

Effective Character-augmented Word Embedding for Machine Reading Comprehension

Zhuosheng Zhang^{1,2}, Yafang Huang^{1,2}, Pengfei Zhu^{1,2,3}, Hai Zhao^{1,2,*}

¹Department of Computer Science and Engineering, Shanghai Jiao Tong University ²Key Laboratory of Shanghai Education Commission for Intelligent Interaction and Cognitive Engineering, Shanghai Jiao Tong University, Shanghai, 200240, China ³School of Computer Science and Software Engineering, East China Normal University, China {zhangzs, huangyafang}@sjtu.edu.cn, 10152510190@stu.ecnu.edu.cn zhaohai@cs.sjtu.edu.cn

The cloze-style task can be described as a triple $\langle D; Q; A \rangle$, where D is a document (context), Q is a query over the contents of D, in which a word or phrase is replaced with a placeholder, and A is the answer to Q.

| Document | 1 早上,青蛙、小白兔、刺猬和大蚂蚁高高兴兴过桥去赶 | 1 In the morning, the frog, the little white rabbit, the hedgehog and the big ant happily crossed the | | |
|----------|------------------------------------|--|--|--|
| | 集。 | bridge for the market. | | |
| | 2 不料,中午下了一场大暴雨,哗啦啦的河水把桥冲走了。 | 2 Unexpectedly, a heavy rain fell at noon, and the water swept away the bridge. | | |
| | 3 天快黑了,小白兔、刺猬和大蚂蚁都不会游泳。 | 3 It was going dark. The little white rabbit, hedgehog and big ant cannot swim. | | |
| | 4 过不了河,急得哭了。 | 4 Unable to cross the river, they were about to cry. | | |
| | 5 这时,青蛙想,我可不能把朋友丢下,自己过河回家呀。 | 5 At that time, the frog made his mind that he could not leave his friend behind and went home alone. | | |
| | 6他一面劝大家不要着急,一面动脑筋。 | 6 Letting his friends take it easy, he thought and thought. | | |
| | 7 嗬,有了! | 7 Well, there you go! | | |
| | 8他说:"我有个朋女住在这儿,我去找他想想办法。 | 8 He said, "I have a friend who lives here, and I'll go and find him for help." | | |
| | 9青蛙找到了他的朋友,请求他说:"大家过不了河 | 9 The frog found his friend and told him, "We cannot get across the river. Please give us a | | |
| | 了,请帮个忙吧! | hand!" | | |
| | 10 鼹鼠说:"可以,请把大家领到我家里来吧。 | 10 The mole said, "That's fine, please bring them to my house." | | |
| | 11 鼹鼠把大家带到一个洞口,打开了电筒,让小白兔、刺 | 11 The mole took everyone to a hole, turned on the flashlight and asked the little white rabbit, the | | |
| | 猬、大蚂蚁和青蛙跟着他,"大家别害怕,一直朝前走。 | hedgehog, the big ant and the frog to follow him, saying, "Don't be afraid, just go ahead." | | |
| | 12 走呀走呀, 只听见上面"哗啦哗啦"的声音, 象唱歌。 | 12 They walked along, hearing the "walla-walla" sound, just like a song. | | |
| | 13 走着走着,突然,大家看见了天空,天上的月亮真亮呀。 | 13 All of a sudden, everyone saw the sky, and the moon was really bright. | | |
| | 14 小白兔回头一瞧,高兴极了:"哈,咱们过了河啦! | 14 The little white rabbit looked back and rejoiced: "ha, the river crossed!". | | |
| | 15 唷,真了不起。 15 "Oh, really great." | | | |
| | 16原来,鼹鼠在河底挖了一条很长的地道,从这头到那头。 | 16 Originally, the mole dug a very long tunnel under the river, from one end to the other. | | |
| | 17 青蛙、小白兔、刺猬和大蚂蚁是多么感激鼹鼠啊! | 17 How grateful the frog, the little white rabbit, the hedgehog and the big ant felt to the mole! | | |
| | 18 第二天,青蛙、小白兔、刺猬和大蚂蚁带来很多很多同 | 18 The next day, the frog, the little white rabbit, the hedgehog, and the big ant with a lot of his fellows, | | |
| | 伴, 杠着木头, 抬着石头, 要求鼹鼠让他们来把地道挖大 | took woods and stones. They asked the mole to dig tunnels bigger, and build a great bridge under the | | |
| | 些,修成河底大"桥"。 | river. | | |
| | 19 不久,他们就把鼹鼠家的地道,挖成了河底的一条大隧 | 19 It was not long before they dug a big tunnel under the river, and they could pass the river from the | | |
| | 道,大家可以从河底过何,还能通车,真有劲哩! | bottom of the river, and it could be open to traffic. It is amazing! | | |
| Query | 青蛙找到了他的朋友,请求他说:"大家过不了河 | The frog found his friend and told him, "We cannot get across the river. Please give us a | | |
| | 了,请帮个忙吧!" | hand!" | | |
| Answer | 鼹 鼠 | the mole | | |

Reading comprehension systems usually suffer from out-of-vocabulary (**OOV**) word issues, especially when the ground-truth answers contain rare words or name entities, which are hardly fully recorded in the vocabulary.



There are over **13,000** characters in Chinese while there are only **26** letters in English without regard to punctuation marks.

If a reading comprehension system can not effectively manage the OOV issues, the performance will not be semantically accurate for the task.

Two levels of embedding

Word-level Embedding 青蛙|和|小白兔|去|赶集 Character-level Embedding 青|蛙|和|小|白|兔|去|赶|集

- Intuitively, word-level representation is good at catching global context and dependency relationships between words. However, rare words are often expressed poorly due to data sparsity.
- Character embedding are more expressive to model sub-word morphologies, which is beneficial to deal with rare words. However, quite a lot of Chinese words, like "吉(auspicious)普(ordinary)" (jeep) are not semantically character-level compositional at all.
- Using extra features, such as named entity recognition (NER) and part-ofspeech (POS) tagging will result in tremendous computational complexity.

Word representation module



Framework

• Given the triple $\langle D; Q; A \rangle$, the system will be built in the following steps.



Trainable Embedding

| Motivation: insufficient training for UNK words | |
|---|--|
| Technique: | |
| • Sort the distionary according to the word | |

- Sort the dictionary according to the word frequency from high to low.
- A frequency filter ratio γ is set to filter out the low-frequency words (rare words) from the lookup table.
- For example, if γ is 0.9, then the last 10% low-frequency words will be mapped into UNK words.
- Thus, *AE*(*w*) can be rewritten as

 $AE(w) = \begin{cases} WE(w) \diamond SE(w) & \text{if } w \in H \\ UNK \diamond SE(w) & \text{otherwise} \end{cases}$



Fine-grained Embedding

- Word embedding WE(w) is indexed from word lookup table
- Characters of each word are successively fed to the forward GRU and backward GRU. The output for each input is the concatenation of the two vectors from both directions: $\overleftarrow{h_t} = \overrightarrow{h_t} \parallel \overleftarrow{h_t}$



 The augmented embedding (AE) is given by concatenating the word embedding and character-level representation. AE(w) = W E(w) || CE(w)

Attention Module

• Contextual representations of the document and query

 $H_q = \operatorname{BiGRU}(Q)$ $H_d = \operatorname{BiGRU}(D)$

• Gated-attention

$$\alpha_{i} = softmax(H_{q}^{\top}d_{i})$$
$$\beta_{i} = Q\alpha_{i}$$
$$x_{i} = d_{i} \odot \beta_{i}$$

• Probability of each candidate word as being the answer

$$p = softmax((q_t)^{\top} H_D)$$
$$P(w|D,Q) \propto \sum_{i \in I(w,D)} p_i$$

• The predicted answer

$$A^* = \mathrm{argmax}_{w \in C} P(w|D,Q)$$

Dataset and hyper-parameters

| | CMRC-2017 | | | PD | | | CFT |
|----------------------|-----------|--------|--------|---------|-------|-------|-------|
| | Train | Valid | Test | Train | Valid | Test | human |
| # Query | 354,295 | 2,000 | 3,000 | 870,710 | 3,000 | 3,000 | 1,953 |
| Max # words in docs | 486 | 481 | 484 | 618 | 536 | 634 | 414 |
| Max # words in query | 184 | 72 | 106 | 502 | 153 | 265 | 92 |
| Avg # words in docs | 324 | 321 | 307 | 379 | 425 | 410 | 153 |
| Avg # words in query | 27 | 19 | 23 | 38 | 38 | 41 | 20 |
| # Vocabulary | 94,352 | 21,821 | 38,704 | 248,160 | 536 | 634 | 414 |

- Three Chinese Machine Reading Comprehension datasets, namely CMRC-2017, People's Daily (PD) and Children Fairy Tales (CFT).
- We also use the Children's Book Test (CBT) dataset (Hill et al., 2015) to test the generalization ability in multi-lingual case.

CMRC-2017 Leaderboard

填空类问题 (Cloze-style Question)

| 最终排名 | 参赛单位 | 单/多系统 | 开发集准确率 | 测试集准确率↓ |
|------|---|-------|--------|---------|
| 8 1 | 6ESTATES PTE LTD | 多系统 | 81.85% | 81.90% |
| 7 2 | 上海交通大学仿脑计算与机器智能研究中心自然语言组 Shanghai Jiao Tong University (SJTU BCMI-NLP) | 多系统 | 78.35% | 80.67% |
| 83 | 南京云思创智信息科技有限公司 | 多系统 | 79.20% | 80.27% |

用户提问类问题 (User-Query Question)

| 最终排名 | 参赛单位 | 单/多系统 | 开发集准确率 | 测试集准确率↓ |
|------|--|-------|--------|---------|
| 8 1 | 华东师范大学 East China Normal University (ECNU) | 多系统 | 90.45% | 69.53% |
| 82 | 山西大学三队 Shanxi University (SXU-3) | 单系统 | 47.80% | 49.07% |
| 83 | 郑州大学 Zhengzhou University (ZZU) | 单系统 | 31.10% | 32.53% |

最佳单系统 (Best Single System)

| 最终排名 | 参赛单位 | 单/多系统 | 开发集准确率 | 测试集准确率↓ |
|------|---|-------|--------|---------|
| 8 1 | 上海交通大学仿脑计算与机器智能研究中心自然语言组 Shanghai Jiao Tong University (SJTU BCMI-NLP) | 单系统 | 76.15% | 77.73% |

Main results

- Our CAW Reader (*mul*) outperforms all other single models
- *mul* might be more informative than *concat* and *sum* operations

| Model | CMRC-2017 | | | |
|----------------------|-----------|-------|--|--|
| WIOUCI | Valid | Test | | |
| Random Guess † | 1.65 | 1.67 | | |
| Top Frequency † | 14.85 | 14.07 | | |
| AS Reader † | 69.75 | 71.23 | | |
| GA Reader | 72.90 | 74.10 | | |
| SJTU BCMI-NLP † | 76.15 | 77.73 | | |
| 6ESTATES PTE LTD † | 75.85 | 74.73 | | |
| Xinktech † | 77.15 | 77.53 | | |
| Ludong University † | 74.75 | 75.07 | | |
| ECNU † | 77.95 | 77.40 | | |
| WHU † | 78.20 | 76.53 | | |
| CAW Reader (WE only) | 69.70 | 70.13 | | |
| CAW Reader (concat) | 71.55 | 72.03 | | |
| CAW Reader (sum) | 72.90 | 74.07 | | |
| CAW Reader (mul) | 77.95 | 78.50 | | |

| Model | Stratagy | PD | | CFT | | - |
|------------------------------|----------|-------|-------------|------------|-------|------|
| Model | Strategy | Valid | Test | Test-human | | |
| AS Reader | - | 64.1 | 67.2 | 33 | - | |
| GA Reader | - | 64.1 | 65.2 | 35 | 5.7 | |
| CAS Reader | - | 65.2 | 68.1 | 35 | 5.0 | |
| | concat | 64.2 | 65.3 | 37.2 | | - |
| CAW Reader | sum | 65.0 | 68.1 | 38.7 | | |
| | mul | 69.4 | 70.5 | 39.7 | | _ |
| | | | CDT | | CDT | |
| Model | | | CBI | | | -CN |
| TT ! | | | Valid | Test | Valid | Test |
| Human ‡ | | | - | 81.6 | - | 81.6 |
| LSTMs ‡ | | | 51.2 | 41.8 | 62.6 | 56.0 |
| MemNets ‡ | | | 70.4 | 66.6 | 64.2 | 63.0 |
| AS Reader ‡ | | | 73.8 | 68.6 | 68.8 | 63.4 |
| Iterative Attentive Reader ‡ | | | 75.2 | 68.2 | 72.1 | 69.2 |
| EpiReader ‡ | | | 75.3 | 69.7 | 71.5 | 67.4 |
| AoA Reader ‡ | | | 77.8 | 72.0 | 72.2 | 69.4 |
| NSE ‡ | | | 78.2 | 73.2 | 74.3 | 71.9 |
| GA Reader ‡ | | | 74.9 | 69.0 | 69.0 | 63.9 |
| GA word char concat ‡ | | | 76.8 | 72.5 | 73.1 | 69.6 |
| GA scalar gate ‡ | | | 78.1 | 72.6 | 72.4 | 69.1 |
| GA fine-grained gate ‡ | | | 78.9 | 74.6 | 72.3 | 70.8 |
| FG Reader ‡ | | | 79.1 | 75.0 | 75.3 | 72.0 |
| CAW Reader | | | 78.4 | 74.9 | 74.8 | 71.5 |
| | | | | | | |

Influence of the short list

- When $\gamma = 0.9$, the models could obtain the best performance.
- It is not optimal to build the vocabulary among the whole training set.
- We can reduce the frequency filter ratio properly to promote the accuracy.



Conclusion

- Multiple embedding enhancement strategies
- Effective embedding architecture by attending character representations to word embedding with a short list to enhance the simple baseline for the reading comprehension task.
- The intensified embeddings can help our model achieve state-of the-art performance on multiple large-scale benchmark datasets.
- Different from most existing works that focus on either complex attention architectures or manual features, our model is more simple but effective.

Thanks! Q&A