# Multi-Perspective Fusion Network for Commonsense Reading Comprehension

Chunhua Liu[1], Yan Zhao[1], Qingyi Si[1], Haiou Zhang[1], Bohan Li[1], and Dong Yu[1,2] ✉

[1] Beijing Language and Culture University
[2] Beijing Advanced Innovation for Language Resources of BLCU
{chunhualiu596,zhaoyan.nlp}@gmail.com
{xk17sqy,hozhange1,yudong_blcu}@126.com
bohanli.lavida@gmail.com

**Abstract.** Commonsense Reading Comprehension (CRC) is a significantly challenging task, aiming at choosing the right answer for the question referring to a narrative passage, which may require commonsense knowledge inference. Most of the existing approaches only fuse the interaction information of choice, passage, and question in a simple combination manner from a *union* perspective, which lacks the comparison information on a deeper level. Instead, we propose a Multi-Perspective Fusion Network (MPFN), extending the single fusion method with multiple perspectives by introducing the *difference* and *similarity* fusion. More comprehensive and accurate information can be captured through the three types of fusion. We design several groups of experiments on MCScript dataset [11] to evaluate the effectiveness of the three types of fusion respectively. From the experimental results, we can conclude that the difference fusion is comparable with union fusion, and the similarity fusion needs to be activated by the union fusion. The experimental result also shows that our MPFN model achieves the state-of-the-art with an accuracy of 83.52% on the official test set.

**Keywords:** Commonsense Reading Comprehension · Fusion Network · Multi-Perspective

## 1 Introduction

Machine Reading Comprehension (MRC) is an extremely challenging topic in natural language processing field. It requires a system to answer the question referring to a given passage. In real reading comprehension, the human reader can fully understand the passage with the prior knowledge to answer the question. To directly relate commonsense knowledge to reading comprehension, SemEval2018 Task 11 defines a new sub-task called Commonsense Reading Comprehension, aiming at answering the questions that requires both commonsense knowledge and the understanding of the passage. The challenge of this task is how to answer questions with the commonsense knowledge that does not appear in the passage explicitly. Table 1 shows an example of CRC.

Most studies on CRC task are neural network based (NN-based) models, which typically have the following characteristics. Firstly, word representations are augmented by additional lexical information. Secondly, the interaction process is usually implemented by the attention mechanism, which can provide the interaction representations

---

**Passage:** It was night time and it was time to go to bed. The boy wanted to keep playing. I told him that after he got ready for bed I would read a story to him. He dawdled a bit but finally started getting ready for bed. First of all he had to take a bath. He splashed in the tub and split water all over the floor. Next he dried off in a big, fluffy blue towel. Then he brushed his teeth with his special Star Wars toothbrush. Next he dressed in his Star Wars underwear and then put on his Star Wars pajamas. His dad and I tucked him into his bed that was made with Star Wars sheets. He said his prayers. Next was story time. I pulled out his favorite book about (you guessed it) Star Wars. He gradually dozed off dreaming about Anakin Skywalker and a galaxy far, far away.

**Q1:** Did they sleep in the same room as their parents?
A. Yes, they all slept in one big loft      B. No they have their own room

**Q2:** Why didn't the child go to bed by themselves?
A. The child wanted to watch a Star Wars movie.      B. The child wanted to continue playing.

---

Table 1: An example of CRC.

like choice-aware passage, choice-aware question, and question-aware passage. Thirdly, the original representations and interaction representations are fused together and then aggregated by a Bidirectional Long Short-Term Memory Network (BiLSTM) [4] to get high-order semantic information. Fourthly, the final output based on their bilinear interactions.

The NN-based models have shown powerfulness on this task. However, there are still some limitations. Firstly, the two fusion processes of passage and question to choice are implemented separately, until producing the final output. Secondly, the existing fusion method used in reading comprehension task is usually implemented by concatenation [24,2], which is monotonous and cannot capture the partial comparison information between two parts. Studies on Natural Language Inference (NLI) have explored more functions [10,1], such as element-wise subtraction and element-wise multiplication, to capture more comparison information, which have been proved to be effective.

In this paper, we introduce a Muti-Perspective Fusion Network (MPFN) to tackle these limitations. The model can fuse the choice with passage and question simultaneously to get a multi-perspective fusion representation. Furthermore, inspired by the element-wise subtraction and element-wise multiplication function used in [1], we define three kinds of fusion functions from multiple perspectives to fuse choice, choice-aware passage, and choice-aware question. The three fusions are union fusion, difference fusion, and similarity fusion. Note that, we name the concatenation fusion method as union fusion in this paper, which collects the global information. The difference fusion and the similarity fusion can discover the different parts and similar parts among choice, choice-aware passage, and choice-aware question respectively.

MPFN comprises an encoding layer, a context fusion layer, and an output layer. In the encoding layer, we employ a BiLSTM as the encoder to obtain context representations. To acquire better semantic representations, we apply union fusion in the word level. In the context fusion layer, we apply union fusion, difference fusion, and similarity fusion to obtain a multi-perspective fusion representation. In the output layer, a self-attention and a feed-forward neural network are used to make the final prediction.

We conduct experiments on MRScript dataset released by [11]. Our single and ensemble model achieve the accuracy of 83.52% and 84.84% on the official test set respectively. Our main contributions are as follows:

- We propose a general fusion framework with two-layer fusion, which can fuse the passage, question, and choice simultaneously.
- To collect multi-perspective fusion representations, we define three types of fusions, consisting of union fusion, difference fusion, and similarity fusion.
- We design several groups of experiments to evaluate the effectiveness of the three types of fusion and prove that our MPFN model outperforms all the other models.

## 2  Related Work

MRC has gained significant popularity over the past few years. Several datasets have been constructed for testing the comprehension ability of a system, such as *MCTest* [15], *SQuAD* [14], *BAbI* [22], *TriviaQA* [6], RACE [8], and NewsQA [17]. Each dataset focuses on one specific aspect of reading comprehension. Particularly, the MCScript [11] dataset concerns answering the question which requires using commonsense knowledge.

Many architectures on MRC follow the process of representation, attention, fusion, and aggregation [16,24,27,5,20,25]. BiDAF [16] fuses the passage-aware question, the question-aware passage, and the original passage in context layer by concatenation, and then uses a BiLSTM for aggregation. The fusion levels in current advanced models are categorized into three types by [5] , including word-level fusion, high-level fusion, and self-boosted fusion. They further propose a FusionNet to fuse the attention information from bottom to top to obtain a fully-aware representation for answer span prediction.

On SemEval2018 Task 11, most of the models use the attention mechanism to build interactions among the passage, the question, and the choice [18,2,23,3]. The most competitive models are [18,2], and both of them employ concatenation fusion to integrate the information. [18] utilizes choice-aware passage and choice-aware question to fuse the choice in word level. In addition, they apply the question-aware passage to fuse the passage in context level. Different from [18], both the choice-aware passage and choice-aware question are fused into choice in the context level in [2] , which is the current state-of-the-art result on the MCSript dataset.

On NLI task, fusing the premise-aware hypothesis into the hypothesis is an effective and commonly-used method. [19,12] leverage the concatenation of the hypothesis and the hypothesis-aware premise to help improve the performance of their model. The element-wise subtraction and element-wise multiplication between the hypothesis and the hypothesis-aware premise are employed in [1] to enhance the concatenation.

Almost all the models on CRC only use the union fusion. In our MPFN model, we design another two fusion methods to extend the perspective of fusion. We evaluate the MPFN model on MRC task and achieve the state-of-the-art result.
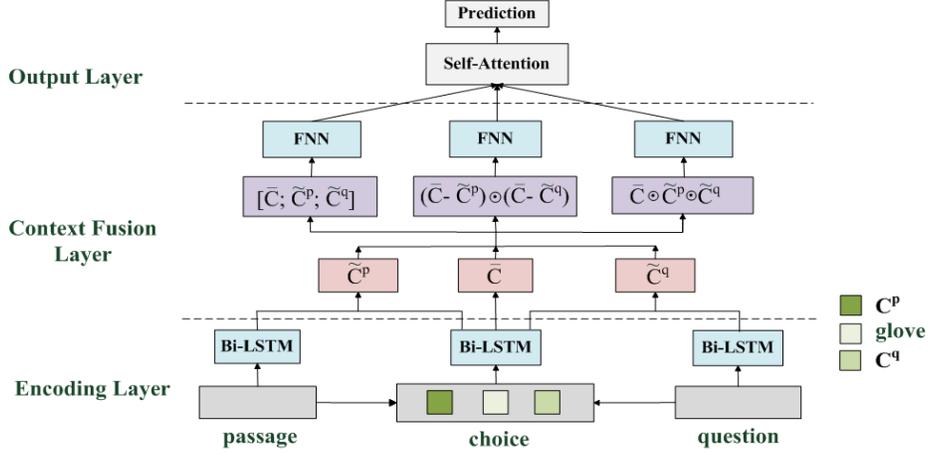
Fig. 1: Architecture of our MPFN Model.

## 3  Model

The overview of our Multi-Perspective Fusion Network (MPFN) is shown in Fig. 1. Given a narrative passage about a series of daily activities and several corresponding questions, a system requires to select a correct choice from two options for each question. In this paper, we denote $\mathbf{p} = \{\mathbf{p_1}, \mathbf{p_2}, ..., \mathbf{p_{|p|}}\}$ as the passage, $\mathbf{q} = \{\mathbf{q_1}, \mathbf{q_2}, ..., \mathbf{q_{|q|}}\}$ as a question, $\mathbf{c} = \{\mathbf{c_1}, \mathbf{c_2}, ..., \mathbf{c_{|c|}}\}$ as one of the candidate choice, and a true label $y^* \in \{0, 1\}$. Our model aims to compute a probability for each choice and take the one with higher probability as the prediction label. Our model consists of three layers: an encoding layer, a context fusion layer, and an output layer. The details of each layer are described in the following subsections.

### 3.1  Encoding Layer

This layer aims to encode the passage embedding $p$, the question embedding $q$, and the choice embedding $c$ into context embeddings. Specially, we use a one-layer BiLSTM as the context encoder.

$$\bar{c}_i = \text{BiLSTM}(c, i), \qquad\qquad i \in [1, 2, \cdots, |c|] \qquad\qquad (1)$$

$$\bar{p}_j = \text{BiLSTM}(p, j), \qquad\qquad j \in [1, 2, \cdots, |p|] \qquad\qquad (2)$$

$$\bar{q}_k = \text{BiLSTM}(q, k), \qquad\qquad k \in [1, 2, \cdots, |q|] \qquad\qquad (3)$$

The embeddings of $p$, $q$ and $c$ are semantically rich word representations consisting of several kinds of embeddings. Specifically, the embeddings of passage and question are the concatenation of the Golve word embedding, POS embedding, NER embedding, Relation embedding and Term Frequency feature. And the embeddings of choice comprise the Golve word embedding, the choice-aware passage embedding, and choice-aware question embedding . The details about each embedding are follows:

**Glove word embedding** We use the 300-dimensional Glove word embeddings trained from 840B Web crawl data [13]. The out-of-vocabulary words are initialized randomly. The embedding matrix are fixed during training.

**POS&NER embedding** We leverage the Part-of-Speech (POS) embeddings and Named-Entity Recognition(NER) embeddings. The two embeddings are randomly initialized and updated during training.

**Relation embedding** Relations are extracted form ConceptNet. For each word in the choice, if it satisfies any relation with another word in the passage or the question, the corresponding relation will be taken out. If the relations between two words are multiple, we just randomly choose one. The relation embeddings are generated in the similar way of POS embeddings.

**Term Frequency** Following [18], we introduce the term frequency feature to enrich the embedding of each word. The calculation is based on English Wikipedia.

**Choice-aware passage embedding** The information in the passage that is relevant to the choice can help encode the choice [21]. To acquire the choice-aware passage embedding $c_i^p$, we utilize dot product between non-linear mappings of word embeddings to compute the attention scores for the passage [9].

$$c_i^p = Attn(c_i, \{p_j\}_1^{|p|}) = \sum_{j=1}^{|p|} \alpha_{ij} p_j \tag{4}$$

$$\alpha_{ij} \propto exp(S(c_i, p_j)), \quad S(c_i, p_j) = ReLU(Wc_i)^T ReLU(Wp_j) \tag{5}$$

**Choice-aware question embedding** The choice relevant question information is also important for the choice. Therefore, we adopt the similar attention way as above to get the choice-aware question embedding $c_i^q = Attn(c_i, \{q_k\}_1^{|q|})$.

The embeddings delivered to the BiLSTM are the concatenation the above components, where $p_j = [p_j^{glove}, p_j^{pos}, p_j^{ner}, p_j^{rel}, p_j^{tf}]$, $c_i = [c_i^{glove}, c_i^p, c_i^q]$, and $q_k = [q_k^{glove}, q_k^{pos}, q_k^{ner}, q_k^{rel}, q_k^{tf}]$.

### 3.2   Context Fusion Layer

This is the core layer of our MPFN model. In this layer, we define three fusion functions, which consider the union information, the different information, and the similar information of the choice, passage, and question.

Since we have obtained the choice context $\bar{c}_i$, the passage context $\bar{p}_j$, and the question context $\bar{q}_k$ in the encoding layer, we can calculate the choice-aware passage contexts $\tilde{c}_i^p$ and choice-aware question contexts $\tilde{c}_i^q$. Then we deliver them together with the choice contexts $\bar{c}_i$ to the three fusion functions.

**Choice-aware passage context** In this part, we calculate the choice-aware passage representations $\tilde{c}_i^p = \sum_j \beta_{ij} \bar{p}_j$. For model simplification, here we use dot product between choice contexts and passage contexts to compute the attention scores $\beta_{ij}$:

$$\beta_{ij} = \frac{exp(\bar{c}_i^T \bar{p}_j)}{\sum_{j'=1}^{|p|} exp(\bar{c}_i^T \bar{p}_{j'})} \tag{6}$$

**Choice-aware question context** In a similar way as above, we get the choice-aware question context $\tilde{c}_i^q = \sum_j \beta_{ik}\bar{q}_k$. The $\beta_{ik}$ is the dot product of the choice context $\bar{c}_i$ and question context $\bar{q}_k$.

**Multi-perspective Fusion** This is the key module in our MPFN model. The goal of this part is to produce multi-perspective fusion representation for the choice $\bar{c}_i$, the choice-aware passage $\tilde{c}_i^p$, and the choice-aware question $\tilde{c}_i^q$. In this paper, we define fusion in three perspectives: *union*, *difference*, and *similarity*. Accordingly, we define three fusion functions to describe the three perspectives. The outputs and calculation of the three functions are as follows:

$$u_i = [\bar{c}_i \, ; \tilde{c}_i^p \, ; \tilde{c}_i^q], \tag{7}$$

$$d_i = (\bar{c}_i - \tilde{c}_i^p) \odot (\bar{c}_i - \tilde{c}_i^q), \tag{8}$$

$$s_i = \bar{c}_i \odot \tilde{c}_i^p \odot \tilde{c}_i^q, \tag{9}$$

where $;$ , $-$, and $\odot$ represent concatenation, element-wise subtraction, and element-wise multiplication respectively. And $u_i$, $d_i$, and $s_i$ are the representations from the union, difference and similarity perspective respectively.

The union perspective is commonly used in a large bulk of tasks [12,5,24]. It can see the whole picture of the passage, the question, and the choice by concatenating the $\tilde{c}_i^p$ and $\tilde{c}_i^q$ together with $c_i$ . While the difference perspective captures the different parts between choice and passage, and the difference parts between choice and question by $\bar{c}_i - \tilde{c}_i^p$ and $\bar{c}_i - \tilde{c}_i^q$ respectively. The $\odot$ in difference perspective can detect the two different parts at the same time and emphasize them. In addition, the similarity perspective is capable of discovering the similar parts among the passage, the question, and the choice.

To map the three fusion representations to lower and same dimension, we apply three different FNNs with the ReLU activation to $u_i$, $d_i$, and $s_i$. The final output $g_i$ is the concatenation of the results of the three FNNs, which represents a global perspective representation.

$$g_i = [f^u(u_i), f^d(d_i), f^s(s_i)] \tag{10}$$

### 3.3   Output Layer

The output layer includes a self-attention layer and a prediction layer. Following [26], we summarize the global perspective representation $\{g_i\}_1^{|c|}$ to a fixed length vector $r$. We compute the $r = \sum_{i=1}^{|c|} b_i g_i$, where $b_j$ is the self-weighted attention score :

$$b_i = \frac{exp(Wg_i)}{\sum_{i'=1}^{|c|} exp(Wg_{i'})} \tag{11}$$

In the prediction layer, we utilize the output of self-attention $r$ to make the final prediction.

## 4   Experiments

### 4.1   Experimental Settings

**Data** We conduct experiments on the MCScript [11], which is used as the official dataset of SemEval2018 Task11. This dataset constructs a collection of text passages

| Model | Test (%acc) |
|---|---|
| SLQA | 79.94 |
| Rusalka | 80.48 |
| HMA Model (single) [2] | 80.94 |
| TriAN (single) [18] | 81.94 |
| **MPFN** (single) | **83.52** |
| (jiangnan) (ensemble) [23] | 80.91 |
| MITRE (ensemble) [3] | 82.27 |
| TriAN (ensemble) [18] | 83.95 |
| HMA Model (ensemble) [2] | 84.13 |
| **MPFN** (ensemble) | **84.84** |

Table 2: Experimental Results of Models

about daily life activities and a series of questions referring to each passage, and each question is equipped with two answer choices. The MCScript comprises 9731, 1411, and 2797 questions in training, development, and test set respectively. For data preprocessing, we use spaCy [1] for sentence tokenization, Part-of-Speech tagging, and Name Entity Recognization. The relations between two words are generated by ConceptNet.

**Parameters** We use the standard cross-entropy function as the loss function. We choose Adam [7] with initial momentums for parameter optimization. As for hyper-parameters, we set the batch size as 32, the learning rate as 0.001, the dimension of BiLSTM and the hidden layer of FNN as 123. The embedding size of Glove, NER, POS, Relation are 300, 8, 12, 10 respectively. The dropout rate of the word embedding and BiLSTM output are 0.386 and 0.40 respectively.

### 4.2 Experimental Results

Table2 shows the results of our MPFN model along with the competitive models on the MCScript dataset. The TriAN achieves 81.94% in terms of test accuracy, which is the best result of the single model. The best performing ensemble result is 84.13%, provided by HMA, which is the voting results of 7 single systems.

Our single MPFN model achieves 83.52% in terms of accuracy, outperforming all the previous models. The model exceeds the HMA and TriAN by approximately 2.58% and 1.58% absolute respectively. Our ensemble model surpasses the current state-of-the-art model with an accuracy of 84.84%. We got the final ensemble result by voting on 4 single models. Every single model uses the same architecture but different parameters.

### 4.3 Discussion of Multi-Perspective

To study the effectiveness of each perspective, we conduct several experiments on the three single perspectives and their combination perspective. Table 3 presents their comparison results. The first group of models are based on the three single perspectives,

---

[1] https://github.com/explosion/spaCy

| Perspective | MPFN | MPFN+BiLSTM |
|---|---|---|
| U | 82.73 | 82.73 |
| D | 82.27 | 81.77 |
| S | 81.55 | 80.59 |
| DU | 82.84 | 82.16 |
| SU | 82.48 | 82.87 |
| SD | 83.12 | 83.09 |
| SDU | **83.52** | 82.70 |

Table 3: Test Accuracy of Multi-Perspective

| Model | Test (%acc) |
|---|---|
| **MPFN** | **83.52** |
| w/o POS | 82.70 |
| w/o NER | 82.62 |
| w/o Rel | 81.98 |
| w/o TF | 81.91 |
| w/o $C^p$ | 81.62 |
| w/o $C^q$ | 82.16 |
| w/o $C^p$&$C^q$ | 81.66 |

Table 4: Encoding Inputs Ablation Study.

and we can observe that the union perspective performs best compared with the difference and similarity perspective. Moreover, the union perspective achieves 82.73% in accuracy, exceeding the TriAN by 0.79% absolute. We can also see that the similarity perspective is inferior to the other two perspectives.

The second group of models are formed from two perspectives. Compared with the single union perspective, combining the difference perspective with the union perspective can improve 0.11%. Composing union and similarity fusion together doesn't help the training. To our surprise, the combination of similarity perspective and difference perspective obtains 83.09% accuracy score.

The last model is our MPFN model, which performing best. The final result indicates that composing the union perspective, difference perspective, and similarity perspective together to train is helpful.

Many advanced models employ a BiLSTM to further aggregate the fusion results. To investigate whether a BiLSTM can assist the model, we apply another BiLSTM to the three fusion representations in Formula 10 respectively and then put them together. The results are shown in the second column in Table 3, which indicate that the BiLSTM does not help improve the performance of the models.

### 4.4   Encoding Inputs Ablation

In the section, we conduct ablation study on the encoding inputs to examine the effectiveness each component. The experiment results are listed in Table 4.

From the best model, if we remove the POS embedding and NER embedding, the accuracy drops by 0.82% and 0.9%. Without Relation embedding, the accuracy drops to 81.98%, revealing that the external relations are helpful to the context fusions. Without Term Frequency, the accuracy drops by approximately 1.61%. This behavior suggests that the Term Frequency feature has a powerful capability to guide the model.

After removing the $C^p$, we find the performance degrades to 81.62%. This demonstrates that information in the passage is significantly important to final performance. If we remove $C^q$ from the MPFN, the accuracy drops to 82.16%. If we remove the word level fusion completely, we will obtain an 81.66% accuracy score. These results
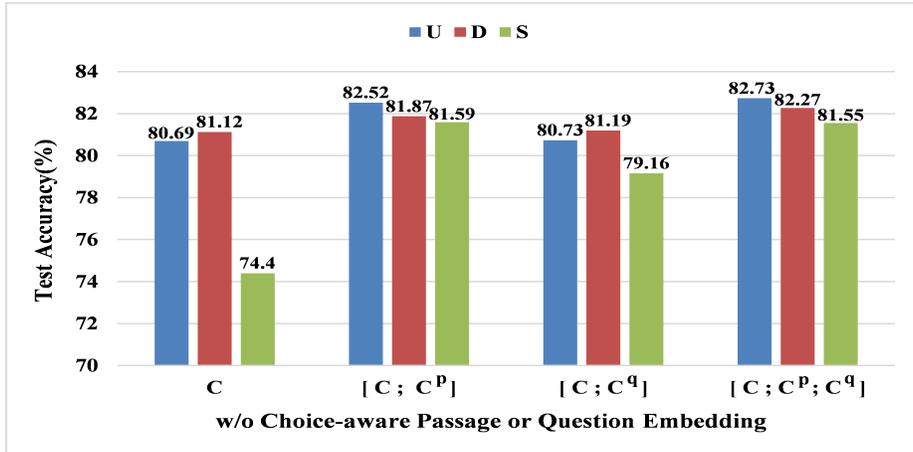
Fig. 2: Influence of Word-level Interaction.

demonstrate that each component is indispensable and the bottom embeddings are the basic foundations of the top layer fusions.

### 4.5 Influence of Word-level Interaction

In this section, we explore the influence of word-level interaction to each perspective. Fig 2 reports the overall results of how each perspective can be affected by the lower level interaction. The $C^p$ and the $C^q$ represent the choice-aware passage embedding and the choice-aware question embedding respectively. We can observe that the results of $[C; C^p]$, $[C; C^q]$, and $[C; C^p; C^q]$ are all higher than the result of $C$ alone, indicating the effectiveness of word embedding interaction.

Both the union fusion and difference fusion can achieve more than 80% accuracy, while the similarity fusion is very unstable. We also observe that the difference fusion is comparable with the union fusion, which even works better than the union fusion when the information of $C^p$ is not introduced into the input of encoding. The similarity fusion performs poorly in $C$ and $[C; C^q]$, while yielding a huge increase in the remaining two groups of experiments, which is an interesting phenomenon. We infer that the similarity fusion needs to be activated by the union fusion.

In summary, we can conclude that integrate the information of $C^p$ into $C$ can greatly improve the performance of the model. Combining $C^q$ together with $C^p$ can further increase the accuracy.

### 4.6 Visualization

In this section, we visualize the union and difference fusion representations and show them in Fig 3. And, we try to analyze their characteristics and compare them to discover some connections. The values of similarity fusion are too small to observe useful information intuitively, so we do not show it here. We use the example presented in Table 1

for visualization, where the question is *Why didn't the child go to bed by themselves?* and the corresponding True choice is *The child wanted to continue playing.*
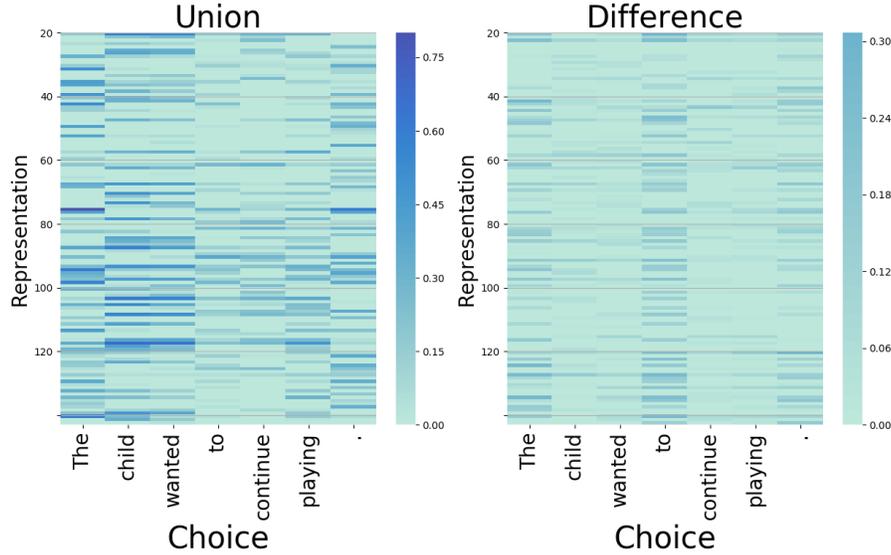


Fig. 3: Visualization of Fusions

The left region in Fig 3 is the union fusion. The most intuitive observation is that it captures comprehensive information. The values of *child*, *wanted*, *playing* are obvious higher than other words. This is consistent with our prior cognition, because the concatenation operation adopted in union fusion does not lose any content. While the difference union shows in the right region in Fig 3 focuses on some specific words. By further comparison, we find that the difference fusion can pay attention to the content ignored by the union fusion. What's more, the content acquired by the union would not be focused by the difference again. In other words, the union fusion and difference fusion indeed can emphasize information from the different perspective.

## 5 Conclusion

In this paper, we propose the Multi-Perspective Fusion Network (MPFN) for the Commonsense Reading Comprehension (CMC) task. We propose a more general framework for CRC by designing the difference and similarity fusion to assist the union fusion. Our MPFN model achieves an accuracy of 83.52% on MCScript, outperforming the previous models. The experimental results show that union fusion based on the choice-aware passage, the choice-aware question, and the choice can surpass the TriAN and HMA model. The difference fusion performs stably, which is comparable with the union fusion. We find that the word-level union fusion can significantly influence the

context-level fusion. The choice-aware passage word embedding can activate the similarity fusion. We find that combining the similar parts and the difference parts together can obtain the best performance among the two-perspective models. By taking the three types of fusion methods into consideration, our MPFN model achieves a state-of-the-art result.

## Acknowledgements

## References

1. Chen, Q., Zhu, X., Ling, Z.H., Wei, S., Jiang, H., Inkpen, D.: Enhanced lstm for natural language inference. In: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 1657–1668. Association for Computational Linguistics, Vancouver, Canada (July 2017), http://aclweb.org/anthology/P17-1152

2. Chen, Z., Cui, Y., Ma, W., Wang, S., Liu, T., Hu, G.: Hfl-rc system at semeval-2018 task 11: Hybrid multi-aspects model for commonsense reading comprehension. arXiv preprint arXiv:1803.05655 (2018)

3. Elizabeth M. Merkhofer, John Henderson, D.B.L.S., Zarrella, G.: Mitre at semeval-2018 task 11: Commonsense reasoning without commonsense knowledge (2018)

4. Hochreiter, S., Schmidhuber, J.: Long short-term memory. Neural Comput. **9**(8), 1735–1780 (Nov 1997). https://doi.org/10.1162/neco.1997.9.8.1735, http://dx.doi.org/10.1162/neco.1997.9.8.1735

5. Huang, H., Zhu, C., Shen, Y., Chen, W.: Fusionnet: Fusing via fully-aware attention with application to machine comprehension. CoRR **abs/1711.07341** (2017)

6. Joshi, M., Choi, E., Weld, D.S., Zettlemoyer, L.: Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. In: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics. Association for Computational Linguistics, Vancouver, Canada (July 2017)

7. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. CoRR **abs/1412.6980** (2014)

8. Lai, G., Xie, Q., Liu, H., Yang, Y., Hovy, E.H.: Race: Large-scale reading comprehension dataset from examinations. In: EMNLP (2017)

9. Lee, K., Kwiatkowski, T., Parikh, A.P., Das, D.: Learning recurrent span representations for extractive question answering. CoRR **abs/1611.01436** (2016), http://arxiv.org/abs/1611.01436

10. Mou, L., Men, R., Li, G., Xu, Y., Zhang, L., Yan, R., Jin, Z.: Natural language inference by tree-based convolution and heuristic matching. In: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers). pp. 130–136. Association for Computational Linguistics, Berlin, Germany (August 2016), http://anthology.aclweb.org/P16-2022

11. Ostermann, S., Modi, A., Roth, M., Thater, S., Pinkal, M.: MCScript: A Novel Dataset for Assessing Machine Comprehension Using Script Knowledge. In: Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018). European Language Resources Association (ELRA), Miyazaki, Japan (May 7-12, 2018 2018)

12. Parikh, A., Täckström, O., Das, D., Uszkoreit, J.: A decomposable attention model for natural language inference. In: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing. pp. 2249–2255. Association for Computational Linguistics (2016). https://doi.org/10.18653/v1/D16-1244, http://www.aclweb.org/anthology/D16-1244

13. Pennington, J., Socher, R., Manning, C.D.: Glove: Global vectors for word representation. In: In EMNLP (2014)

14. Rajpurkar, P., Zhang, J., Lopyrev, K., Liang, P.: Squad: 100, 000+ questions for machine comprehension of text. CoRR **abs/1606.05250** (2016)

15. Richardson, M., Burges, C.J.C., Renshaw, E.: Mctest: A challenge dataset for the open-domain machine comprehension of text. In: EMNLP (2013)

16. Seo, M.J., Kembhavi, A., Farhadi, A., Hajishirzi, H.: Bidirectional attention flow for machine comprehension. CoRR **abs/1611.01603** (2016)

17. Trischler, A., Wang, T., Yuan, X., Harris, J., Sordoni, A., Bachman, P., Suleman, K.: Newsqa: A machine comprehension dataset. In: Rep4NLP@ACL (2017)

18. Wang, L., Sun, M., Zhao, W., Shen, K., Liu, J.: Yuanfudao at semeval-2018 task 11: Three-way attention and relational knowledge for commonsense machine comprehension. In: SemEval@NAACL-HLT. pp. 758–762. Association for Computational Linguistics (2018)

19. Wang, S., Jiang, J.: Learning natural language inference with lstm. In: Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. pp. 1442–1451. Association for Computational Linguistics, San Diego, California (June 2016), http://www.aclweb.org/anthology/N16-1170

20. Wang, W., Yang, N., Wei, F., Chang, B., Zhou, M.: Gated self-matching networks for reading comprehension and question answering. In: Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). pp. 189–198. Association for Computational Linguistics (2017). https://doi.org/10.18653/v1/P17-1018, http://www.aclweb.org/anthology/P17-1018

21. Weissenborn, D., Wiese, G., Seiffe, L.: Fastqa: A simple and efficient neural architecture for question answering. CoRR **abs/1703.04816** (2017), http://arxiv.org/abs/1703.04816

22. Weston, J., Bordes, A., Chopra, S., Mikolov, T.: Towards ai-complete question answering: A set of prerequisite toy tasks. CoRR **abs/1502.05698** (2015)

23. Xia, J.: Jiangnan at semeval-2018 task 11: Deep neural network with attention method for machine comprehension task (2018)

24. Xiong, C., Zhong, V., Socher, R.: DCN+: Mixed objective and deep residual coattention for question answering. In: International Conference on Learning Representations (2018), https://openreview.net/forum?id=H1meywxRW

25. Xu, Y., Liu, J., Gao, J., Shen, Y., Liu, X.: Towards human-level machine reading comprehension: Reasoning and inference with multiple strategies. CoRR **abs/1711.04964** (2017)

26. Yang, Z., Yang, D., Dyer, C., He, X., Smola, A., Hovy, E.: Hierarchical attention networks for document classification. In: NAACL. pp. 1480–1489. Association for Computational Linguistics, San Diego, California (June 2016), http://www.aclweb.org/anthology/N16-1174

27. Zhu, H., Wei, F., Qin, B., Liu, T.: Hierarchical attention flow for multiple-choice reading comprehension. In: AAAI (2018)