

Incorporating Relation Knowledge into Commonsense Reading Comprehension with Multi-task Learning

Jiangnan Xia, Chen Wu, Ming Yan
 Alibaba DAMO Academy
 Hangzhou, China
 {jiangnan.xjn,wuchen.wc,ym119608}@alibaba-inc.com

ABSTRACT

This paper focuses on how to take advantage of external relational knowledge to improve machine reading comprehension (MRC) with multi-task learning. Most of the traditional methods in MRC assume that the knowledge used to get the correct answer generally exists in the given documents. However, in real-world task, part of knowledge may not be mentioned and machines should be equipped with the ability to leverage external knowledge. In this paper, we integrate relational knowledge into MRC model for commonsense reasoning. Specifically, based on a pre-trained language model (LM), we design two auxiliary relation-aware tasks to predict if there exists any commonsense relation and what is the relation type between two words, in order to better model the interactions between document and candidate answer option. We conduct experiments on two multi-choice benchmark datasets: the SemEval-2018 Task 11 and the Cloze Story Test. The experimental results demonstrate the effectiveness of the proposed method, which achieves superior performance compared with the comparable baselines on both datasets.

CCS CONCEPTS

• **Computing methodologies** → **Information extraction**; Semantic networks; Intelligent agents.

KEYWORDS

machine reading comprehension, commonsense reasoning, multi-task learning

ACM Reference Format:

Jiangnan Xia, Chen Wu, Ming Yan. 2019. Incorporating Relation Knowledge into Commonsense Reading Comprehension with Multi-task Learning. In *The 28th ACM International Conference on Information and Knowledge Management (CIKM '19)*, November 3–7, 2019, Beijing, China. ACM, New York, NY, USA, 4 pages. <https://doi.org/10.1145/3357384.3358165>

1 INTRODUCTION

Machine reading comprehension (MRC) enables machines with the ability to answer questions with given corresponding documents. Recent years have witnessed the bloom of various well-designed

MRC models [13, 14, 16], which achieve promising performance when provided with adequate manually labeled instances [9]. However, these models generally assume that the knowledge required to answer the questions has already existed in the documents, which does not hold at some time. How to leverage the commonsense knowledge for better reading comprehension remains largely unexplored.

Recently, some preliminary studies have begun to incorporate certain side information (e.g., triplets from external knowledge base) into the model design of various NLP tasks, such as question answering [1] and conversation generation [18]. Generally, there are two lines of this work. The first line focuses on designing task-specific model structures [1, 15], which exploit the retrieved concepts from external knowledge base for enhancing the representation. Recently, the other line has studied to pre-train a language model over large corpus to learn the inherent word-level knowledge in an unsupervised way [4, 8], which achieves very promising performance.

The first line of work is usually carefully designed for the target task, which is not widely applicable. The second line can only learn the co-occurrence of words or entities in the context, while it may not be that robust for some complex scenarios such as *reasoning* task. For example, to answer the question "Was the light bulb still hot?" when the document is given as "I went into my bedroom and flipped the light switch. Oh, I see that the ceiling lamp is not turning on...", machines should have the commonsense knowledge that "the bulb is not hot when turning off" to correctly answer the question. The explicit relation information can act as a bridge to connect the scattered context, which may not be easily captured. Therefore, the aim of this paper is to take the advantage of both the pre-trained language model and the explicit relation knowledge from the external knowledge base for commonsense reading comprehension.

Specifically, we first extract the triplets from the popular ConceptNet knowledge base [11] and design two auxiliary relation-aware tasks to predict if there exists any relation and what is the relation type between two concepts. To make the model be aware of the commonsense relations between concepts, we propose a multi-task learning framework to jointly learn the prediction of the target MRC task and the two relation-aware tasks in a unified model. We conduct experiments on two multi-choice commonsense reading comprehension datasets: Story Cloze Test [6] and SemEval-2018 Task 11 [7]. Experimental results demonstrate the effectiveness of our method, which achieves superior performance compared with the comparable baselines on both datasets.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

CIKM '19, November 3–7, 2019, Beijing, China

© 2019 Association for Computing Machinery.

ACM ISBN 978-1-4503-6976-3/19/11...\$15.00

<https://doi.org/10.1145/3357384.3358165>

2 RELATED WORK

Previous studies mainly focused on developing effective model structures to improve the reading ability of the systems [14, 16], which have achieved promising performance. However, the success on these tasks is not adequate considering the model’s ability of commonsense reasoning. Recently, a number of efforts have been invested in developing datasets for commonsense reading comprehension such as Story Cloze Test and SemEval-2018 Task 11 [6, 7]. In these datasets, part of required knowledge may not be mentioned in the document and machines should be equipped with commonsense knowledge to make correct prediction. There exists increasing interest in incorporating commonsense knowledge into commonsense reading comprehension tasks. Most of previous studies focused on developing special model structures to introduce external knowledge into neural network models [1, 15], which have achieved promising results. For example, Yang and Mitchell [15] use concepts from WordNet and weighted average vectors of the retrieved concepts to calculate a new LSTM state. These methods relied on task-specific model structures which are difficult to adapt to other tasks. Pre-trained language model such as BERT and GPT[4, 8] is also used as a kind of commonsense knowledge source. However, the LM method mainly captures the co-occurrence of words and phrases and cannot address some more complex problems which may require the reasoning ability.

Unlike previous work, we incorporate external knowledge by jointly training MRC model with two auxiliary tasks which are relevant to commonsense knowledge. The model can learn to fill in the knowledge gap without changing the original model structure.

3 KNOWLEDGE-ENRICHED MRC MODEL

3.1 Task Definition

Here we formally define the task of multi-choice commonsense reading comprehension. Given a reference document D (a question q if possible), a set of N answer options $\{O_1, O_2, \dots, O_N\}$ and an external knowledge base $F = \{f_1, \dots, f_M\}$, the goal is to choose the correct answer option according to their probabilities $\{p_1, p_2, \dots, p_N\}$ given by MRC model, where N is the total number of options.

In this paper, we use ConceptNet knowledge base [11], which is a large semantic network of commonsense knowledge with a total of about 630k facts. Each fact f_i is represented as a triplet $f_i = (subject, relation, object)$, where *subject* and *object* can be a word or phrase and *relation* is a relation type. An example is: $([Car]_{subj}, [UsedFor]_{rel}, [Driving]_{obj})$

3.2 Overall Framework

The proposed method can be roughly divided into three parts: a pre-trained LM encoder, a task-specific prediction layer for multi-choice MRC and two relation-aware auxiliary tasks. The overall framework is shown in Figure 1.

The pre-trained LM encoder acts as the foundation of the model, which is used to capture the relationship between question, document and answer options. Here we utilize BERT [4] as the pre-trained encoder for its superior performance in a range of natural language understanding tasks. Specially, we concatenate the given document, question (as sentence A) and each option (as sentence B)

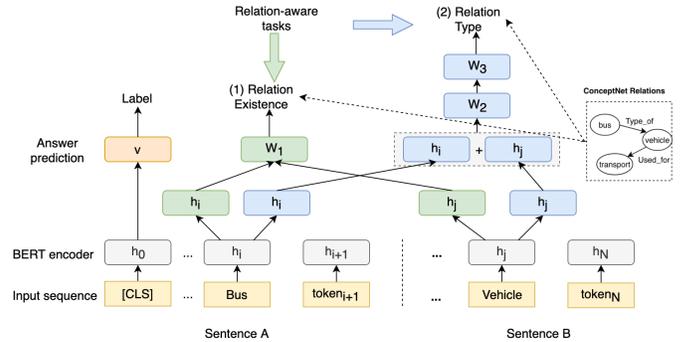


Figure 1: MRC model with two relation-aware tasks

with special delimiters as one segment, which is then fed into BERT encoder. The input sequence is packed as “[CLS]D(Q)[SEP]O[SEP]”¹, where [CLS] and [SEP] are the special delimiters. After BERT encoder, we obtain the contextualized word representation $h_i^L \in \mathbb{R}^H$ for the i -th input token from the final layer of BERT. H is the dimension of hidden state.

Next, on top of BERT encoder, we add a task-specific output layer and view the multi-choice MRC as a multi-class classification task. Specifically, we apply a linear head layer plus a softmax layer on the final contextualized word representation of [CLS] token h_0^L . We minimize the Negative Log Likelihood (NLL) loss with respect to the correct class label, as:

$$L^{AP}(\hat{O}|D) = -\log \frac{\exp(\mathbf{v}^T \hat{h}_0)}{\sum_{k=1}^N \exp(\mathbf{v}^T h_0^k)} \quad (1)$$

where \hat{h}_0 is the final hidden state of the correct option \hat{O} , N is the number of options and $\mathbf{v}^T \in \mathbb{R}^H$ is a learnable vector.

Finally, to make the model be aware of certain implicit commonsense relations between concepts, we further introduce two auxiliary relation-aware prediction tasks for joint learning. Since it may be difficult to directly predict the actual relation type between two concepts without adequate training data, we split the relation prediction problem into two related tasks: i.e., *relation-existence* task and *relation-type* task. In *relation-existence*, we basically add an auxiliary task to predict if there exists any relation between two concepts, which is a relatively easy task. Then, we take one step further to decide what is the right type of the relation in *relation-existence*. The basic premise is that by guiding the MRC model training with extra relation information, the proposed model can be equipped with the ability to capture some underlying commonsense relationships. The two auxiliary tasks are jointly trained with the multi-choice answer prediction task. In the following, we will describe the two auxiliary tasks in detail.

3.3 Incorporating Relation Knowledge

Task 1 is the *relation-existence* task. Following [4], we first convert the concept to a set of BPE tokens tokens A and tokens B, with beginning index i and j in the input sequence respectively. The probability of whether there is a relation in each pair (tokens A,

¹We directly concatenate the question after the document to form the BERT input if a question is given.

Table 1: Number of examples in datasets.

Dataset	Train	Dev	Test
SemEval-2018 Task11	9,731	1,411	2,797
Story Cloze Test	1,433	347	1,871

tokens B) is computed as:

$$p_{ij}^{RE} = \text{sigmoid}(\mathbf{h}_i^T \mathbf{W}_1 \mathbf{h}_j) \quad (2)$$

where $\mathbf{W}_1 \in \mathbb{R}^{H \times H}$ is a trainable matrix.

We define the pair (tokens A, tokens B) that has relation in ConceptNet as a positive example and others as negative examples. We down-sample the negative examples and keep ratio of positive vs negative is $1 : \gamma$. We define the *relation-existence* loss as the average binary cross-entropy (BCE) loss:

$$\mathcal{L}^{RE} = \frac{1}{|A|} \frac{1}{|B|} \sum_{i=1}^{|A|} \sum_{j=1}^{|B|} \text{BCE}(y_{ij}^{RE}, p_{ij}^{RE}) \quad (3)$$

where $|A|$, $|B|$ are the number of sampled concepts in sentence A and sentence B respectively, y_{ij}^{RE} is the label of whether there is a relation between concepts.

Task 2 is the *relation-type* task. We predict the relation type between tokens A and tokens B. The *relation-type* probabilities are computed as:

$$p_{ij}^{RT} = \text{softmax}(\mathbf{W}_3 \text{ReLU}(\mathbf{W}_2[\mathbf{h}_i; \mathbf{h}_j])) \quad (4)$$

where $\mathbf{W}_2 \in \mathbb{R}^{H \times 2H}$ and $\mathbf{W}_3 \in \mathbb{R}^{R \times H}$ are new trainable matrices, R is the number of selected relation types².

The *relation-type* loss is computed as:

$$\mathcal{L}^{RT} = -\frac{1}{|S|} \sum_{i=1}^{|A|} \sum_{j=1}^{|B|} s_{ij} \log[p_{ij}^{RT}]_k \quad (5)$$

We define s_{ij} as the label whether there is a relation from sentence A to sentence B. k is the index of ground-truth relation in ConceptNet, $|S|$ is the number of relations among the tokens in two sentences.

As the three tasks share the same BERT architecture with only different linear head layers, we propose to train them together. The joint objective function is formulated as follows:

$$\mathcal{L} = \mathcal{L}^{AP} + \frac{1}{N} \sum_{l=1}^N (\lambda_1 \mathcal{L}_l^{RE} + \lambda_2 \mathcal{L}_l^{RT}) \quad (6)$$

where λ_1 and λ_2 are two hyper-parameters that control the weight of the tasks, N is number of options.

4 EXPERIMENTS

4.1 Dataset

We conduct experiments on two commonsense reading comprehension tasks: SemEval-2018 shared task11 [7] and Story Cloze Test [6]. The statistics of the datasets are shown in Table 1³.

²We don't use the entire set of relation types because not all the types are defined clearly. We choose 34 kind of relations except "RelatedTo", "ExternalURL" and "dbpedia".
³Story Cloze Test consists 98,161 unlabeled examples 1,871 labeled examples. Following [8], we divide the labeled examples to a new training set and development set with 1,433 and 347 examples respectively.

Table 2: Performance on SemEval-2018 Task 11.

Model	ACC.(Test) %
NN-T [5]	80.23
HMA [3]	80.94
TriAN [12]	81.94
TriAN + $f_{cs}^{dir} + f_{cs}^{ins}$ [17]	81.80
TriAN + relation-aware tasks	82.84
BERT(base)	87.53
BERT(base) + relation-aware tasks	88.23

Table 3: Performance on Cloze Story Test.

Model	ACC.(Test) %
Style+RNLM [10]	75.2
HCM [2]	77.6
GPT [8]	86.5
BERT(base)	86.7
BERT(base) + relation-aware tasks	87.4

Table 4: Performance with different tasks

Model	ACC.(Dev)%	Δ
Basic Model	87.51	-
+ \mathcal{L}_{RE}	88.02	+0.51
+ \mathcal{L}_{RT}	87.87	+0.36
+ $\mathcal{L}_{RE} + \mathcal{L}_{RT}$	88.25	+0.74
+ $\mathcal{L}_{RT} + \text{"No Relation"}$	87.78	+0.41

4.2 Implementation Details

We use the uncased BERT(base) [4] as pre-trained language model. We set the batch size to 24, learning rate to $2e-5$. The maximum sequence length is 384 for SemEval-2018 Task 11 and 512 for Story Cloze Test. We fine-tune for 3 epochs on each dataset. The task-specific hyper-parameters λ_1 and λ_2 are set to 0.5 and the ratio γ is set to 4.0.

4.3 Experimental Results and Analysis

Model Comparison The performances of our model on two datasets are shown in Table 2 and Table 3. We compare our single model with other existing systems (single model). We also adopt relation-aware tasks on the co-attention layer of the TriAN model [12]. From the result, we can observe that: (1) Our method achieves better performance on both datasets compared with previous methods and Bert(base) model. (2) By adopt relation-aware tasks on the attention layer of the TriAN model [12] on SemEval, the model performance can also be improved. The results show that the relation-aware tasks can help to better align sentences due to knowledge gap.

Effectiveness of Relation-aware Tasks To get better insight into our model, we analyze the benefit brought by using relation-aware tasks on the development set of SemEval-2018 Task 11. The performance of jointly training the basic answer prediction model

Table 5: Examples that require commonsense relations between concepts

Document	Question	Options	Commonsense Facts
[SemEval] ... We organized it from the fruit to the dairy in an organized order. We put the shopping list on our fridge so that we wouldn't forget it the next day when we went to go buy the food.	What did they write the list on?	(A) Paper . (B) Fridge. Correct: A	[paper , RelatedTo, write]
[Cloze Story Test] My kitchen had too much trash in it. I cleaned it up and put it into bags. I took the bags outside of my house. I then carried the bags down my driveway to the trash can.	-	(A) I missed the trash in the kitchen. (B) I was glad to get rid of the trash. Correct: B	[get rid of, RelatedTo, clean up]
[SemEval] ... I settled on Earl Gray, which is a black tea flavored with bergamot orange. I filled the kettle with water and placed it on the stove, turning on the burner...	Why did they use a kettle?	(A) To drink from. (B) To boil water . Correct: B	[kettle , UsedFor , boil water]; [kettle , RelatedTo, drinking]

with different tasks is shown in table 4. From the result we can see that by incorporating the auxiliary *relation-existence* task (\mathcal{L}_{RE}) or *relation-type* task (\mathcal{L}_{RT}) in the joint learning framework, the performance can always improved. The result shows the advantage of incorporating auxiliary tasks. Besides, the performance gain by adding *relation-existence* task is larger, which shows *relation-existence* task can incorporate more knowledge into the model. We also attempt to merge two relation-aware tasks into one task by simply taking "No-Relation" as a special type of relation. The model performance is just slightly higher than using *relation-type* task and lower than using two tasks separately. The result is due to the number of "No-Relation" labels is much more than other relation types, which makes the task hard to train.

Analysis Table 5 shows the examples that are incorrectly predicted by BERT(base), while correctly solved by incorporating relation-aware tasks. The first two examples can benefit from *relation-existence* knowledge. From the first example we can see that the retrieved relation between concepts from ConceptNet provide useful evidence to connect the question to the correct option (A). The second example is from Cloze Story Test dataset, we can see that the retrieved relation is also helpful in making correct prediction. The third example from SemEval requires *relation-type* knowledge. but the relation type (*kettle*, *UsedFor*, *boilwater*) in option (A), is more relevant to the question, which shows that relation type knowledge can be used as side information to do the prediction.

5 CONCLUSION

In this paper, we aim to enrich the neural model with external knowledge to improve commonsense reading comprehension. We use two auxiliary relation-aware tasks to incorporate ConceptNet knowledge into the MRC model. Experimental results demonstrate the effectiveness of our method which achieves improvements compared with the pre-trained language model baselines on both datasets.

REFERENCES

- [1] Lisa Bauer, Yicheng Wang, and Mohit Bansal. 2018. Commonsense for Generative Multi-Hop Question Answering Tasks. *arXiv preprint arXiv:1809.06309* (2018).
- [2] Snigdha Chaturvedi, Haoruo Peng, and Dan Roth. 2017. Story comprehension for predicting what happens next. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. 1603–1614.
- [3] Zhipeng Chen, Yiming Cui, Wentao Ma, Shijin Wang, Ting Liu, and Guoping Hu. 2018. HFL-RC System at SemEval-2018 Task 11: Hybrid Multi-Aspects Model for Commonsense Reading Comprehension. *arXiv preprint arXiv:1803.05655* (2018).
- [4] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [5] Elizabeth Merkhofer, John Henderson, David Bloom, Laura Strickhart, and Guido Zarrella. 2018. MITRE at SemEval-2018 Task 11: Commonsense Reasoning without Commonsense Knowledge. In *Proceedings of The 12th International Workshop on Semantic Evaluation*. 1078–1082.
- [6] Nasrin Mostafazadeh, Michael Roth, Annie Louis, Nathanael Chambers, and James Allen. 2017. Lsdsem 2017 shared task: The story cloze test. In *Proceedings of the 2nd Workshop on Linking Models of Lexical, Sentential and Discourse-level Semantics*. 46–51.
- [7] Simon Ostermann, Michael Roth, Ashutosh Modi, Stefan Thater, and Manfred Pinkal. 2018. SemEval-2018 Task 11: Machine Comprehension using Commonsense Knowledge. In *Proceedings of The 12th International Workshop on Semantic Evaluation*. 747–757.
- [8] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training. (2018).
- [9] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250* (2016).
- [10] Roy Schwartz, Maarten Sap, Ioannis Konstas, Leila Zilles, Yejin Choi, and Noah A Smith. 2017. Story cloze task: Uw nlp system. In *Proceedings of the 2nd Workshop on Linking Models of Lexical, Sentential and Discourse-level Semantics*. 52–55.
- [11] Robert Speer, Joshua Chin, and Catherine Havasi. 2016. ConceptNet 5.5: An Open Multilingual Graph of General Knowledge. (2016).
- [12] Liang Wang. 2018. Yuanfudao at SemEval-2018 Task 11: Three-way Attention and Relational Knowledge for Commonsense Machine Comprehension. *arXiv preprint arXiv:1803.00191* (2018).
- [13] Wei Wang, Ming Yan, and Chen Wu. 2018. Multi-Granularity Hierarchical Attention Fusion Networks for Reading Comprehension and Question Answering. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 1705–1714.
- [14] Wenhui Wang, Nan Yang, Furu Wei, Baobao Chang, and Ming Zhou. 2017. Gated self-matching networks for reading comprehension and question answering. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, Vol. 1. 189–198.
- [15] Bishan Yang and Tom Mitchell. 2017. Leveraging knowledge bases in lstms for improving machine reading. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, Vol. 1. 1436–1446.
- [16] Adams Wei Yu, David Dohan, Minh-Thang Luong, Rui Zhao, Kai Chen, Mohammad Norouzi, and Quoc V Le. 2018. QANet: Combining Local Convolution with Global Self-Attention for Reading Comprehension. *arXiv preprint arXiv:1804.09541* (2018).
- [17] Wanjun Zhong, Duyu Tang, Nan Duan, Ming Zhou, Jiahai Wang, and Jian Yin. 2018. Improving question answering by commonsense-based pre-training. *arXiv preprint arXiv:1809.03568* (2018).
- [18] Hao Zhou, Tom Young, Minlie Huang, Haizhou Zhao, Jingfang Xu, and Xiaoyan Zhu. 2018. Commonsense Knowledge Aware Conversation Generation with Graph Attention.. In *IJCAI*. 4623–4629.