

# Bilateral Multi-Perspective Matching for Natural Language Sentences

Zhiguo Wang, Wael Hamza, Radu Florian

IBM T.J. Watson Research Center

{zhigwang,whamza,raduf}@us.ibm.com

## Abstract

Natural language sentence matching is a fundamental technology for a variety of tasks. Previous approaches either match sentences from a single direction or only apply single granular (word-by-word or sentence-by-sentence) matching. In this work, we propose a bilateral multi-perspective matching (BiMPM) model. Given two sentences  $P$  and  $Q$ , our model first encodes them with a BiLSTM encoder. Next, we match the two encoded sentences in two directions  $P$  against  $Q$  and  $Q$  against  $P$ . In each matching direction, each time step of one sentence is matched against all time-steps of the other sentence from multiple perspectives. Then, another BiLSTM layer is utilized to aggregate the matching results into a fixed-length matching vector. Finally, based on the matching vector, a decision is made through a fully connected layer. We evaluate our model on three tasks: paraphrase identification, natural language inference and answer sentence selection. Experimental results on standard benchmark datasets show that our model achieves the state-of-the-art performance on all tasks.

## 1 Introduction

Natural language sentence matching (NLSM) is the task of comparing two sentences and identifying the relationship between them. It is a fundamental technology for a variety of tasks. For example, in a paraphrase identification task, NLSM is used to determine whether two sentences are paraphrase or not [Yin *et al.*, 2015]. For a natural language inference task, NLSM is utilized to judge whether a hypothesis sentence can be inferred from a premise sentence [Bowman *et al.*, 2015]. For question answering and information retrieval tasks, NLSM is employed to assess the relevance between query-answer pairs and rank all the candidate answers [Wang *et al.*, 2016d]. For machine comprehension tasks, NLSM is used for matching a passage with a question and pointing out the correct answer span [Wang *et al.*, 2016b].

With the renaissance of neural network models [LeCun *et al.*, 2015; Peng *et al.*, 2015a; Peng *et al.*, 2016], two types of deep learning frameworks were proposed for NLSM.

The first framework is based on the “Siamese” architecture [Bromley *et al.*, 1993]. In this framework, the same neural network encoder (e.g., a CNN or a RNN) is applied to two input sentences individually, so that both of the two sentences are encoded into sentence vectors in the same embedding space. Then, a matching decision is made solely based on the two sentence vectors [Bowman *et al.*, 2015; Tan *et al.*, 2015]. The advantage of this framework is that sharing parameters makes the model smaller and easier to train, and the sentence vectors can be used for visualization, sentence clustering and many other purposes [Wang *et al.*, 2016c]. However, a disadvantage is that there is no explicit interaction between the two sentences during the encoding procedure, which may lose some important information. To deal with this problem, a second framework “matching-aggregation” has been proposed [Wang and Jiang, 2016; Wang *et al.*, 2016d]. Under this framework, smaller units (such as words or contextual vectors) of the two sentences are firstly matched, and then the matching results are aggregated (by a CNN or a LSTM) into a vector to make the final decision. The new framework captures more interactive features between the two sentences, therefore it acquires significant improvements. However, the previous “matching-aggregation” approaches still have some limitations. First, some of the approaches only explored the word-by-word matching [Rocktäschel *et al.*, 2015], but ignored other granular matchings (e.g., phrase-by-sentence); Second, the matching is only performed in a single direction (e.g., matching  $P$  against  $Q$ ) [Wang and Jiang, 2015], but neglected the reverse direction (e.g., matching  $Q$  against  $P$ ).

In this paper, to tackle these limitations, we propose a bilateral multi-perspective matching (BiMPM) model for NLSM tasks. Our model essentially belongs to the “matching-aggregation” framework. Given two sentences  $P$  and  $Q$ , our model first encodes them with a bidirectional Long Short-Term Memory Network (BiLSTM). Next, we match the two encoded sentences in two directions  $P \rightarrow Q$  and  $P \leftarrow Q$ . In each matching direction, let’s say  $P \rightarrow Q$ , each time step of  $Q$  is matched against all time-steps of  $P$  from multiple perspectives. Then, another BiLSTM layer is utilized to aggregate the matching results into a fixed-length matching vector. Finally, based on the matching vector, a decision is made through a fully connected layer. We evaluate our model on three NLSM tasks: paraphrase identification, natural lan-

guage inference and answer sentence selection. Experimental results on standard benchmark datasets show that our model achieves the state-of-the-art performance on all tasks.

In following parts, we start with a brief definition of the NLSM task (Section 2), followed by the details of our model (Section 3). Then we evaluate our model on standard benchmark datasets (Section 4). We talk about related work in Section 5, and conclude this work in Section 6.

## 2 Task Definition

Formally, we can represent each example of the NLSM task as a triple  $(P, Q, y)$ , where  $P = (p_1, \dots, p_j, \dots, p_M)$  is a sentence with a length  $M$ ,  $Q = (q_1, \dots, q_i, \dots, q_N)$  is the second sentence with a length  $N$ ,  $y \in \mathcal{Y}$  is the label representing the relationship between  $P$  and  $Q$ , and  $\mathcal{Y}$  is a set of task-specific labels. The NLSM task can be represented as estimating a conditional probability  $\Pr(y|P, Q)$  based on the training set, and predicting the relationship for testing examples by  $y^* = \arg \max_{y \in \mathcal{Y}} \Pr(y|P, Q)$ . Concretely, for a paraphrase identification task,  $P$  and  $Q$  are two sentences,  $\mathcal{Y} = \{0, 1\}$ , where  $y = 1$  means that  $P$  and  $Q$  are paraphrase of each other, and  $y = 0$  otherwise. For a natural language inference task,  $P$  is a premise sentence,  $Q$  is a hypothesis sentence, and  $\mathcal{Y} = \{entailment, contradiction, neutral\}$  where *entailment* indicates  $Q$  can be inferred from  $P$ , *contradiction* indicates  $Q$  cannot be true condition on  $P$ , and *neutral* means  $P$  and  $Q$  are irrelevant to each other. In an answer sentence selection task,  $P$  is a question,  $Q$  is a candidate answer, and  $\mathcal{Y} = \{0, 1\}$  where  $y = 1$  means  $Q$  is a correct answer for  $P$ , and  $y = 0$  otherwise.

## 3 Method

In this section, we first give a high-level overview of our model in Sub-section 3.1, and then give more details about our novel multi-perspective matching operation in Sub-section 3.2.

### 3.1 Model Overview

We propose a bilateral multi-perspective matching (BiMPM) model to estimate the probability distribution  $\Pr(y|P, Q)$ . Our model belongs to the “matching-aggregation” framework [Wang and Jiang, 2016]. Contrarily to previous “matching-aggregation” approaches, our model matches  $P$  and  $Q$  in two directions ( $P \rightarrow Q$  and  $P \leftarrow Q$ ). In each individual direction, our model matches the two sentences from multiple perspectives. Figure 1 shows the architecture of our model. Given a pair of sentences  $P$  and  $Q$ , the BiMPM model estimates the probability distribution  $\Pr(y|P, Q)$  through the following five layers.

**Word Representation Layer.** The goal of this layer is to represent each word in  $P$  and  $Q$  with a  $d$ -dimensional vector. We construct the  $d$ -dimensional vector with two components: a word embedding and a character-composed embedding. **The word embedding is a fixed vector** for each individual word, which is pre-trained with GloVe [Pennington et al., 2014] or word2vec [Mikolov et al., 2013]. **The character-composed embedding is calculated by feeding each character** (represented as a character embedding) **within a word**

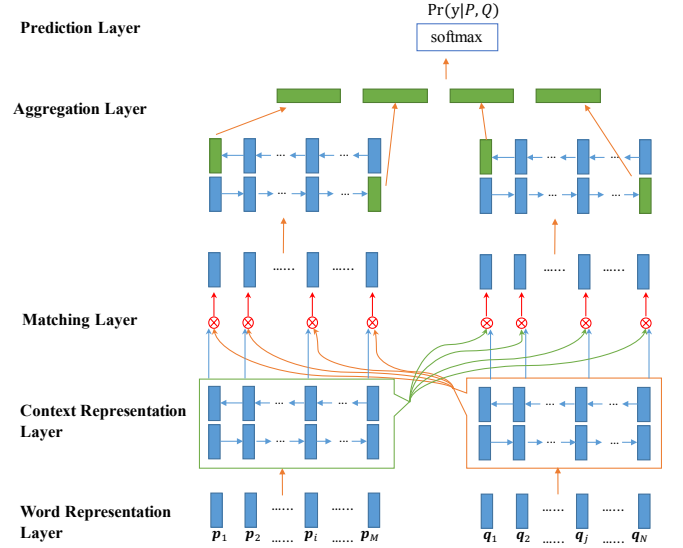


Figure 1: Architecture for Bilateral Multi-Perspective Matching (BiMPM) Model, where  $\otimes$  is the multi-perspective matching operation described in sub-section 3.2.

into a Long Short-Term Memory Network (LSTM) [Hochreiter and Schmidhuber, 1997], where the character embeddings are randomly initialized and learned jointly with other network parameters from NLSM tasks. The output of this layer are two sequences of word vectors  $P : [p_1, \dots, p_M]$  and  $Q : [q_1, \dots, q_N]$ .

**Context Representation Layer.** The purpose of this layer is to incorporate contextual information into the representation of each time step of  $P$  and  $Q$ . We utilize a **bi-directional LSTM (BiLSTM)** to encode contextual embeddings for each time-step of  $P$ .

$$\begin{aligned} \vec{h}_i^p &= \overrightarrow{\text{LSTM}}(\vec{h}_{i-1}^p, p_i) & i = 1, \dots, M \\ \overleftarrow{h}_i^p &= \overleftarrow{\text{LSTM}}(\overleftarrow{h}_{i+1}^p, p_i) & i = M, \dots, 1 \end{aligned} \quad (1)$$

Meanwhile, we apply **the same BiLSTM** to encode  $Q$ :

$$\begin{aligned} \vec{h}_j^q &= \overrightarrow{\text{LSTM}}(\vec{h}_{j-1}^q, q_j) & j = 1, \dots, N \\ \overleftarrow{h}_j^q &= \overleftarrow{\text{LSTM}}(\overleftarrow{h}_{j+1}^q, q_j) & j = N, \dots, 1 \end{aligned} \quad (2)$$

**Matching Layer.** This is the core layer within our model. The goal of this layer is to compare each contextual embedding (time-step) of one sentence against all contextual embeddings (time-steps) of the other sentence. As shown in Figure 1, we will match the two sentences  $P$  and  $Q$  in two directions: match each time-step of  $P$  against all time-steps of  $Q$ , and match each time-step of  $Q$  against all time-steps of  $P$ . To match one time-step of a sentence against all time-steps of the other sentence, we design a multi-perspective matching operation  $\otimes$ . We will give more details about this operation in Sub-section 3.2. The output of this layer are two sequences of matching vectors (right above the operation  $\otimes$  in Figure 1), where each matching vector corresponds to the matching result of one time-step against all time-steps of the other sentence.

**Aggregation Layer.** This layer is employed to aggregate the two sequences of matching vectors into a fixed-length matching vector. We utilize another BiLSTM model, and apply it to the two sequences of matching vectors individually. Then, we construct the fixed-length matching vector by concatenating (the four green) vectors from the last time-step of the BiLSTM models.

**Prediction Layer.** The purpose of this layer is to evaluate the probability distribution  $\Pr(y|P, Q)$ . To this end, we employ a two layer feed-forward neural network to consume the fixed-length matching vector, and apply the *softmax* function in the output layer. The number of nodes in the output layer is set based on each specific task described in Section 2.

### 3.2 Multi-perspective Matching Operation

We define the multi-perspective matching operation  $\otimes$  in following two steps:

**First**, we define a multi-perspective cosine matching function  $f_m$  to compare two vectors

$$\mathbf{m} = f_m(\mathbf{v}_1, \mathbf{v}_2; \mathbf{W}) \quad (3)$$

where  $\mathbf{v}_1$  and  $\mathbf{v}_2$  are two  $d$ -dimensional vectors,  $\mathbf{W} \in \mathbb{R}^{l \times d}$  is a trainable parameter with the shape  $l \times d$ ,  $l$  is the number of perspectives, and the returned value  $\mathbf{m}$  is a  $l$ -dimensional vector  $\mathbf{m} = [m_1, \dots, m_k, \dots, m_l]$ . Each element  $m_k \in \mathbf{m}$  is a matching value from the  $k$ -th perspective, and it is calculated by the cosine similarity between two weighted vectors

$$m_k = \text{cosine}(W_k \circ \mathbf{v}_1, W_k \circ \mathbf{v}_2) \quad (4)$$

where  $\circ$  is the element-wise multiplication, and  $W_k$  is the  $k$ -th row of  $\mathbf{W}$ , which controls the  $k$ -th perspective and assigns different weights to different dimensions of the  $d$ -dimensional space.

**Second**, based on  $f_m$ , we define four matching strategies to compare each time-step of one sentence against all time-steps of the other sentence. To avoid repetition, we only define these matching strategies for one matching direction  $P \rightarrow Q$ . The readers can infer equations for the reverse direction easily.

(1) **Full-Matching.** Figure 2 (a) shows the diagram of this matching strategy. In this strategy, each forward (or backward) contextual embedding  $\vec{\mathbf{h}}_i^p$  (or  $\overleftarrow{\mathbf{h}}_i^p$ ) is compared with the last time step of the forward (or backward) representation of the other sentence  $\vec{\mathbf{h}}_N^q$  (or  $\overleftarrow{\mathbf{h}}_1^q$ ).

$$\begin{aligned} \vec{\mathbf{m}}_i^{\text{full}} &= f_m(\vec{\mathbf{h}}_i^p, \vec{\mathbf{h}}_N^q; \mathbf{W}^1) \\ \overleftarrow{\mathbf{m}}_i^{\text{full}} &= f_m(\overleftarrow{\mathbf{h}}_i^p, \overleftarrow{\mathbf{h}}_1^q; \mathbf{W}^2) \end{aligned} \quad (5)$$

(2) **Maxpooling-Matching.** Figure 2 (b) gives the diagram of this matching strategy. In this strategy, each forward (or backward) contextual embedding  $\vec{\mathbf{h}}_i^p$  (or  $\overleftarrow{\mathbf{h}}_i^p$ ) is compared with every forward (or backward) contextual embeddings of the other sentence  $\vec{\mathbf{h}}_j^q$  (or  $\overleftarrow{\mathbf{h}}_j^q$ ) for  $j \in (1 \dots N)$ , and only the maximum value of each dimension is retained.

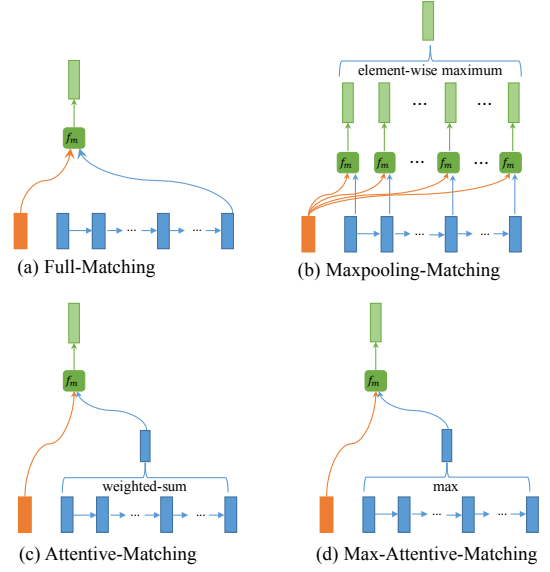


Figure 2: Diagrams for different matching strategies, where  $f_m$  is the multi-perspective cosine matching function in Eq.(3), the input includes one time step of one sentence (left orange block) and all the time-steps of the other sentence (right blue blocks), and the output is a vector of matching values (top green block) calculated by Eq.(3).

$$\vec{\mathbf{m}}_i^{\text{max}} = \max_{j \in (1 \dots N)} f_m(\vec{\mathbf{h}}_i^p, \vec{\mathbf{h}}_j^q; \mathbf{W}^3)$$

$$\overleftarrow{\mathbf{m}}_i^{\text{max}} = \max_{j \in (1 \dots N)} f_m(\overleftarrow{\mathbf{h}}_i^p, \overleftarrow{\mathbf{h}}_j^q; \mathbf{W}^4) \quad (6)$$

where  $\max_{j \in (1 \dots N)}$  is element-wise maximum.

(3) **Attentive-Matching.** Figure 2 (c) shows the diagram of this matching strategy. We first calculate the cosine similarities between each forward (or backward) contextual embedding  $\vec{\mathbf{h}}_i^p$  (or  $\overleftarrow{\mathbf{h}}_i^p$ ) and every forward (or backward) contextual embeddings of the other sentence  $\vec{\mathbf{h}}_j^q$  (or  $\overleftarrow{\mathbf{h}}_j^q$ ):

$$\begin{aligned} \vec{\alpha}_{i,j} &= \text{cosine}(\vec{\mathbf{h}}_i^p, \vec{\mathbf{h}}_j^q) & j = 1, \dots, N \\ \overleftarrow{\alpha}_{i,j} &= \text{cosine}(\overleftarrow{\mathbf{h}}_i^p, \overleftarrow{\mathbf{h}}_j^q) & j = 1, \dots, N \end{aligned} \quad (7)$$

Then, we take  $\vec{\alpha}_{i,j}$  (or  $\overleftarrow{\alpha}_{i,j}$ ) as the weight of  $\vec{\mathbf{h}}_j^q$  (or  $\overleftarrow{\mathbf{h}}_j^q$ ), and calculate an attentive vector for the entire sentence  $Q$  by weighted summing all the contextual embeddings of  $Q$ :

$$\begin{aligned} \vec{\mathbf{h}}_i^{\text{mean}} &= \frac{\sum_{j=1}^N \vec{\alpha}_{i,j} \cdot \vec{\mathbf{h}}_j^q}{\sum_{j=1}^N \vec{\alpha}_{i,j}} \\ \overleftarrow{\mathbf{h}}_i^{\text{mean}} &= \frac{\sum_{j=1}^N \overleftarrow{\alpha}_{i,j} \cdot \overleftarrow{\mathbf{h}}_j^q}{\sum_{j=1}^N \overleftarrow{\alpha}_{i,j}} \end{aligned} \quad (8)$$

Finally, we match each forward (or backward) contextual embedding of  $\vec{h}_i^p$  (or  $\overleftarrow{h}_i^p$ ) with its corresponding attentive vector:

$$\begin{aligned}\vec{m}_i^{att} &= f_m(\vec{h}_i^p, \vec{h}_i^{mean}; \mathbf{W}^5) \\ \overleftarrow{m}_i^{att} &= f_m(\overleftarrow{h}_i^p, \overleftarrow{h}_i^{mean}; \mathbf{W}^6)\end{aligned}\quad (9)$$

(4) **Max-Attentive-Matching**. Figure 2 (d) shows the diagram of this matching strategy. This strategy is similar to the Attentive-Matching strategy. However, instead of taking the weighed sum of all the contextual embeddings as the attentive vector, we **pick the contextual embedding with the highest cosine similarity as the attentive vector**. Then, we match each contextual embedding of the sentence  $P$  with its new attentive vector.

We apply all these four matching strategies to each time-step of the sentence  $P$ , and **concatenate the generated eight vectors as the matching vector for each time-step of  $P$** . We also perform the same process for the reverse matching direction.

## 4 Experiments

In this section, we evaluate our model on three tasks: paraphrase identification, natural language inference and answer sentence selection. We will first introduce the general setting of our BiMPM models in Sub-section 4.1. Then, we demonstrate the properties of our model through some ablation studies in Sub-section 4.2. Finally, we compare our model with state-of-the-art models on some standard benchmark datasets in Sub-section 4.3, 4.4 and 4.5.

### 4.1 Experiment Settings

We initialize word embeddings in the word representation layer with the **300-dimensional GloVe** word vectors pre-trained from the 840B Common Crawl corpus [Pennington *et al.*, 2014]. For the out-of-vocabulary (OOV) words, we initialize the word embeddings randomly. For the character-composed embeddings, we initialize **each character as a 20-dimensional vector**, and **compose each word into a 50-dimensional vector with a LSTM layer**. We set **the hidden size as 100 for all BiLSTM layers**. We **apply dropout to every layers** in Figure 1, and **set the dropout ratio as 0.1**. To train the model, we minimize the cross entropy of the training set, and **use the ADAM optimizer** [Kingma and Ba, 2014] to update parameters. We set **the learning rate as 0.001**. During training, we **do not update the pre-trained word embeddings**. For all the experiments, we pick the model which works the best on the dev set, and then evaluate it on the test set.

### 4.2 Model Properties

To demonstrate the properties of our model, we choose the paraphrase identification task, and experiment on the ‘‘Quora Question Pairs’’ dataset <sup>1</sup>. This dataset consists of over

<sup>1</sup><https://data.quora.com/First-Quora-Dataset-Release-Question-Pairs>

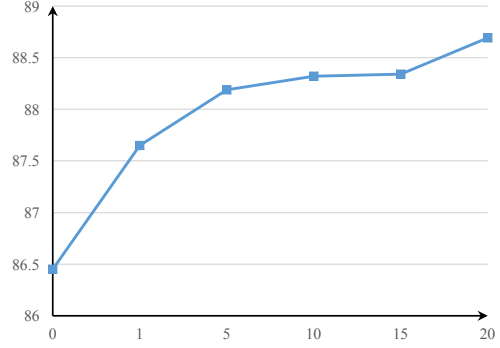


Figure 3: Influence of the multi-perspective cosine matching function in Eq.(3) .

400,000 question pairs, and each question pair is annotated with a binary value indicating whether the two questions are paraphrase of each other. We **randomly select 5,000 paraphrases and 5,000 non-paraphrases as the dev set**, and sample another 5,000 paraphrases and 5,000 non-paraphrases as the test set. We keep the remaining instances as the training set <sup>2</sup>.

First, we study the influence of our multi-perspective cosine matching function in Eq.(3). We **vary the number of perspectives  $l$  among  $\{1, 5, 10, 15, 20\}$** <sup>3</sup>, and keep the other options unchanged. We also build a baseline model by replacing Eq.(3) with the vanilla cosine similarity function. Figure 3 shows the performance curve on the dev set, where  $l = 0$  corresponds to the performance of our baseline model. We can see that, even if we only utilize one perspective ( $l = 1$ ), our model gets a significant improvement. When increasing the number of perspectives, the performance improves significantly. Therefore, our multi-perspective cosine matching function is really effective for matching vectors.

Second, to check the effectiveness of bilateral matching, we build two ablation models to matching sentences in only a single direction: 1) ‘‘Only  $P \rightarrow Q$ ’’ which only matches  $P$  against  $Q$ ; 2) ‘‘Only  $P \leftarrow Q$ ’’ which only matches  $Q$  against  $P$ . Table 1 shows the performance on the dev set. Comparing the two ablation models with the ‘‘Full Model’’, we can observe that single direction matching hurts the performance for about 1 percent. Therefore, matching sentences in two directions is really necessary for acquiring better performance.

Third, we evaluate the effectiveness of different matching strategies. To this end, we construct four ablation models (w/o Full-Matching, w/o Maxpooling-Matching, w/o Attentive-Matching, w/o Max-Attentive-Matching) by eliminating a matching strategy at each time. Table 1 shows the performance on the dev set. We can see that eliminating any of the matching strategies would hurt the performance significantly.

<sup>2</sup>We will release our source code and the dataset partition at <https://zhiguowang.github.io/>.

<sup>3</sup>Due to practical limitations, we did not experiment with more perspectives.

Models	Accuracy
Only $P \rightarrow Q$	87.74
Only $P \leftarrow Q$	87.47
w/o Full-Matching	87.86
w/o Maxpooling-Matching	87.64
w/o Attentive-Matching	87.87
w/o MaxAttentive-Matching	87.98
Full Model	<b>88.69</b>

Table 1: Ablation studies on the dev set.

Models	Accuracy
Siamese-CNN	79.60
Multi-Perspective-CNN	81.38
Siamese-LSTM	82.58
Multi-Perspective-LSTM	83.21
L.D.C.	85.55
BiMPM	<b>88.17</b>

Table 2: Performance for paraphrase identification on the Quora dataset.

### 4.3 Experiments on Paraphrase Identification

In this Sub-section, we compare our model with state-of-the-art models on the paraphrase identification task. We still experiment on the “Quora Question Pairs” dataset, and use the same dataset partition as Sub-section 4.2. This dataset is a brand-new dataset, and no previous results have been published yet. Therefore, we implemented three types of baseline models.

**First**, under the Siamese framework, we implement two baseline models: “Siamese-CNN” and “Siamese-LSTM”. Both of the two models encode two input sentences into sentence vectors with a neural network encoder, and make a decision based on the cosine similarity between the two sentence vectors. But they implement the sentence encoder with a CNN and a LSTM respectively. We design the CNN and the LSTM model according to the architectures in [Wang *et al.*, 2016c].

**Second**, based on the two baseline models, we implement two more baseline models “Multi-Perspective-CNN” and “Multi-Perspective-LSTM”. In these two models, we change the cosine similarity calculation layer with our multi-perspective cosine matching function in Eq.(3), and apply a fully-connected layer (with *sigmoid* function on the top) to make the prediction.

**Third**, we re-implement the “L.D.C.” model proposed by [Wang *et al.*, 2016d], which is a model under the “matching-aggregation” framework and acquires the state-of-the-art performance on several tasks.

Table 2 shows the performances of all baseline models and our “BiMPM” model. We can see that “Multi-Perspective-CNN” (or “Multi-Perspective-LSTM”) works much better than “Siamese-CNN” (or “Siamese-LSTM”), which further indicates that our multi-perspective cosine matching func-

Models	Accuracy
[Bowman <i>et al.</i> , 2015]	77.6
[Vendrov <i>et al.</i> , 2015]	81.4
[Mou <i>et al.</i> , 2015]	82.1
[Rocktäschel <i>et al.</i> , 2015]	83.5
[Liu <i>et al.</i> , 2016b]	85.0
[Liu <i>et al.</i> , 2016a]	85.1
[Wang and Jiang, 2015]	86.1
[Cheng <i>et al.</i> , 2016]	86.3
[Parikh <i>et al.</i> , 2016]	86.8
[Munkhdalai and Yu, 2016]	87.3
[Sha <i>et al.</i> , 2016]	87.5
[Chen <i>et al.</i> , 2016] (Single)	87.7
[Chen <i>et al.</i> , 2016] (Ensemble)	88.3
Only $P \rightarrow Q$	85.6
Only $P \leftarrow Q$	86.3
BiMPM	86.9
BiMPM (Ensemble)	<b>88.8</b>

Table 3: Performance for natural language inference on the SNLI dataset.

tion (Eq.(3)) is very effective for matching vectors. Our “BiMPM” model outperforms the “L.D.C.” model by more than two percent. Therefore, our model is very effective for the paraphrase identification task.

### 4.4 Experiments on Natural Language Inference

In this Sub-section, we evaluate our model on the natural language inference task over the SNLI dataset [Bowman *et al.*, 2015]. We test four variations of our model on this dataset, where “Only  $P \rightarrow Q$ ” and “Only  $P \leftarrow Q$ ” are the single direction matching models described in Sub-section 4.2, “BiMPM” is our full model, and “BiMPM (Ensemble)” is an ensemble version of our “BiMPM” model. We design the ensemble model by simply averaging the probability distributions [Peng *et al.*, 2015b; Peng *et al.*, 2017] of four “BiMPM” models, and each of the “BiMPM” model has the same architecture, but is initialized with a different seed.

Table 3 shows the performances of the state-of-the-art models and our models. First, we can see that “Only  $P \leftarrow Q$ ” works significantly better than “Only  $P \rightarrow Q$ ”, which tells us that, for natural language inference, matching the hypothesis against the premise is more effective than the other way around. Second, our “BiMPM” model works much better than “Only  $P \leftarrow Q$ ”, which reveals that matching premise against the hypothesis can also bring some benefits. Finally, comparing our models with all the state-of-the-art models, we can observe that our single model “BiMPM” is on par with the state-of-the-art single models, and our “BiMPM (Ensemble)” works much better than “[Chen *et al.*, 2016] (Ensemble)”. Therefore, our models achieve the state-of-the-art performance in both single and ensemble scenarios for the natural language inference task.



Models	TREC-QA		WikiQA	
	MAP	MRR	MAP	MRR
[Yang <i>et al.</i> , 2015]	0.695	0.763	0.652	0.665
[Tan <i>et al.</i> , 2015]	0.728	0.832	–	–
Wang and Itty. [2015]	0.746	0.820	–	–
[Santos <i>et al.</i> , 2016]	0.753	0.851	0.689	0.696
[Yin <i>et al.</i> , 2015]	–	–	0.692	0.711
[Miao <i>et al.</i> , 2016]	–	–	0.689	0.707
[Wang <i>et al.</i> , 2016d]	0.771	0.845	0.706	0.723
[He and Lin, 2016]	0.777	0.836	0.709	0.723
[Rao <i>et al.</i> , 2016]	0.801	<b>0.877</b>	0.701	0.718
[Wang <i>et al.</i> , 2016a]	–	–	0.734	0.742
[Wang and Jiang, 2016]	–	–	<b>0.743</b>	<b>0.755</b>
BiMPM	<b>0.802</b>	0.875	0.718	0.731

Table 4: Performance for answer sentence selection on TREC-QA and WikiQA datasets.

## 4.5 Experiments on Answer Sentence Selection

In this Sub-section, we study the effectiveness of our model for answer sentence selection tasks. The answer sentence selection task is to rank a list of candidate answer sentences based on their similarities to the question, and the performance is measured by the mean average precision (MAP) and mean reciprocal rank (MRR). We experiment on two datasets: TREC-QA [Wang *et al.*, 2007] and WikiQA [Yang *et al.*, 2015]. Experimental results of the state-of-the-art models <sup>4</sup> and our “BiMPM” model are listed in Table 4, where the performances are evaluated with the standard `trec_eval-8.0` script <sup>5</sup>. We can see that the performance from our model is on par with the state-of-the-art models. Therefore, our model is also effective for answer sentence selection tasks.

## 5 Related Work

Natural language sentence matching (NLSM) has been studied for many years. Early approaches focused on designing hand-craft features to capture n-gram overlapping, word re-ordering and syntactic alignments phenomena [Heilman and Smith, 2010; Wang and Ittycheriah, 2015]. This kind of method can work well on a specific task or dataset, but it’s hard to generalize well to other tasks.

With the availability of large-scale annotated datasets [Bowman *et al.*, 2015], many deep learning models were proposed for NLSM. The first kind of framework is based the Siamese architecture [Bromley *et al.*, 1993], where sentences are encoded into sentence vectors based on some neural network encoders, and then the relationship between two sentences was decided solely based on the two sentence vectors [Bowman *et al.*, 2015; Yang *et al.*, 2015; Tan *et al.*, 2015]. However, this kind of framework ignores the fact that the lower level interactive features between two

<sup>4</sup>[Rao *et al.*, 2016] pointed out that there are two versions of TREC-QA dataset: raw-version and clean-version. In this work, we utilized the clean-version. Therefore, we only compare with approaches reporting performance on this dataset.

<sup>5</sup>[http://trec.nist.gov/trec\\_eval/](http://trec.nist.gov/trec_eval/)

sentences are indispensable. Therefore, many neural network models were proposed to match sentences from multiple level of granularity [Yin *et al.*, 2015; Wang and Jiang, 2016; Wang *et al.*, 2016d]. Experimental results on many tasks have proved that the new framework works significantly better than the previous methods. Our model also belongs to this framework, and we have shown its effectiveness in Section 4.

## 6 Conclusion

In this work, we propose a bilateral multi-perspective matching (BiMPM) model under the “matching-aggregation” framework. Different from the previous “matching-aggregation” approaches, our model matches sentences  $P$  and  $Q$  in two directions ( $P \rightarrow Q$  and  $P \leftarrow Q$ ). And, in each individual direction, our model matches the two sentences from multiple perspectives. We evaluated our model on three tasks: paraphrase identification, natural language inference and answer sentence selection. Experimental results on standard benchmark datasets show that our model achieves the state-of-the-art performance on all tasks.

## References

- [Bowman *et al.*, 2015] Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. A large annotated corpus for learning natural language inference. *arXiv preprint arXiv:1508.05326*, 2015.
- [Bromley *et al.*, 1993] Jane Bromley, James W. Bentz, Léon Bottou, Isabelle Guyon, Yann LeCun, Cliff Moore, Eduard Säckinger, and Roopak Shah. Signature verification using a “siamese” time delay neural network. *IJPRAI*, 7(4):669–688, 1993.
- [Chen *et al.*, 2016] Qian Chen, Xiaodan Zhu, Zhenhua Ling, Si Wei, and Hui Jiang. Enhancing and combining sequential and tree lstm for natural language inference. *arXiv preprint arXiv:1609.06038*, 2016.
- [Cheng *et al.*, 2016] Jianpeng Cheng, Li Dong, and Mirella Lapata. Long short-term memory-networks for machine reading. *arXiv preprint arXiv:1601.06733*, 2016.
- [He and Lin, 2016] Hua He and Jimmy Lin. Pairwise word interaction modeling with deep neural networks for semantic similarity measurement. In *NAACL*, 2016.
- [Heilman and Smith, 2010] Michael Heilman and Noah A Smith. Tree edit models for recognizing textual entailments, paraphrases, and answers to questions. In *NAACL*, 2010.
- [Hochreiter and Schmidhuber, 1997] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [Kingma and Ba, 2014] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [LeCun *et al.*, 2015] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.

- [Liu *et al.*, 2016a] Pengfei Liu, Xipeng Qiu, and Xuanjing Huang. Modelling interaction of sentence pair with coupled-lstms. *arXiv preprint arXiv:1605.05573*, 2016.
- [Liu *et al.*, 2016b] Yang Liu, Chengjie Sun, Lei Lin, and Xiaolong Wang. Learning natural language inference using bidirectional lstm model and inner-attention. *arXiv preprint arXiv:1605.09090*, 2016.
- [Miao *et al.*, 2016] Yishu Miao, Lei Yu, and Phil Blunsom. Neural variational inference for text processing. In *ICML*, 2016.
- [Mikolov *et al.*, 2013] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositional-ity. In *Advances in neural information processing systems*, pages 3111–3119, 2013.
- [Mou *et al.*, 2015] Lili Mou, Rui Men, Ge Li, Yan Xu, Lu Zhang, Rui Yan, and Zhi Jin. Natural language inference by tree-based convolution and heuristic matching. *arXiv preprint arXiv:1512.08422*, 2015.
- [Munkhdalai and Yu, 2016] Tsendsuren Munkhdalai and Hong Yu. Neural tree indexers for text understanding. *arXiv preprint arXiv:1607.04492*, 2016.
- [Parikh *et al.*, 2016] Ankur P Parikh, Oscar Täckström, Dipanjan Das, and Jakob Uszkoreit. A decomposable attention model for natural language inference. *arXiv preprint arXiv:1606.01933*, 2016.
- [Peng *et al.*, 2015a] Xi Peng, Junzhou Huang, Qiong Hu, Shaoting Zhang, Ahmed Elgammal, and Dimitris Metaxas. From circle to 3-sphere: Head pose estimation by instance parameterization. *Computer Vision and Image Understanding*, 136:92–102, 2015.
- [Peng *et al.*, 2015b] Xi Peng, Shaoting Zhang, Yu Yang, and Dimitris N Metaxas. Piefa: Personalized incremental and ensemble face alignment. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3880–3888, 2015.
- [Peng *et al.*, 2016] Xi Peng, Rogerio S Feris, Xiaoyu Wang, and Dimitris N Metaxas. A recurrent encoder-decoder network for sequential face alignment. In *European Conference on Computer Vision*, pages 38–56. Springer International Publishing, 2016.
- [Peng *et al.*, 2017] Xi Peng, Shaoting Zhang, Yang Yu, and Dimitris N Metaxas. Toward personalized modeling: Incremental and ensemble alignment for sequential faces in the wild. *International Journal of Computer Vision*, pages 1–14, 2017.
- [Pennington *et al.*, 2014] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *EMNLP*, 2014.
- [Rao *et al.*, 2016] Jinfeng Rao, Hua He, and Jimmy Lin. Noise-contrastive estimation for answer selection with deep neural networks. In *CIKM*, 2016.
- [Rocktäschel *et al.*, 2015] Tim Rocktäschel, Edward Grefenstette, Karl Moritz Hermann, Tomáš Kočiský, and Phil Blunsom. Reasoning about entailment with neural attention. *arXiv preprint arXiv:1509.06664*, 2015.
- [Santos *et al.*, 2016] Cicero dos Santos, Ming Tan, Bing Xiang, and Bowen Zhou. Attentive pooling networks. *arXiv preprint arXiv:1602.03609*, 2016.
- [Sha *et al.*, 2016] Lei Sha, Baobao Chang, Zhifang Sui, and Sujian Li. Reading and thinking: Re-read lstm unit for textual entailment recognition. In *COLING*, 2016.
- [Tan *et al.*, 2015] Ming Tan, Cicero dos Santos, Bing Xiang, and Bowen Zhou. Lstm-based deep learning models for non-factoid answer selection. *arXiv preprint arXiv:1511.04108*, 2015.
- [Vendrov *et al.*, 2015] Ivan Vendrov, Ryan Kiros, Sanja Fidler, and Raquel Urtasun. Order-embeddings of images and language. *arXiv preprint arXiv:1511.06361*, 2015.
- [Wang and Ittycheriah, 2015] Zhiguo Wang and Abraham Ittycheriah. Faq-based question answering via word alignment. *arXiv preprint arXiv:1507.02628*, 2015.
- [Wang and Jiang, 2015] Shuohang Wang and Jing Jiang. Learning natural language inference with lstm. *arXiv preprint arXiv:1512.08849*, 2015.
- [Wang and Jiang, 2016] Shuohang Wang and Jing Jiang. A compare-aggregate model for matching text sequences. *arXiv preprint arXiv:1611.01747*, 2016.
- [Wang *et al.*, 2007] Mengqiu Wang, Noah A Smith, and Teruko Mitamura. What is the jeopardy model? a quasi-synchronous grammar for qa. In *EMNLP*, 2007.
- [Wang *et al.*, 2016a] Bingning Wang, Kang Liu, and Jun Zhao. Inner attention based recurrent neural networks for answer selection. In *ACL*, 2016.
- [Wang *et al.*, 2016b] Zhiguo Wang, Haitao Mi, Wael Hamza, and Radu Florian. Multi-perspective context matching for machine comprehension. *arXiv preprint arXiv:1612.04211*, 2016.
- [Wang *et al.*, 2016c] Zhiguo Wang, Haitao Mi, and Abraham Ittycheriah. Semi-supervised clustering for short text via deep representation learning. In *CoNLL*, 2016.
- [Wang *et al.*, 2016d] Zhiguo Wang, Haitao Mi, and Abraham Ittycheriah. Sentence similarity learning by lexical decomposition and composition. In *COLING*, 2016.
- [Yang *et al.*, 2015] Yi Yang, Wen-tau Yih, and Christopher Meek. Wikiqa: A challenge dataset for open-domain question answering. In *EMNLP*, 2015.
- [Yin *et al.*, 2015] Wenpeng Yin, Hinrich Schütze, Bing Xiang, and Bowen Zhou. Abcnn: Attention-based convolutional neural network for modeling sentence pairs. *arXiv preprint arXiv:1512.05193*, 2015.