducted under the commonsense knowledge base completion task, ECBRF could potentially be utilized by other commonsense reasoning tasks, such as temporal commonsense reasoning (Yang et al., 2020), commonsense based sentiment analysis (Cambria et al., 2022), and metaphor processing (Mao et al., 2018) (similar to CBR, metaphor is also highly related to analogical reasoning).

A.6 Human Evaluation Details

Inter-annotator agreement Since we have three annotators for human evaluation, and the spearman correlation is to compare two rank lists, we calculate the averaged spearman correlation for each metric. Specifically, the average spearman correlation for “validness” is 0.71, and the average spearman correlation for “informativeness” is 0.66; For the “preference score”, since it’s to compare each ECBRF’s generation with one baseline (COMET) generation, we can’t rank generations according to the value of the preference score. So following Pan et al. (2011), we use Cohen’s kappa coefficient instead. The kappa coefficient for “preference score” is 0.81.

Instructions for human evaluators There are 3 evaluation metrics: preference score, validness score and informativeness score.

Preference score: a valid generation with more information provided will be assigned 1.0 point, and a generation that is not valid or with less information with being assigned 0.0 instead. However, if the two generations perform comparably, both generations will be assigned 0.5 points.

Validness score: Validness score is given in a 5-point scale (1 5) obviously true (5) generally true (4) plausible (3) neutral or unclear or basically a repetition (sub-sentence) of the query (2) doesn’t make sense (1). Examples for validness: (5): PersonX wants to learn piano, others want, to teach personX; (4): PersonX wants to learn piano, others want, to hear peronX play; (3): PersonX wants to learn piano, others want, to sell piano to peronX; (2): PersonX wants to learn piano, others want, to turn their music on; (1): PersonX wants to learn

Example 1	<i>sub</i> : PersonX spends ___ working <i>rel</i> : As a result, others feel ground truth: ['happy', 'happy to have x in their life'] COMET's generation: happy (BLEU score: 31.62) ECBRF's generation: satisfied with personx's work (BLEU score: 0.00)
Example 2	<i>sub</i> : PersonX expects another ___ <i>rel</i> : This person then ground truth: ['prepares themselves', 'gains knowledge'] COMET's generation: gains knowledge (BLEU score: 100.00) ECBRF's generation: wants to find out what it is that they are going to do next (BLEU score: 0.00)
Example 3	<i>sub</i> : PersonX spends the ___ working <i>rel</i> : This person is seen as ground truth: ['diligent', 'tired', 'hardworking'] COMET's generation: hardworking (BLEU score: 31.62) ECBRF's generation: dedicated (BLEU score: 0.00)

Table 8: Examples to show that BLEU is not a perfect metric for CKBC. This table shows that sometimes a good generation is assigned with a low BLEU score, especially when the generation is novel and unseen from the ground truth set. Table 5 shows that ECBRF produces more novel generations compared with the COMET baseline, which might make ECBRF suffer more from the imperfectness of the BLEU score.

Example 1	<i>sub</i> : PersonX wants to play with PersonY <i>rel</i> : Before, this person needed COMET's generation: to have a game Retrieved cases by ECBRF include: (PersonX plays tennis with PersonY's friend, Before, this person needed, get a tennis racket) ECBRF's generation: to find a tennis court
Example 2	<i>sub</i> : PersonX advances another ___ <i>rel</i> : This person then COMET's generation: PersonX gains knowledge Retrieved cases by ECBRF include: (PersonX marries the king's ___, This person then, he becomes king) ECBRF's generation: becomes more powerful
Example 3	<i>sub</i> : PersonX smiles broadly <i>rel</i> : As a result, this person wants to COMET's generation: to smile back Retrieved cases by ECBRF include: (PersonX grins like a cheshire cat, As a result, this person wants to, express their feelings) ECBRF's generation: to tell others about the good time they had

Table 9: Examples to show how the retrieved cases influence the generated *obj*. This table shows that many ECBRF's generations are related to the retrieved case, exhibiting the case-based reasoning ability of *reusing* the *retrieved* old experience to solve new problems.

piano, others want, to wait for the dinner.

Informativeness score: Informativeness score is given in a 5-point scale (1-5) rich in relevant details (5) has relevant details (4) it seems some details is provided (3) not related information or basically a repetition (sub-sentence) of the query (2) unfinished generation or doesn't make sense (1) Examples for informativeness: (5): PersonX wants to learn piano, others want, to teach personX the basic usage of piano and how to buy a suitable piano; (4): PersonX wants to learn piano, others want, to teach personX on how to use piano; (3): PersonX wants to learn piano, others want, to teach

personX; (2): PersonX wants to learn piano, others want, to turn their music on; (1): PersonX wants to learn piano, others want, to teach.

A.7 Other Details on Data Pre-processing

We observe that in Table 2, our GPT experiment results are lower than Bosselut et al. (2019). We attribute the reason to our different data pre-processing method — we filter all the “.” in the *obj* in ATOMIC and ConceptNet datasets since we observe that only a part of *obj* are equipped with “.”, which might confuse the model on whether to generate “.” or not.

Without “.”, a fixed token in generation, it would be harder for a model to reach higher perplexity and BLEU. In our preliminary experiments that do not especially filter “.”, our re-implemented COMET’s results are comparable to [Bosselut et al. \(2019\)](#).

A.8 Other Details on Retrieval

For ECBRF, we empirically find that not retrieving cases that share the same sub_r with sub is beneficial for the performance. The intuition behind this mechanism is that when one of the retrieved cases has the same sub_r with sub , then it should be natural to copy the corresponding obj_r as generation. However, it is common that obj is different from obj_r , which might confuse the model.