

# Improving open domain question answering with Knowledge Base and Wikipedia graph

Ruiyu Lin  
SUN YAT-SEN UNIVERSITY  
linry23@mail2.sysu.edu.cn

## ABSTRACT

A clear and well-documented  $\text{\LaTeX}$  document is presented as an

## KEYWORDS

neural networks

## ACM Reference Format:

Ruiyu Lin. 2018. Improving open domain question answering with Knowledge Base and Wikipedia graph. In *Woodstock '18: ACM Symposium on Neural Gaze Detection, June 03–05, 2018, Woodstock, NY*. ACM, New York, NY, USA, 2 pages. <https://doi.org/10.1145/1122445.1122456>

## 1 INTRODUCTION

### Input

- Dataset:(question, answer) pair.
- Knowledge Base:(entity, relation, entity) triple. Knowledge Base is a multi-relational graphs, each edge has a label and direction associated with it. every node in the graph is an entity.
- Wikipedia graph[1]:(text, text) pair. Wikipedia graph is constructed by Hyperlinks and within-document links, every node in the graph is an article.

**Output** representation of all the retrieved passage, as the input to a reader model to extract answer.

**Goal** To better embed the retrieved passage, which based on Wikipedia graph, with Knowledge Base knowledge.

**Method** Fuse Knowledge Base knowledge into Wikipedia graph to formulate  $(P_a, \text{relation}, P_b)$  triple.

- (1) identify the relations between two passage node

Assuming that :

$rSet_{(a,b)}$  : the set of relation between  $P_a$  and  $P_b$ , initialized empty

$(P_a, P_b)$  exists in Wikipedia graph.

$P_a$  contains n entities(  $e_{a1}, \dots, e_{an}$  ),

$P_b$  contains m entities(  $e_{ab}, \dots, e_{mb}$  )

If  $(e_{ai}, r, e_{bj})(1 \leq i, j \leq n)$  exists in Knowledge Base, add r into

$rSet_{(a,b)}$  . Finally,  $rSet_{a,b} = (r_1, \dots, r_k)$

- (2) corporate identified relations into passage node embedding

$$\alpha_{r_i} = \text{score}(h_a, e_{ai}, r_i, q)$$

$$h_{inter-b} = FNN(h_a, \sum_{i=1}^k \alpha_{r_i} * r_i)$$

$h_{inter-b}$  is the intermediate representation of passage b.  $\alpha_{r_i}$  is the relation score.

Suppose that  $h_b$  is link by  $[h_{a1}, \dots, h_{at}]$ , then update  $h_b$  by

$$h_b = GCN(h_b, \sum_1^t h_{inter-b})$$

- (3) answer extraction

We adopt Multi-passage BERT [7] as our reader model, which use Shared normalization[2], specifically to process passages independently, but compute the span probability across spans in all passages in every mini-batch. Globally normalizing answer scores across all passages of the same question enables to find better answers by utilizing more passages.

Denote the passage score as  $Pr(P_i|Q, P)$ , which reranks all retrieved passages

$$Pr(P_i|Q, P) = \text{softmax}(h_i^T W)$$

W is a trainable parameter.

The score of an answer span from passage  $P_i$  will be

$$Pr(a|Q, P) = Pr(P_i|Q, P)P_s(a_s|Q, P)P_e(a_e|Q, P).$$

### Some details remain specific

- (1) how to get the representation of relation.

In GRAFT-Nets[5], they average word vectors to compute a relation vector from the surface form of the relation.

In PullNet[4], embedding of relations are pretrained, and can be looked up from an embedding table.

In [3], they only consider the most frequent 100 relations, and pretrain their embedding.

[8] tokenize the relation, and then encode it by a shared LSTM with the question.

In [6], it handles multi-relational graphs representation where each edge has a label and direction associated with it, and jointly embeds both nodes and relations in a relational graph.

- (2) the representation of question

In GRAFT-Nets[5], the question representation is updated as  $h_q = FNN(\sum_{v \in S_q} h_v)$ , where  $S_q$  denotes the seed entities mentioned in the question.

In PullNet[4] and [8], the question representation is acquired by a LSTM.

In [3], the question is not directly encoded, but with passage jointly by BERT.

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](mailto:permissions@acm.org).

Woodstock '18, June 03–05, 2018, Woodstock, NY

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/10.1145/1122445.1122456>

## (3) how to score relation attention

Most work take the dot product between the embedding of relation and question. Considering the different framework here, we reformulate the score as:

$$\alpha_{r_i} = \text{score}(h_a, e_{ai}, r_i, q)$$

## REFERENCES

- [1] Akari Asai, Kazuma Hashimoto, Hannaneh Hajishirzi, Richard Socher, and Caiming Xiong. ICLR 2020. Learning to retrieve reasoning paths over wikipedia graph for question answering. (ICLR 2020).
- [2] C. Clark and M. Gardner. ACL 2018. Simple and Effective Multi-Paragraph Reading Comprehension. (ACL 2018).
- [3] Sewon Min, Danqi Chen, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2019. Knowledge guided text retrieval and reading for open domain question answering. *arXiv preprint arXiv:1911.03868* (2019).
- [4] Haitian Sun, Tania Bedrax-Weiss, and William W Cohen. 2019. Pullnet: Open domain question answering with iterative retrieval on knowledge bases and text. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing*, pages 2380–2390, 7.
- [5] Haitian Sun, Bhuwan Dhingra, Manzil Zaheer, Kathryn Mazaitis, Ruslan Salakhutdinov, and William W Cohen. 2018. Open domain question answering using early fusion of knowledge bases and text. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4231–4242.
- [6] S. Vashishth, S. Sanyal, V. Nitin, and P. Talukdar. ICLR 2020. Composition-based Multi-Relational Graph Convolutional Networks. (ICLR 2020).
- [7] Z. Wang, P. Ng, X. Ma, R. Nallapati, and B. Xiang. EMNLP 2019. Multi-passage BERT: A Globally Normalized BERT Model for Open-domain Question Answering. (EMNLP 2019).
- [8] Wenhan Xiong, Mo Yu, Shiyu Chang, Xiaoxiao Guo, and William Yang Wang. 2019. Improving question answering over incomplete kbs with knowledge-aware reader. *arXiv preprint arXiv:1905.07098* (2019).