

Is Knowledge Embedding Fully Exploited in Language Understanding? An Empirical Study

Anonymous ACL submission

Abstract

The recent development of knowledge embedding (KE) enables machines to represent knowledge graphs (KGs) with low-dimensional embeddings, which facilitates utilizing KGs for various downstream natural language understanding (NLU) tasks. However, less work has been done on systematically evaluating the impact of KE on NLU. In this work, we conduct a comprehensive analysis of utilizing KE on four downstream knowledge-driven NLU tasks using two representative knowledge-guided frameworks, including knowledge augmentation and knowledge attention. From the experimental results, we find that: (1) KE models that have better performance on knowledge graph completion do not necessarily help knowledge-driven NLU tasks better in the knowledge-guided frameworks; (2) KE could effectively benefit NLU tasks from two aspects including entity similarity and entity relation information; (3) KE could further benefit pre-trained language models which have already learned rich knowledge from pre-training. We hope the results could help and guide future studies to utilize KE in NLU tasks. Our source code will be released to support further exploration.

1 Introduction

Knowledge graphs (KGs) organize entity knowledge and concept knowledge into structured relational data, potentially providing rich information for a variety of NLP tasks, such as information retrieval (Hu et al., 2009), information extraction (Hoffmann et al., 2011), and question answering (Bordes et al., 2014a,c). Both the research community and the industry have built various large-scale KGs¹ and intend to exploit the rich information in KGs to help natural language understanding.

KG is a typical kind of non-Euclidean data, which is difficult for deep learning models to use di-

¹E.g., YAGO (Suchanek et al., 2007), Freebase (Bollacker et al., 2008), and Wikidata (Vrandečić and Krötzsch, 2014).

rectly (Bronstein et al., 2017), while deep learning has become the standard technique of NLP. Knowledge embedding (KE) represents entities and relations in KGs as low-dimensional semantic embeddings in a Euclidean space, which clears the way for injecting KGs into deep learning models. Recently, many efforts have been devoted to KE (Minervini et al., 2017; Guo et al., 2018; Padia et al., 2019) and KE has shown its strong ability to represent knowledge. Hence, it is feasible to integrate KE in downstream NLP tasks.

Although some recent work has explored utilizing KE for NLP, these studies usually only focus on a single task with a single KE (Weston et al., 2013; Bordes et al., 2014a; Zhang et al., 2016; Xin et al., 2018). Less work has been done to systematically evaluate the impact of KE on NLP. To advance the utilization of KE, we need to understand how and to what extent KE contributes to downstream NLP tasks.

In this paper, we focus on the impact of KE on language understanding. First, we summarize two mainstream knowledge-guided frameworks based on existing work: knowledge augmentation and knowledge attention. Then, we perform a comprehensive analysis of utilizing KE on four knowledge-driven NLU tasks. Specifically, we evaluate these frameworks on the following two types of tasks: (1) **Entity-oriented tasks**: relation extraction and entity typing; (2) **General NLU tasks**: information retrieval and fact verification. Besides, to investigate KE’s effect with various text encoders, we implement three representative text encoders: CNN (Kim, 2014), LSTM (Hochreiter and Schmidhuber, 1997), and BERT (Devlin et al., 2019). We have the following observations:

(1) In most cases, KE models can improve the performance of the models that only use texts. However, a KE model with better performance on KGC does not necessarily better help NLU tasks.

(2) For what information of KE could help lan-

082 guage understanding, our experiments show that
083 the models using entity embeddings as external
084 knowledge could effectively capture entity similar-
085 ity and entity relation information.

086 (3) Although previous work (Petroni et al., 2019)
087 has revealed that pre-trained language models
088 (PLMs) such as BERT could learn rich factual
089 knowledge from the pre-training on large-scale cor-
090 pora, our experiments indicate that KE is still valu-
091 able for enhancing PLMs, and how to design a
092 feasible way to combine KE and PLMs remains an
093 exciting research direction.

094 Hopefully, the results of our analysis would pro-
095 vide some insights about how to better utilize KE
096 for language understanding in the future.

097 2 Background

098 2.1 Knowledge Embedding

099 In this subsection, we first introduce several repre-
100 sentative KE models and then summarize the KE
101 models chosen in our experiments.

102 **Linear Models** utilize a linear combination of
103 the relation embedding and head/tail entity em-
104 beddings to model the probability of the relational
105 fact (Bordes et al., 2011, 2012, 2014b). LFM (Je-
106 natton et al., 2012; Sutskever et al., 2009) is a rep-
107 resentative linear model, which employs a relation-
108 specific bilinear form to consider the relatedness
109 between entities and relations. DistMult (Yang
110 et al., 2014) further reduces the number of relation
111 parameters in LFM via simply restricting relation
112 matrices to be diagonal matrices, resulting in a less
113 complicated model and better performance.

114 **Translation Models** regard the relation embed-
115 ding as a translation between the head and tail en-
116 tities’ embeddings. Bordes et al. (2013) propose
117 the first translation model TransE, which is simple
118 but effective. Although TransE achieves promising
119 results, it cannot handle the complex relations in
120 KGs well due to its simple structure. Various trans-
121 lation models have been proposed to address this
122 issue, such as TransH (Wang et al., 2014), Tran-
123 sR/CTransR (Lin et al., 2015), TransD (Ji et al.,
124 2015), TransSparse (Ji et al., 2016), KG2E (He et al.,
125 2015), and ManifoldE (Xiao et al., 2016).

126 **Neural Models** utilize neural networks to model
127 the probability of the relational fact by taking the
128 head/tail entity and relation embeddings as inputs.
129 NTN (Socher et al., 2013) employs a bilinear ten-
130 sor to combine two entities’ embeddings via multi-
131 ple aspects. Moreover, HoIE (Nickel et al., 2016)

132 uses the circular correlation of vectors to repre-
133 sent pairs of entities, which could combine the
134 expressive power of the tensor product with the
135 efficiency and simplicity of TransE. In other work,
136 both NAM (Liu et al., 2016) and ConvE (Dettmers
137 et al., 2018) utilize multi-layer networks to capture
138 the interactions among entities and relations. We
139 find that most of neural models are designed for the
140 small-scale KG like FB15K-237. We have tried to
141 conduct experiments on neural KE and find that the
142 GPU memory cannot place our large-scale KG. We
143 think the feasibility for large-scale KGs is impor-
144 tant for KE algorithms in the application scenario,
145 and thus the neural KE models are not included.

146 **Complex-Valued Models** exploit complex em-
147 beddings to represent the entities and relations.
148 ComplEx (Trouillon et al., 2016) first considers
149 complex embeddings in KE models by employing
150 an eigenvalue decomposition model. Moreover,
151 RotatE (Sun et al., 2019b) defines each relation
152 as a rotation from the head entity to the tail entity
153 in a complex vector space. Benefiting from the
154 strong modeling ability of complex embeddings,
155 complex-valued models achieve quite good perfor-
156 mance compared with other KE models.

157 Notably, some hyperbolic models (Chami et al.,
158 2020; Wang et al., 2020) have been proposed re-
159 cently, which are mainly designed for extremely
160 low-dimensional embeddings and cannot be easily
161 used by neural networks. Hence, we do not choose
162 this kind of models.

163 In this paper, we compare the most typical
164 KE models from each type, including (1) Linear
165 model: DistMult; (2) Translation model: TransE;
166 (3) Complex-valued model: RotatE.

167 2.2 Utilizing External Knowledge for NLP

168 This subsection introduces previous work on utiliz-
169 ing external knowledge in different tasks, including
170 entity-oriented and general NLP tasks.

171 What we call “Entity-oriented tasks” includes
172 most of the information extraction tasks (Chang
173 et al., 2006). These tasks naturally benefit from
174 external knowledge about entities, and thus there
175 are several methods using KGs for these tasks, such
176 as entity typing (Xin et al., 2018; Liu et al., 2019a),
177 and relation extraction (Weston et al., 2013; Han
178 et al., 2018a; Li et al., 2019).

179 There are also several general NLP tasks that do
180 not focus on entities but could effectively benefit
181 from the information of KGs, such as question an-

swering (Bordes et al., 2014a; Miller et al., 2016; Yang and Mitchell, 2017; Huang et al., 2019; Sun et al., 2019a; Verga et al., 2020; Yasunaga et al., 2021), fact verification (Thorne et al., 2018), information retrieval (Xiong et al., 2017; Liu et al., 2018), recommendation systems (Zhang et al., 2016; Wang et al., 2018, 2019a,c; Xian et al., 2019; Wang et al., 2019b; Dhingra et al., 2020), language modeling (Ahn et al., 2016; Gu et al., 2018; Parvez et al., 2018), and dialog systems (He et al., 2017; Ghazvininejad et al., 2018). In general NLP tasks, KGs can provide external background knowledge to understand the context, such as in question answering and fact verification; or serve as external interactions between two texts for similarity measuring, such as for information retrieval and recommendation systems.

3 Knowledge-Guided Frameworks

We derive two general knowledge-guided frameworks based on previous work mentioned in the last section: knowledge augmentation framework and knowledge attention framework.

In the scenario of knowledge-driven NLU, we aim to obtain the representations of the word sequence O_w and entity sequence O_e and fuse them for prediction, where the entity sequence consists of entities appearing in text. The token sequence is denoted by $\{w_i\}_{i=1}^n$, where n is the sequence length. Meanwhile, the entity sequence is denoted by $\{e_i\}_{i=1}^m$, where m is the number of entities.

For classification tasks, such as relation classification, the representations of O_w and O_e are two vectors summarizing all information in the sequence. To use them, we input the concatenation of these two vectors into a multi-layer perceptron (MLP) to predict labels. For matching tasks, such as information retrieval, the representations of O_w and O_e are two sequences of embeddings for each word or entity. To use them, we follow the kernel method proposed by Dai et al. (2018).

Knowledge Augmentation Framework aims to directly integrate entity knowledge by treating entity sequences as external features. It could be generalized to a variety of existing work utilizing KE (Weston et al., 2013; Han et al., 2018a; Xiong et al., 2017; Liu et al., 2018). This framework formulates the entity representation O_e as

$$O_e = \text{Enc}_e(e_1, \dots, e_m), \quad (1)$$

where Enc_e is the entity encoder, which is usually an MLP, and e_i is the entity embedding of e_i . For

classification tasks, Enc_e takes the concatenation of entity embeddings as input. For matching tasks, Enc_e is applied to each entity embedding.

Knowledge Attention Framework is expected to capture semantic correlations of context and entity knowledge. It utilizes entity information to gather different aspects of semantic meanings in the text sequence. This framework is also generalized to another part of knowledge-guided language understanding models (Xin et al., 2018; Kumar et al., 2018; Li et al., 2019). It treats entity embeddings as attention queries and word representations as attention key-value pairs. The process to compute the attention output h_{e_i} of entity e_i and the general representation O_e is formulated as

$$\begin{aligned} h_{e_i} &= W^T \text{softmax}(W A e_i), \\ O_e &= \text{Enc}_e(h_{e_1}, \dots, h_{e_m}), \end{aligned} \quad (2)$$

where A is a bi-linear matrix, Enc_e is identical to that of knowledge augmentation framework, and $W = \{w_1, \dots, w_n\}$ is the word representation matrix. Note that the word representations can be contextualized, such as outputs of CNN, LSTM, or uncontextualized, such as GloVe.

4 Experimental Setup

KG Details. We adopt a sub-graph of Wikidata to train the KE models. There are 5,039,998 entities, 927 relations, and 24,248,796 fact triples. Note that the triples appearing in the relation classification task are removed from this KG.

Training Details. The training details of KE models and text encoders are introduced in Section A of the Appendix due to the space limitation.

Frameworks. We denote the knowledge augmentation framework by **+Aug** and the knowledge attention framework by **+Att**. We can combine the names of text encoders and frameworks to represent instantiations of these frameworks, e.g., CNN+Aug denotes the instance of CNN in the knowledge augmentation framework. We also report the results of only using KE or text in downstream tasks, which is denoted by **KE-Only** or **Text-Only**.

Evaluation Datasets. We choose four typical knowledge-driven NLU tasks, which can be divided into two types: **Entity-oriented tasks** including relation classification and entity typing, and **General NLU tasks** including information retrieval and fact verification. Examples of these tasks are shown in Figure 1. Unlike previous work designing specific models for each task, this work systematically evaluates two general knowledge-

Relation Classification Text: <u>Newton</u> served as the president of the <u>Royal Society</u> . Relation: <u>member_of</u>
Entity Typing Text: Newton served as the president of the <u>Royal Society</u> . Type: <u>organization</u>
Information Retrieval Query: How large was <u>Medusa</u> ? Document: <u>Medusa</u> , a reticulated <u>python</u> , clocked in at 7.67 meters (25 feet, 2 inches) long in its official world record measurement. Relevance: <u>high</u>
Fact Verification Statement: <u>Home Alone</u> was written by <u>Barack Obama</u> . Correctness: <u>false</u>

Figure 1: Examples for the evaluation tasks. The underlined mentions are the entities appearing in the inputs. The last line in each box is the corresponding label.

guided frameworks on these tasks.

(1) *Relation Classification* aims to determine the correct relation between two entities in a given sentence, which is an important task for information extraction. In this work, we choose a large-scale human-annotated relation classification dataset FewRel (Han et al., 2018b), which consists of 56,000 instances and 80 relation classes.

(2) *Entity Typing* aims to infer the semantic type of the entity mention by its context. In this work, we adopt the large-scale entity typing dataset used by Xin et al. (2018), which contains 68 types, 860,011 training instances, 66,860 development instances and 68,242 testing instances.

(3) *Information Retrieval* aims to capture the query-document relevancy by calculating the similarities between queries and documents. We use ClueWeb09 as the dataset since Xiong et al. (2017) have shown that the understanding of its many cases needs external knowledge. There are 200 queries and we adopt the five-fold cross-validation.

(4) *Fact Verification* aims to verify the correctness of a given statement regarding entities. Here, we verify the statement without evidence and keep the statements with more than two entities in FEVER (Thorne et al., 2018) to evaluate the help of KE. There are 17,918 instances for training, 2,238 instances for development and testing, respectively.

We select these tasks for two reasons. First, there have been many works and datasets in these tasks for knowledge integration, making the comprehensive comparison available. Second, these four tasks are representative: entity typing focuses on a single entity; relation classification focuses on the relation between two entities; information retrieval focuses on the similarities between entities; and fact verification focuses on the reasoning among entities. Notably, We exclude language modeling and dialog

KE	MRR	HITS@1	HITS@3	HITS@10
DistMult	0.226	0.173	0.252	0.327
TransE	0.279	0.196	0.334	0.416
RotatE	0.302	0.234	0.345	0.418

Table 1: Performance on knowledge graph completion.

system because we focus on NLU tasks here.

Evaluation Metrics. For relation classification and fact verification, which are multi-class classification tasks, we report the prediction accuracy. For entity typing, which is a multi-label classification task, we adopt micro averaged metrics to measure the model performance. For information retrieval, which is a ranking task, we adopt precision@20 (P@20) and NDCG@20 as evaluation metrics².

5 Experimental Results

5.1 Effects of KE Models

We first investigate **whether KE models can help language understanding**. To this end, we evaluate the performance of different KE models on KGC and the effects of these KE models with different text encoders and knowledge-guided frameworks.

The performance on KGC is shown in Table 1. From the table, we observe that RotatE achieves the best results on all evaluation metrics.

For downstream NLU tasks, we report the results of KE-Only, Text-Only and two knowledge-guided frameworks based on three text encoders in Table 2. The best performance of each text encoder is in boldface. From the table, we find that:

(1) For CNN and LSTM, both knowledge augmentation and knowledge attention frameworks achieve better results compared to the Text-Only models on almost every task. It shows the generality and effectiveness of two knowledge-guided frameworks and the usefulness of KE models for downstream NLU tasks. Besides, knowledge augmentation works better than knowledge attention for three text encoders in most of the tasks. This suggests that directly using entity embeddings as features is more suitable for integrating KE’s information into conventional text encoders.

(2) Good performance on KGC does not correlate with good performance on NLU tasks. On the one hand, RotatE, which achieves the best results in KGC, does not have consistent superior performance when applied to these NLU tasks. On the

²The evaluation toolkit provided by TREC (Van Gysel and de Rijke, 2018) is used.

Text Enc.	Framework	KE	RC Acc	P	ET R	F1	P@20	IR NDCG@20	FV Acc	
-	KE-Only	DistMult	0.724	0.738	0.695	0.716	0.167	0.172	0.564	
		TransE	0.803	0.649	0.741	0.692	0.157	0.165	0.576	
		RotatE	0.683	0.364	0.704	0.480	0.168	0.180	0.580	
CNN	Text-Only	-	0.668	0.768	0.626	0.690	0.258	0.276	0.737	
		+Aug	DistMult	0.772	0.811	0.714	0.759	0.243	0.283	0.740
			TransE	0.857	0.828	0.733	0.778	0.282	0.328	0.740
	RotatE		0.796	0.812	0.675	0.738	0.271	0.320	0.743	
	+Att	DistMult	0.670	0.783	0.685	0.731	0.268	0.320	0.752	
		TransE	0.722	0.806	0.737	0.770	0.280	0.326	0.747	
RotatE		0.673	0.797	0.721	0.757	0.276	0.317	0.754		
LSTM	Text-Only	-	0.619	0.754	0.668	0.708	0.228	0.241	0.733	
		+Aug	DistMult	0.753	0.797	0.714	0.753	0.228	0.272	0.723
			TransE	0.848	0.830	0.720	0.771	0.274	0.322	0.736
	RotatE		0.774	0.780	0.703	0.740	0.267	0.326	0.732	
	+Att	DistMult	0.645	0.795	0.722	0.757	0.235	0.255	0.750	
		TransE	0.660	0.809	0.735	0.770	0.260	0.301	0.747	
RotatE		0.610	0.781	0.699	0.737	0.237	0.274	0.748		
BERT	Text-Only	-	0.849	0.769	0.755	0.762	0.294	0.332	0.831	
		+Aug	DistMult	0.858	0.767	0.758	0.762	0.295	0.330	0.831
			TransE	0.859	0.764	0.751	0.758	0.307	0.348	0.832
	RotatE		0.858	0.787	0.747	0.766	0.296	0.333	0.826	
	+Att	DistMult	0.843	0.748	0.749	0.748	0.288	0.321	0.827	
		TransE	0.845	0.745	0.731	0.738	0.280	0.326	0.826	
RotatE		0.845	0.763	0.743	0.753	0.297	0.338	0.831		

Table 2: Performance on four NLU tasks with different KE models and text encoders. RC: Relation Classification; ET: Entity Typing; IR: Information Retrieval; FV: Fact Verification.

other hand, the performance of KE-Only is also inconsistent with two knowledge-guided frameworks. For example, in entity typing, DistMult performs best in KE-Only but TransE achieves the best result when applied in knowledge-guided frameworks. These observations indicate that the knowledge-driven frameworks may not be able to utilize the information of KE well.

(3) For fact verification, there is a tiny difference between the performance of knowledge-guided models and Text-Only models. It suggests that only using entity embeddings could not benefit this task. We will further study this phenomenon in Section 5.4 to discuss how to combine language and knowledge information.

To further investigate the performance mismatch between KGC and downstream tasks, we compare the distance of the semantic spaces between each KE model and GloVE (GloVE is the word embedding used by CNN and LSTM). We suppose that for a model using two sources (GloVE and KE), if two sources are closer to each other, the model will use them easier. Specifically, for each matched entity e , which has both word embedding and knowledge embedding, we use a unified lin-

ear matrix \mathbf{M} to transform its entity embedding \mathbf{e} into its corresponding word embedding \mathbf{w}_e in GloVE, and define the semantic distance of e as $\|\mathbf{e}\mathbf{M} - \mathbf{w}_e\|_2$. The semantic space distance between KE and GloVE is defined as the average distance of all matched entities. We show the semantic space distance of three KE models and visualize the embeddings of 40 entities with t-SNE (Maaten and Hinton, 2008) in Figure 2. From this figure, we can see that: (1) The semantic space distances of three KE models are more consistent with their performance on NLU tasks than their performance on KGC. (2) DistMult has the largest semantic space distance with GloVE, which may be one reason for its large performance gap between the KE-only framework and two knowledge-guided frameworks. According to this observation, to better utilize KE in NLU tasks, a feasible solution is to build connections between KE and text representation by joint training or designing specific fusion architectures. We will discuss this more in Section 5.4.

5.2 Analysis on KE’s Helpful Information

Based on the promising results of KE models in Table 2, we further raise a question: **What informa-**

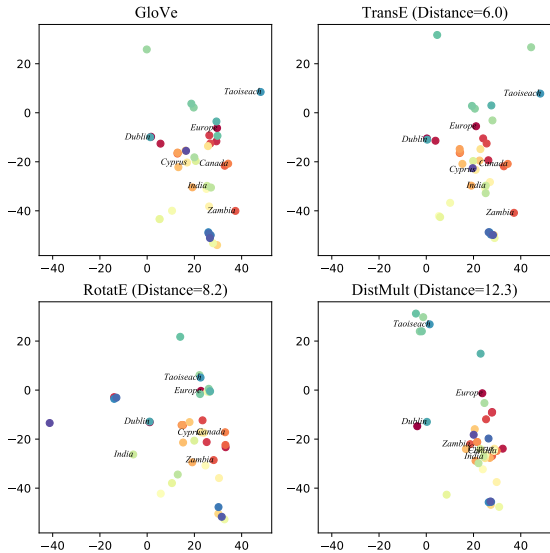


Figure 2: Visualization of GloVe and KE models. Distance: the semantic space distance. An entity in different plots has the same color. The KE with a smaller semantic distance with GloVe will plot more similarly.

tion of KE could help language understanding?

This will help us figure out the possible directions to improve the utilization of KE. Due to the space limit, we only report the results using TransE, while the conclusion of our analysis is consistent among all KE models. For the results of DistMult and RotatE, please refer to Section B of the Appendix. From our study, there are two main kinds of information in KE benefiting language understanding:

Entity Similarity Information. Intuitively, the similarities between different entities are the most important information provided by entity embeddings. Based on the similarities, we can cluster similar entities together, which could be beneficial for entity typing, and directly using the similarities could benefit information retrieval.

For entity clustering, we cluster entities with K-means (MacQueen et al., 1967), and assume the entities in the same cluster share the same information. To evaluate the effect of this information, we replace input entity embeddings with their corresponding cluster embeddings (the average of all entity embeddings in the cluster). Here, we set the number of clusters as 100. From the results in Figure 3, we can see that:

(1) In the entity typing (ET) task, the knowledge-guided frameworks using cluster embeddings perform very closely to those with original entity embeddings and the KE-Only model using cluster embeddings even achieves better performance than the KE-Only model with original embeddings. It

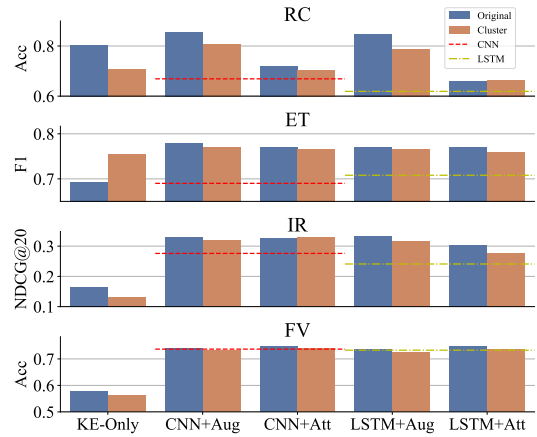


Figure 3: Comparisons of the models using entity embeddings (Origin) and cluster embeddings (Cluster).

Text Enc.	Framework	Top 1	Top 5	NDCG@20
-	KE-Only	12.71	12.95	0.165
CNN	Text-Only	11.42	11.54	0.276
	+Aug	13.50	13.30	0.328
	+Att	12.10	12.80	0.326
LSTM	Text-Only	11.79	12.29	0.241
	+Aug	13.16	12.63	0.322
	+Att	10.17	11.28	0.301
Ground Truth		13.33	-	-

Table 3: Average entity similarities of the query-document pairs having high relevance scores. Groud Truth is the entity similarities of ground truth pairs.

reveals that after removing the other information from inputs, the models may further make full use of the cluster information, which is related to the entity type and entity typing mainly benefits from the cluster information of KE.

(2) For the other tasks, using cluster embeddings can also bring improvements over Text-Only (dashed lines), while there is a performance degradation compared to the models using entity embeddings. It indicates the cluster information is useful for these tasks, but there still exists other information of KE that could help these tasks.

Directly using entity similarities may play an important role in information retrieval, which emphasizes capturing the similarities between queries and documents. To verify this, we calculate the entity similarities of the top-5 query-document pairs retrieved by the models. Specifically, given a query-document pair, we calculate the cosine similarities of all entity pairs between the query and document, and average them out as the entity similarities. We report the average results of query-document pairs in Table 3. We observe that the entity similarities of the ground truth are higher than those of most

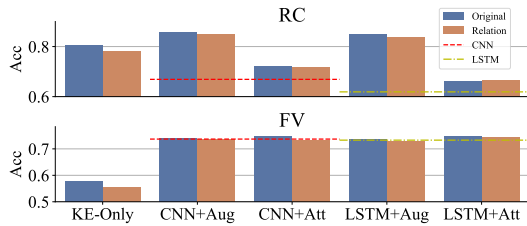


Figure 4: Comparisons of the models using entity embeddings (Origin) and relation embeddings (Relation).

models. The knowledge augmentation framework, having the highest entity similarities among models, achieves the best performance. It indicates that entity similarity is useful for information retrieval. However, we also need to consider both texts and entities’ information because KE-Only has high entity similarities but the worst performance.

Entity Relation Information. Since KE is learned from relational data, the relation information of KE should be important when utilizing KE. For example, relation classification (RC) and fact verification (FV) require modeling the relations among entities in text. To extract the relation information of KE, we calculate the relation embeddings according to the scoring function of KE with the entity embeddings (e.g., the relation embedding in TransE is the difference between head and tail entity embeddings). We replace the entity embeddings with the corresponding relation embeddings in the input. The results on relation classification and fact verification are reported in Figure 4. Figure 3 and 4 show that the relation information of KE is more useful than the cluster information for relation classification. However, for fact verification, the benefit of the relation information is similar to that of the cluster information. The reason is that fact verification requires a more complex utilization of the information, which will be further discussed in Section 5.4.

5.3 Utilizing KE for PLMs

From Table 2, we notice that BERT, which is a representative pre-trained language model (PLM) having powerful representation ability, benefits little from the KE models, and is even slightly degraded in fact verification. The reason is perhaps that PLMs such as BERT have learned rich factual knowledge through pre-training from large-scale corpora (Petroni et al., 2019). Hence, we consider a question: **Could KE still benefit PLMs in language understanding?** In other words, we explore how to effectively inject KE into PLMs.

Text Enc.	Framework	RC Acc	ET F1	IR NDCG	FV Acc
RoBERTa	Text-Only	0.852	0.765	0.350	0.841
	+Aug	0.845	0.768	0.329	0.842
	+Att	0.842	0.764	0.326	0.838
KEPLER	Text-Only	0.851	0.767	0.344	0.841
	+Aug	0.845	0.772	0.342	0.840
	+Att	0.848	0.780	0.339	0.836
BERT	Text-Only	0.849	0.762	0.332	0.831
	+Aug	0.859	0.758	0.348	0.832
	+Att	0.845	0.738	0.326	0.826
	ERNIE	0.878	0.799	0.340	0.842

Table 4: The results of PLMs with different frameworks on downstream tasks.

Firstly, we evaluate whether the pre-training task will influence the ability to utilize KE. Hence, we choose RoBERTa (Liu et al., 2019b), which adopts a better pre-training paradigm than BERT, and KEPLER (Wang et al., 2019d), which adds a new pre-training task based on KGs to RoBERTa. Secondly, we evaluate ERNIE (Zhang et al., 2019), which injects KE into PLMs via designing specified model architectures. Note that ERNIE is based on BERT and can be treated as a new knowledge-guided framework. From the results in Table 4, we observe that: (1) Similar to BERT, RoBERTa also cannot benefit from the knowledge-guided frameworks in all tasks. Besides, for entity typing, KEPLER+Att achieves more than 1% improvement over KEPLER. Although the improvement is not consistent across different tasks, it still suggests the possibility of enhancing the ability to utilize knowledge by pre-training. (2) ERNIE works well in all tasks, which shows the effectiveness of designing specific modules for knowledge-guided PLMs.

In summary, the designs of both pre-training tasks and injection modules are promising for better utilizing KE for PLMs and they still need further research for new tasks and frameworks.

5.4 Error Analysis and Discussion

In this section, we analyze the errors of the knowledge guided frameworks to discover their weaknesses for further research. According to the observation on the error cases, we categorize the errors into three kinds: knowledge representation, knowledge selection, knowledge utilization. The detailed descriptions of these errors are shown in Table 5. We also provide some examples in Section C of the Appendix. We randomly sample 100 error cases from four downstream NLU tasks used in this work

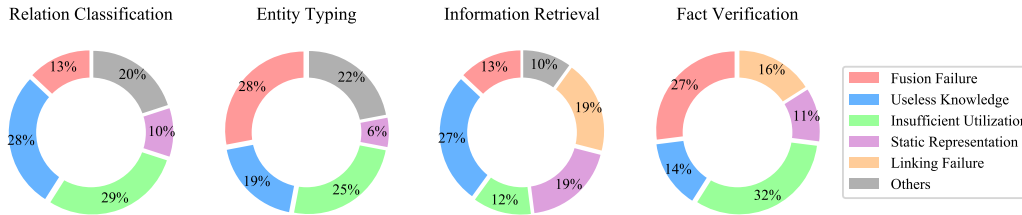


Figure 5: Types of error cases in four knowledge-driven NLU tasks.

Knowledge Representation

1. *Static Representation* - The static representations of entities, which ignore the text context, cannot satisfy the demand of tasks.

Knowledge Selection

2. *Linking Failure* - The results of entity linking contain some errors, which mislead the knowledge-guided model.

Knowledge Utilization

3. *Fusion Failure* - The KE-Only model makes a correct prediction while the fusion model does not.
4. *Useless Knowledge* - The model makes the correct prediction with text while KE causes extra noise.
5. *Insufficient Utilization* - The instances needs both text and KE information but the fusion model makes a incorrect prediction.

Table 5: Descriptions for errors. (5 types, 3 categories.)

and report the statistics of the errors in Figure 5.

From the statistics, we observe that these three kinds of errors account for a great portion of the error cases. For fact verification, where knowledge-guided frameworks do not work well, the knowledge information is still needed. Based on the results, we discuss several promising directions requiring further efforts for each kind of errors:

(1) For knowledge representation, the error of static representation appears in all four downstream tasks. Existing work (Wang et al., 2014; Zhang et al., 2015; Xu et al., 2016) have preliminarily verified the effectiveness of joint learning, which can build connections between knowledge and language. Nevertheless, how to represent knowledge based on the context is an important problem for further research, which is similar to contextualized word representation (Peters et al., 2018).

(2) For knowledge selection, linking failure appears in information retrieval and fact verification where linking results are not human-annotated. This emphasizes the importance of entity linking. Inspired by end-to-end relation extraction (Li and Ji, 2014; Miwa and Bansal, 2016), which jointly extracts entity mentions and relations, we believe entity linking can be integrated into knowledge-

guided frameworks for better results. KnowBERT (Peters et al., 2019) is pioneering work, which introduces a soft entity linking mechanism.

(3) For knowledge utilization, in each task, this kind of errors accounts for more than 50% and three sub-types of error have similar portion. Although we have shown that text encoders can benefit from KE, they cannot make full use of KE and sometimes fail in knowledge fusion. What’s worse, some cases indeed need external knowledge but the insufficient utilization makes it work not well. Meanwhile, directly using KE will introduce useless knowledge to the model in some cases. Hence, we need to explore how to better encode both knowledge and text information simultaneously. Recently, ERNIE (Zhang et al., 2019) and KnowBERT (Peters et al., 2019) provide us a novel perspective to fuse knowledge and language in pre-training. Besides designing novel pre-training objectives, we could also design more suitable model architectures for utilizing KE. There exist some works investigating novel model architectures to encode relational knowledge, such as memory-based models (Yang and Mitchell, 2017; Mihaylov and Frank, 2018), graph neural network-based models (Sun et al., 2018), retrieval-based models (Guu et al., 2020), etc. Nevertheless, the problem of how to effectively fuse knowledge in language understanding still remains unsolved.

6 Conclusion

In this work, we seek to better understand how KE could benefit language understanding in four knowledge-driven NLU tasks. Our comprehensive evaluation reveals (1) the performance inconsistency between knowledge graph completion and downstream NLU tasks; (2) two main kinds of useful information of KE in downstream NLU tasks; (3) how KE could benefit powerful PLMs. These observations can provide some insights for the follow-up researchers to better exploit KE in language understanding tasks.

608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660

References

Sungjin Ahn, Heeyoul Choi, Tanel Pärnamaa, and Yoshua Bengio. 2016. [A neural knowledge language model](#). *CoRR*, abs/1608.00318.

Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. 2008. [Freebase: A collaboratively created graph database for structuring human knowledge](#). In *Proceedings of SIGKDD*, pages 1247–1250.

Antoine Bordes, Sumit Chopra, and Jason Weston. 2014a. [Question answering with subgraph embeddings](#). In *Proceedings of EMNLP*, pages 615–620.

Antoine Bordes, Xavier Glorot, Jason Weston, and Yoshua Bengio. 2012. [Joint learning of words and meaning representations for open-text semantic parsing](#). In *Proceedings of AISTATS*, pages 127–135.

Antoine Bordes, Xavier Glorot, Jason Weston, and Yoshua Bengio. 2014b. [A semantic matching energy function for learning with multi-relational data](#). *Machine Learning*, 94(2):233–259.

Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. 2013. [Translating embeddings for modeling multi-relational data](#). In *Proceedings of NIPS*, pages 2787–2795.

Antoine Bordes, Jason Weston, Ronan Collobert, Yoshua Bengio, et al. 2011. [Learning structured embeddings of knowledge bases](#). In *Proceedings of AAAI*, pages 301–306.

Antoine Bordes, Jason Weston, and Nicolas Usunier. 2014c. [Open question answering with weakly supervised embedding models](#). In *Proceedings of ECML PKDD*, pages 165–180.

Michael M. Bronstein, Joan Bruna, Yann LeCun, Arthur Szlam, and Pierre Vandergheynst. 2017. [Geometric deep learning: Going beyond euclidean data](#). *IEEE Signal Process. Mag.*, 34(4):18–42.

Ines Chami, Adva Wolf, Da-Cheng Juan, Frederic Sala, Sujith Ravi, and Christopher Ré. 2020. [Low-dimensional hyperbolic knowledge graph embeddings](#). In *Proceedings of ACL*, pages 6901–6914.

Chia-Hui Chang, Mohammed Kayed, Moheb R. Girgis, and Khaled F. Shaalan. 2006. [A survey of web information extraction systems](#). *IEEE Trans. Knowl. Data Eng.*, 18(10):1411–1428.

Zhuyun Dai, Chenyan Xiong, Jamie Callan, and Zhiyuan Liu. 2018. [Convolutional neural networks for soft-matching n-grams in ad-hoc search](#). In *Proceedings of WSDM*, pages 126–134.

Tim Dettmers, Pasquale Minervini, Pontus Stenetorp, and Sebastian Riedel. 2018. [Convolutional 2D knowledge graph embeddings](#). In *Proceedings of AAAI*, pages 1811–1818.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of NAACL-HLT*, pages 4171–4186.

Bhuwan Dhingra, Manzil Zaheer, Vidhisha Balachandran, Graham Neubig, Ruslan Salakhutdinov, and William W. Cohen. 2020. [Differentiable reasoning over a virtual knowledge base](#). In *Proceedings of ICLR*.

Marjan Ghazvininejad, Chris Brockett, Ming-Wei Chang, Bill Dolan, Jianfeng Gao, Wen-tau Yih, and Michel Galley. 2018. [A knowledge-grounded neural conversation model](#). In *Proceedings of AAAI*, pages 5110–5117.

Yihong Gu, Jun Yan, Hao Zhu, Zhiyuan Liu, Ruobing Xie, Maosong Sun, Fen Lin, and Leyu Lin. 2018. [Language modeling with sparse product of sememe experts](#). In *Proceedings of EMNLP*, pages 4642–4651.

Shu Guo, Quan Wang, Lihong Wang, Bin Wang, and Li Guo. 2018. [Knowledge graph embedding with iterative guidance from soft rules](#). In *Proceedings of AAAI*, pages 4816–4823.

Kelvin Guu, Kenton Lee, Zora Tung, Panupong Papatat, and Ming-Wei Chang. 2020. [Retrieval augmented language model pre-training](#). In *Proceedings of ICML*.

Xu Han, Zhiyuan Liu, and Maosong Sun. 2018a. [Neural knowledge acquisition via mutual attention between knowledge graph and text](#). In *Proceedings of AAAI*, pages 4832–4839.

Xu Han, Hao Zhu, Pengfei Yu, Ziyun Wang, Yuan Yao, Zhiyuan Liu, and Maosong Sun. 2018b. [FewRel: A large-scale supervised few-shot relation classification dataset with state-of-the-art evaluation](#). In *Proceedings of EMNLP*, pages 4803–4809.

He He, Anusha Balakrishnan, Mihail Eric, and Percy Liang. 2017. [Learning symmetric collaborative dialogue agents with dynamic knowledge graph embeddings](#). In *Proceedings of ACL*, pages 1766–1776.

Shizhu He, Kang Liu, Guoliang Ji, and Jun Zhao. 2015. [Learning to represent knowledge graphs with Gaussian embedding](#). In *Proceedings of CIKM*, pages 623–632.

Sepp Hochreiter and Jürgen Schmidhuber. 1997. [Long short-term memory](#). *Neural computation*, 9(8):1735–1780.

Raphael Hoffmann, Congle Zhang, Xiao Ling, Luke Zettlemoyer, and Daniel S Weld. 2011. [Knowledge-based weak supervision for information extraction of overlapping relations](#). In *Proceedings of ACL-HLT*, pages 541–550.

714	Jian Hu, Gang Wang, Fred Lochovsky, Jian-tao Sun, and Zheng Chen. 2009. Understanding user’s query intent with Wikipedia . In <i>Proceedings of the WWW</i> , pages 471–480.	768
715		769
716		770
717		771
718	Xiao Huang, Jingyuan Zhang, Dingcheng Li, and Ping Li. 2019. Knowledge graph embedding based question answering . In <i>Proceedings of WSDM</i> , page 105–113.	772
719		773
720		774
721		775
722	Rodolphe Jenatton, Nicolas L Roux, Antoine Bordes, and Guillaume R Obozinski. 2012. A latent factor model for highly multi-relational data . In <i>Proceedings of NIPS</i> , pages 3167–3175.	776
723		777
724		778
725		
726	Guoliang Ji, Shizhu He, Liheng Xu, Kang Liu, and Jun Zhao. 2015. Knowledge graph embedding via dynamic mapping matrix . In <i>Proceedings of ACL</i> , pages 687–696.	779
727		780
728		781
729		782
730	Guoliang Ji, Kang Liu, Shizhu He, and Jun Zhao. 2016. Knowledge graph completion with adaptive sparse transfer matrix . In <i>Proceedings of AAAI</i> , pages 5997–6004.	783
731		784
732		785
733		786
734	Yoon Kim. 2014. Convolutional neural networks for sentence classification . In <i>Proceedings of EMNLP</i> , pages 1746–1751.	787
735		
736		
737	Diederik Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. <i>arXiv preprint arXiv:1412.6980</i> .	788
738		789
739		790
740	Abhishek Kumar, Daisuke Kawahara, and Sadao Kurohashi. 2018. Knowledge-enriched two-layered attention network for sentiment analysis . In <i>Proceedings of NAACL</i> , pages 253–258.	791
741		792
742		
743		
744	Pengfei Li, Kezhi Mao, Xuefeng Yang, and Qi Li. 2019. Improving relation extraction with knowledge-attention . In <i>Proceedings of EMNLP</i> , pages 229–239.	793
745		794
746		795
747		796
748	Qi Li and Heng Ji. 2014. Incremental joint extraction of entity mentions and relations . In <i>Proceedings of ACL</i> , pages 402–412.	797
749		798
750		799
751	Yankai Lin, Zhiyuan Liu, Maosong Sun, Yang Liu, and Xuan Zhu. 2015. Learning entity and relation embeddings for knowledge graph completion . In <i>Proceedings of AAAI</i> , pages 2181–2187.	800
752		801
753		802
754		803
755	Angli Liu, Jingfei Du, and Veselin Stoyanov. 2019a. Knowledge-augmented language model and its application to unsupervised named-entity recognition . In <i>Proceedings of NAACL-HLT</i> , pages 1142–1150.	804
756		805
757		806
758		807
759	Quan Liu, Hui Jiang, Andrew Evdokimov, Zhen-Hua Ling, Xiaodan Zhu, Si Wei, and Yu Hu. 2016. Probabilistic reasoning via deep learning: Neural association models . <i>arXiv preprint arXiv:1603.07704</i> .	808
760		809
761		810
762		811
763	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019b. RoBERTa: A robustly optimized bert pretraining approach . <i>arXiv preprint arXiv:1907.11692</i> .	812
764		813
765		814
766		815
767		816
	Zhenghao Liu, Chenyan Xiong, Maosong Sun, and Zhiyuan Liu. 2018. Entity-Duet neural ranking: Understanding the role of knowledge graph semantics in neural information retrieval . In <i>Proceedings of ACL</i> , pages 2395–2405.	817
		818
		819
		820
		821
	Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE . <i>Journal of machine learning research</i> , 9(Nov):2579–2605.	
	James MacQueen et al. 1967. Some methods for classification and analysis of multivariate observations . In <i>Proceedings of BSMSP</i> , pages 281–297.	
	Todor Mihaylov and Anette Frank. 2018. Knowledgeable reader: Enhancing cloze-style reading comprehension with external commonsense knowledge . In <i>Proceedings of ACL</i> , pages 821–832.	
	Alexander Miller, Adam Fisch, Jesse Dodge, Amir-Hossein Karimi, Antoine Bordes, and Jason Weston. 2016. Key-value memory networks for directly reading documents . In <i>Proceedings of EMNLP</i> , pages 1400–1409.	
	Pasquale Minervini, Luca Costabello, Emir Muñoz, Vít Nováček, and Pierre-Yves Vandenbussche. 2017. Regularizing knowledge graph embeddings via equivalence and inversion axioms . In <i>Proceedings of ECML/PKDD</i> , volume 10534, pages 668–683.	
	Makoto Miwa and Mohit Bansal. 2016. End-to-end relation extraction using LSTMs on sequences and tree structures . In <i>Proceedings of ACL</i> , pages 1105–1116.	
	Maximilian Nickel, Lorenzo Rosasco, and Tomaso Poggio. 2016. Holographic embeddings of knowledge graphs . In <i>Proceedings of AAAI</i> , pages 1955–1961.	
	Ankur Padia, Konstantinos Kalpakis, Francis Ferraro, and Tim Finin. 2019. Knowledge graph fact prediction via knowledge-enriched tensor factorization . <i>J. Web Semant.</i> , 59.	
	Md. Rizwan Parvez, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. 2018. Building language models for text with named entities . In <i>Proceedings of ACL</i> , pages 2373–2383.	
	Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation . In <i>Proceedings of EMNLP</i> , pages 1532–1543.	
	Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations . In <i>Proceedings of NAACL</i> .	
	Matthew E. Peters, Mark Neumann, Robert L Logan, Roy Schwartz, Vidur Joshi, Sameer Singh, and Noah A. Smith. 2019. Knowledge enhanced contextual word representations . In <i>Proceedings of EMNLP-IJCNLP</i> , pages 43–54.	

822	Fabio Petroni, Tim Rocktäschel, Sebastian Riedel,	Hongwei Wang, Fuzheng Zhang, Miao Zhao, Wen-	874
823	Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and	jie Li, Xing Xie, and Minyi Guo. 2019a. Multi-	875
824	Alexander Miller. 2019. Language models as knowl-	task feature learning for knowledge graph enhanced	876
825	edge bases? In <i>Proceedings of EMNLP</i> , pages	recommendation . In <i>Proceedings of WWW</i> , page	877
826	2463–2473.	2000–2010.	878
827	Richard Socher, Danqi Chen, Christopher D Manning,	Shen Wang, Xiaokai Wei, Cícero Nogueira dos Santos,	879
828	and Andrew Ng. 2013. Reasoning with neural tensor	Zhiguo Wang, Ramesh Nallapati, Andrew O.	880
829	networks for knowledge base completion . In <i>Pro-</i>	Arnold, Bing Xiang, and Philip S. Yu. 2020.	881
830	<i>ceedings of NIPS</i> , pages 926–934.	H2KGAT: hierarchical hyperbolic knowledge graph	882
831	Fabian M Suchanek, Gjergji Kasneci, and Gerhard	attention network . In <i>Proceedings of EMNLP</i> , pages	883
832	Weikum. 2007. YAGO: A core of semantic knowl-	4952–4962.	884
833	edge . In <i>Proceedings of WWW</i> , pages 697–706.	Xiang Wang, Xiangnan He, Yixin Cao, Meng Liu, and	885
834	Haitian Sun, Tania Bedrax-Weiss, and William W. Co-	Tat-Seng Chua. 2019b. KGAT: Knowledge graph	886
835	hen. 2019a. Pullnet: Open domain question an-	attention network for recommendation . In <i>Proceed-</i>	887
836	swering with iterative retrieval on knowledge bases	<i>ings of SIGKDD</i> , page 950–958.	888
837	and text . In <i>Proceedings of EMNLP-IJCNLP</i> , pages	Xiang Wang, Dingxian Wang, Canran Xu, Xiangnan	889
838	2380–2390.	He, Yixin Cao, and Tat-Seng Chua. 2019c. Explain-	890
839	Haitian Sun, Bhuwan Dhingra, Manzil Zaheer, Kathryn	able reasoning over knowledge graphs for recom-	891
840	Mazaitis, Ruslan Salakhutdinov, and William Cohen.	mendation . In <i>Proceedings of AAAI</i> , pages 5329–	892
841	2018. Open domain question answering using early	5336.	893
842	fusion of knowledge bases and text . In <i>Proceedings</i>	Xiaozhi Wang, Tianyu Gao, Zhaocheng Zhu, Zhiyuan	894
843	<i>of EMNLP</i> , pages 4231–4242.	Liu, Juanzi Li, and Jian Tang. 2019d. KEPLER:	895
844	Zhiqing Sun, Zhi-Hong Deng, Jian-Yun Nie, and Jian	A unified model for knowledge embedding and pre-	896
845	Tang. 2019b. RotatE: Knowledge graph embedding	trained language representation . <i>arXiv preprint</i>	897
846	by relational rotation in complex space . In <i>Proceed-</i>	<i>arXiv:1911.06136</i> .	898
847	<i>ings of ICLR</i> .	Zhen Wang, Jianwen Zhang, Jianlin Feng, and Zheng	899
848	Ilya Sutskever, Joshua B Tenenbaum, and Ruslan	Chen. 2014. Knowledge graph embedding by trans-	900
849	Salakhutdinov. 2009. Modelling relational data us-	lating on hyperplanes . In <i>Proceedings of AAAI</i> ,	901
850	ing Bayesian clustered tensor factorization . In <i>Pro-</i>	pages 1112–1119.	902
851	<i>ceedings of NIPS</i> , pages 1821–1828.	Jason Weston, Antoine Bordes, Oksana Yakhnenko,	903
852	James Thorne, Andreas Vlachos, Christos	and Nicolas Usunier. 2013. Connecting language	904
853	Christodoulopoulos, and Arpit Mittal. 2018.	and knowledge bases with embedding models for re-	905
854	FEVER: A large-scale dataset for fact extraction	lation extraction . In <i>Proceedings of EMNLP</i> , pages	906
855	and verification . pages 809–819.	1366–1371.	907
856	Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric	Yikun Xian, Zuohui Fu, S. Muthukrishnan, Gerard	908
857	Gaussier, and Guillaume Bouchard. 2016. Complex	de Melo, and Yongfeng Zhang. 2019. Reinforce-	909
858	embeddings for simple link prediction . In <i>Proceed-</i>	ment knowledge graph reasoning for explainable	910
859	<i>ings of ICML</i> , pages 2071–2080.	recommendation . In <i>Proceedings of SIGIR</i> , page	911
860	Christophe Van Gysel and Maarten de Rijke. 2018.	285–294.	912
861	Pytrex_eval: An extremely fast python interface to	Han Xiao, Minlie Huang, and Xiaoyan Zhu. 2016.	913
862	trec_eval . In <i>Proceedings of SIGIR</i> , pages 873–876.	From one point to a manifold: Orbit models for	914
863	Pat Verga, Haitian Sun, Livio Baldini Soares, and	knowledge graph embedding . In <i>Proceedings of IJ-</i>	915
864	William W. Cohen. 2020. Facts as experts: Adapt-	<i>CAI</i> , pages 1315–1321.	916
865	able and interpretable neural memory over symbolic	Ji Xin, Yankai Lin, Zhiyuan Liu, and Maosong Sun.	917
866	knowledge . <i>CoRR</i> , abs/2007.00849.	2018. Improving neural fine-grained entity typing	918
867	Denny Vrandečić and Markus Krötzsch. 2014. Wiki-	with knowledge attention . In <i>Proceedings of AAAI</i> .	919
868	data: A free collaborative knowledge base . <i>Commu-</i>	Chenyan Xiong, Jamie Callan, and Tie-Yan Liu. 2017.	920
869	<i>nications of the ACM</i> , 57(10):78–85.	Word-entity duet representations for document rank-	921
870	Hongwei Wang, Fuzheng Zhang, Xing Xie, and Minyi	ing . In <i>Proceedings of SIGIR</i> , pages 763–772.	922
871	Guo. 2018. DKN: Deep knowledge-aware network	Jiacheng Xu, Kan Chen, Xipeng Qiu, and Xuanjing	923
872	for news recommendation . In <i>Proceedings of WWW</i> ,	Huang. 2016. Knowledge graph representation with	924
873	page 1835–1844.	jointly structural and textual encoding . In <i>Proceed-</i>	925
		<i>ings of IJCAI</i> , pages 1318–1324.	926

927 Bishan Yang and Tom Mitchell. 2017. [Leveraging](#)
 928 [knowledge bases in LSTMs for improving machine](#)
 929 [reading](#). In *Proceedings of ACL*, pages 1436–1446.

930 Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng
 931 Gao, and Li Deng. 2014. [Embedding entities and](#)
 932 [relations for learning and inference in knowledge](#)
 933 [bases](#). In *Proceedings of ICLR*.

934 Michihiro Yasunaga, Hongyu Ren, Antoine Bosse-
 935 lut, Percy Liang, and Jure Leskovec. 2021. [QA-](#)
 936 [GNN: reasoning with language models and knowl-](#)
 937 [edge graphs for question answering](#). *CoRR*,
 938 abs/2104.06378.

939 Dongxu Zhang, Bin Yuan, Dong Wang, and Rong Liu.
 940 2015. [Joint semantic relevance learning with text](#)
 941 [data and graph knowledge](#). In *Proceedings of ACL-*
 942 *IJCNLP workshop*, pages 32–40.

943 Fuzheng Zhang, Nicholas Jing Yuan, Defu Lian, Xing
 944 Xie, and Wei-Ying Ma. 2016. [Collaborative knowl-](#)
 945 [edge base embedding for recommender systems](#). In
 946 *Proceedings of SIGKDD*, pages 353–362.

947 Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang,
 948 Maosong Sun, and Qun Liu. 2019. [ERNIE: En-](#)
 949 [hanced language representation with informative en-](#)
 950 [tities](#). In *Proceedings of ACL*, pages 1441–1451.

951 Zhaocheng Zhu, Shizhen Xu, Jian Tang, and Meng Qu.
 952 2019. [Graphvite: A high-performance cpu-gpu hy-](#)
 953 [brid system for node embedding](#). In *Proceedings of*
 954 *WWW*, pages 2494–2504.

955 A Experimental Setup

956 **Training Details of KE models.** The sub-graph
 957 of Wikidata ³ is extracted by (Zhang et al.,
 958 2019). We divide these fact triples into two parts:
 959 24, 247, 796 triples for training and 1,000 triples
 960 for validation.

961 In this work, we evaluate three typical KE
 962 models: DistMult, TransE, and RotatE. We use
 963 GraphVite (Zhu et al., 2019), a high-performance
 964 KE system, to train these models. We follow
 965 most of hyper-parameters provided by GraphVite
 966 for large-scale KE and only search for the best
 967 learning rate based on the result of validation, 0.6
 968 from {0.2, 0.4, 0.6, 0.8} for DistMult, 0.008 from
 969 {0.004, 0.008, 0.01, 0.02} for TransE, 0.01 from
 970 {0.008, 0.01, 0.02, 0.04} for RotatE. We set the di-
 971 mension of KE as 128, which achieves the best
 972 performance on downstream NLP tasks.

973 **Training Details of Text Encoders.** In this
 974 work, we evaluate three typical text encoders. All
 975 models are optimized by Adam (Kingma and Ba,
 976 2014) except for CNN and LSTM in information

³<https://www.wikidata.org>

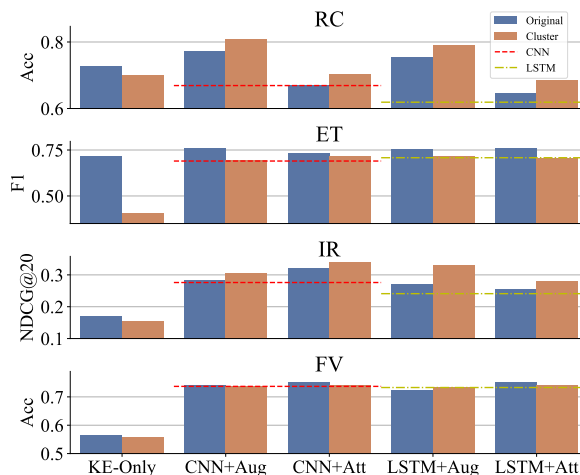


Figure 6: Comparisons of the models using entity embeddings (Origin) and cluster embeddings (Cluster). The KE model is DistMult.

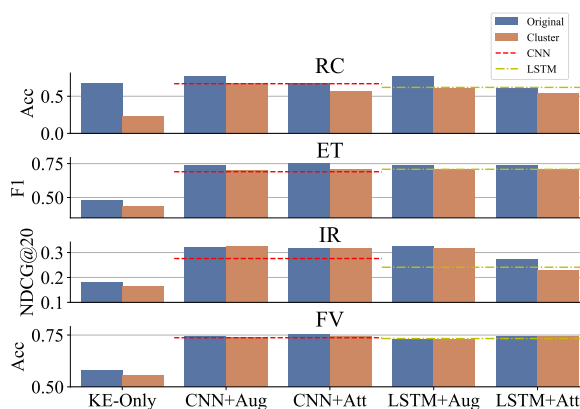


Figure 7: Comparisons of the models using entity embeddings (Origin) and cluster embeddings (Cluster). The KE model is RotatE.

977 retrieval, which use SGD. The hyperparameters for
 978 these models are as follows: (1) **CNN**. We adopt a
 979 single layer CNN. The hidden size of CNN is 100.
 980 We set the batchsize as 32 for relation classification
 981 and 100 for the others. We train the models with the
 982 learning rate of 0.001 for Adam and 0.1 for SGD.
 983 The input word embeddings are GloVe (Pennington
 984 et al., 2014) with the dimension of 50. (2) **LSTM**.
 985 We adopt a single layer bi-directional LSTM. The
 986 hyper-parameters of LSTM are as the same as those
 987 of CNN. (3) **BERT**. We use BERT_{BASE} released
 988 by Google and follow most of hyper-parameters
 989 provided by (Devlin et al., 2019) except that the
 990 training epoch, which varies in different tasks (10
 991 for relation classification, 3 for entity typing, 2 for
 992 information retrieval, 30 for fact verification). The
 993 learning rate is searched from {2e-5, 3e-5, 5e-5}.
 994 The hyper-parameters of RoBERTa and KEPLER

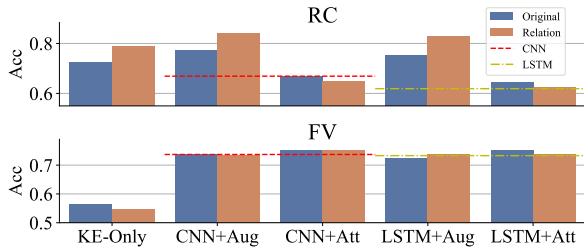


Figure 8: Comparisons between the models using the original entity embeddings and the models using the relation information of DistMult.

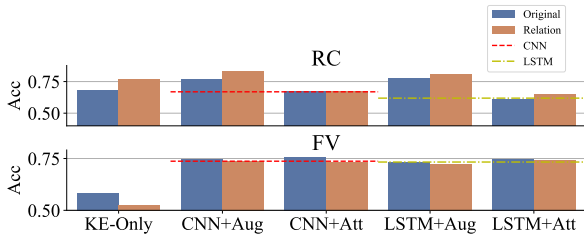


Figure 9: Comparisons between the models using the original entity embeddings and the models using the relation information of RotatE.

also follow the setting of BERT.

Computing Infrastructure and Runtime. We use NVIDIA RTX 2080Ti GPUs and each experiment uses one GPU. The average runtime of LSTM and CNN varies from several minutes to several tens of minutes according to different tasks. And the average runtime of BERT varies from several tens of minutes to several hours.

Entity Linking. For information retrieval and fact verification, which only provide the raw texts, we use TAGME⁴ to link the entities mentioned in text to KGs. Meanwhile, we use the entity linking provided by the datasets for RE and ET. To avoid information leakage, we exclude the triples in the test set of RE from the KG in the training of KE.

B Analysis of KE

The results about cluster information for DistMult and RotatE are shown in Figure 6 and 7 respectively. The results about entity similarity information for DistMult and RotatE are shown in Table 6 and 7 respectively. The results about relation information for DistMult and RotatE are shown in Figure 8 and 9 respectively.

⁴<https://tagme.d4science.org/tagme/>

Text Enc.	Framework	Top 1	Top 5	NDCG@20
-	KE-Only	13.10	12.40	0.172
CNN	Text-Only	9.56	9.94	0.276
	+Aug	12.77	12.86	0.283
	+Att	9.99	10.53	0.320
LSTM	Text-Only	10.51	10.60	0.241
	+Aug	11.6	11.7	0.272
	+Att	10.07	10.45	0.255
Groud Truth		11.92		

Table 6: The entity similarity of the query-document pairs having high relevance score for these models. Groud Truth is the entity similarity of the ground truth pairs. The KE model is DistMult.

Text Enc.	Framework	Top 1	Top 5	NDCG@20
-	KE-Only	14.60	13.30	0.180
CNN	Text-Only	8.25	8.53	0.276
	+Aug	9.88	9.55	0.301
	+Att	9.23	9.86	0.322
LSTM	Text-Only	8.79	9.12	0.241
	+Aug	11.3	10.8	0.326
	+Att	9.54	8.89	0.274
Groud Truth		10.30		

Table 7: The entity similarity of the query-document pairs having high relevance score for these models. Groud Truth is the entity similarity of the ground truth pairs. The KE model is RotatE.

C Error Cases

We provide some example of error cases for each task.

Type of Error	Text	Label	Prediction	KE Prediction	Text Prediction
Static Representation	<i>Lin Liheng</i> is the daughter of Lin Biao and <i>Ye Qun</i> , nicknamed "Dou Dou".	mother	sibling	spouse	spouse
Fusion Failure	The company 's first completed game was " <i>Odin Sphere</i> " for the PlayStation 2 , which was published by <i>Atlus</i> .	publisher	developer	publisher	developer
Useless Knowledge	His next two films "Kutty" and " <i>Uthama Puthiran</i> ", were both collaborations with director <i>Mithran Jawahar</i> .	director	screenwriter	screenwriter	director
Inefficient Utilization	<i>Alphonse John Smith</i> was a 20th-century bishop in the <i>Catholic Church</i> in the United States.	religion	main_subject	participant	language_of_work

Figure 10: Error cases of relation classification.

Type of Error	Text	Label	Prediction	KE Prediction	Text Prediction
Static Representation	The song begins as an <i>acoustic guitar</i> driven pop song and then shifts into a slower bridge section.	art; genre;	instrument	instrument	instrument
Fusion Failure	VET studies are offered Xavier is one of only fifteen schools in Victoria to offer <i>Latin</i> .	language	art	language	art
Useless Knowledge	Denis Smith was born in <i>Meir Stoke</i> on Trent the second youngest of seven siblings.	citytown; administrative_region	sports_team	administrative_region	citytown; administrative_region
Inefficient Utilization	On 9 January 2012 Donadoni was unveiled as head coach of <i>Serie A</i> club Parma replacing Franco Colomba.	organization; sports_league	organization	art	profession

Figure 11: Error cases of entity typing.

Type of Error	Text	Label	Prediction	KE Prediction	Text Prediction
Static Representation	Query: <i>Idaho</i> state flower Document: List of <i>U.S.</i> state flowers – Wikipedia ...	high relevance	low relevance	low relevance	low relevance
Linking Failure	Query: Dangers of <i>asbestos</i> Document: ... <i>South Dakota</i> ... <i>Idaho</i> ... <i>South Carolina</i> ... <i>Hawaii</i> ... <i>asbestos</i> (missing) removal should only be performed by qualified professionals since the risks associate with an improperly conducted <i>asbestos</i> (missing) removal are quite high ...	high relevance	low relevance	low relevance	low relevance
Fusion Failure	Query: <i>Poker tournaments</i> Document: Free <i>Poker Tournaments</i> - Free <i>Poker Tournaments</i> Freerolls play in a free poker tournament...	high relevance	low relevance	high relevance	low relevance
Useless Knowledge	Query: <i>Website design</i> hosting Document: Taos Web Design , Taos <i>Website Design</i> , Taos Web site Design , Taos Web Hosting home news about us web design cms website seo blog web hosting ecommerce ...	high relevance	low relevance	low relevance	high relevance
Inefficient Utilization	Query: <i>Mothers day</i> songs Document: Children 's <i>Lullabies</i> : ... lullabies that start with a all the pretty little horses all the pretty little ponies all the world loves to hear mothers sing all through the night version 1 ...	low relevance	high relevance	high relevance	high relevance

Figure 12: Error cases of information retrieval.

Type of Error	Text	Label	Prediction	KE Prediction	Text Prediction
Static Representation	<i>Connie Britton</i> played a role in the first season of <i>American Crime Story</i> .	false	true	true	true
Linking Failure	<i>Hansel</i> and <i>Gretel</i> is of Mexican (missing) origin.	false	true	true	true
Fusion Failure	The <i>New York Knicks</i> compete in the <i>National Basketball Association</i> (NBA).	true	false	true	false
Useless Knowledge	<i>Live Through This</i> has sold over 1.6 million copies in the <i>United States</i> .	true	false	false	true
Inefficient Utilization	<i>Theodore Roosevelt</i> attended <i>Harvard College</i> in 1824.	false	true	true	true

Figure 13: Error cases of fact verification.