

# Sentence Modeling via Multiple Word Embeddings and Multi-level Comparison for Semantic Textual Similarity

Nguyen Huy Tien<sup>1</sup>, Nguyen Minh Le<sup>1</sup>, Yamasaki Tomohiro<sup>2</sup>, Izuha Tatsuya<sup>2</sup>

<sup>1</sup>Japan Advanced Institute of Science and Technology (JAIST)

<sup>2</sup>Toshiba Research & Development Center, Japan

ntienhuy@jaist.ac.jp, nguyenml@jaist.ac.jp,

tomohiro2.yamasaki@toshiba.co.jp, tatsuya.izuha@toshiba.co.jp

## Abstract

Different word embedding models capture different aspects of linguistic properties.

This inspired us to propose a model (M-MaxLSTM-CNN) for employing multiple sets of word embeddings for evaluating sentence similarity/relation.

Representing each word by multiple word embeddings, the MaxLSTM-CNN encoder generates a novel sentence embedding. We then learn the similarity/relation between our sentence embeddings via Multi-level comparison. Our method M-MaxLSTM-CNN consistently shows strong performances in several tasks (i.e., measure textual similarity, identify paraphrase, recognize textual entailment). According to the experimental results on STS Benchmark dataset and SICK dataset from SemEval, M-MaxLSTM-CNN outperforms the state-of-the-art methods for textual similarity tasks. Our model does not use hand-crafted features (e.g., alignment features, Ngram overlaps, dependency features) as well as does not require pre-trained word embeddings to have the same dimension.

## 1 Introduction

Measuring the semantic similarity/relation of two pieces of short text plays a fundamental role in a variety of language processing tasks (i.e., plagiarism detection, question answering, and machine translation). Semantic textual similarity (STS) task is challenging because of the diversity of linguistic expression. For example, two sentences with different lexicons could have a similar meaning. Moreover, the task requires to measure similarity at several levels (e.g., word level, phrase

level, sentence level). These challenges give difficulties to conventional approaches using hand-crafted features.

Recently, the emergence of word embedding techniques, which encode the semantic properties of a word into a low dimension vector, leads to the successes of many learning models in natural language processing (NLP). For example, [Kalchbrenner et al. \(2014\)](#) randomly initialize word vectors, then tunes them during the training phase of a sentence classification task. By contrast, [Huy Tien and Minh Le \(2017\)](#) initialize word vectors via the pre-train word2vec model trained on Google News ([Mikolov et al., 2013b](#)). [Wieting et al. \(2015\)](#) train a word embedding model on the paraphrase dataset PPDB, then apply the word representation for word and bi-gram similarity tasks.

Several pre-trained word embeddings are available, which are trained on various corpora under different models. [Levy and Goldberg \(2014\)](#) observed that different word embedding models capture different aspects of linguistic properties: a Bag-of-Words contexts based model tends to reflect the domain aspect (e.g., scientist and research) while a paraphrase-relationship based model captures semantic similarities of words (e.g., boy and kid). From experiments, we also observed that the performance of a word embedding model is usually inconsistent over different datasets. This inspired us to develop a model taking advantages of various pre-trained word embeddings for measuring textual similarity/relation.

In this paper, we propose a convolutional neural network (CNN) to learn a multi-aspect word embedding from various pre-trained word embeddings. We then apply the max-pooling scheme and Long Short Term Memory (LSTM) on this embedding to form a sentence representation. In STS tasks, [Shao \(2017\)](#) shows the efficiency of the max-pooling scheme in modeling sentences

from word embedding representations refined via CNN. However, the max-pooling scheme lacks the property of word order (e.g., *sentence*("Bob likes Mary") = *sentence*("Mary likes Bob")). To address this weakness, we use LSTM as an additional scheme for modeling sentences with word order characteristics. For measuring the similarity/relation between two sentence representations, we propose Multi-level comparison which consists of word-word level, sentence-sentence level, and word-sentence level. Through these levels, our model comprehensively evaluates the similarity/relation between two sentences.

We evaluate our M-MaxLSTM-CNN model on three tasks: STS, textual entailment recognition, paraphrase identification. The advantages of M-MaxLSTM-CNN are: i) simple but efficient for combining various pre-trained word embeddings with different dimensions; ii) using Multi-level comparison shows better performances compared to using only sentence-sentence comparison; iii) does not require hand-crafted features (e.g., alignment features, Ngram overlaps, syntactic features, dependency features) compared to the state-of-the-art ECNU (Tian et al., 2017) on STS Benchmark dataset.

Our main contributions are as follows:

- Propose MaxLSTM-CNN encoder for efficiently encoding sentence embeddings from multiple word embeddings.
- Propose Multi-level comparison (M-MaxLSTM-CNN) to learn the similarity/relation between two sentences. The model achieves strong performances over various tasks. Especially in STS tasks, the method obtains the state-of-the-art results.

The remainder of this paper is organized as follows: Section 2 reviews the previous research, Section 3 introduces the architecture of our model, Section 4 describes the three tasks and datasets, Section 5 describes the experiment setting, Section 6 reports and discusses the results of the experiments, and Section 7 concludes our work.

## 2 Related work

Most prior research on modeling textual similarity relied on feature engineering. Wan et al. (2006) extract  $n$ -gram overlap features and dependency-based features, while Madnani et al. (2012) em-

ploy features based on machine translation metrics. Mihalcea et al. (2006) propose a method using corpus-based and knowledge-based measures of similarity. Das and Smith (2009) design a model which incorporates both syntax and lexical semantics using dependency grammars. Ji and Eisenstein (2013) combine the fine-grained  $n$ -gram overlap features with the latent representation from matrix factorization. Xu et al. (2014) develop a latent variable model which jointly learns paraphrase relations between word and sentence pairs. Using Dependency trees, Sultan et al. (2014) propose a robust monolingual aligner and successfully applied it for STS tasks.

The recent emergence of deep learning models has provided an efficient way to learn continuous vectors representing words/sentences. By using a neural network in the context of a word prediction task, Bengio et al. (2003) and (Mikolov et al., 2013a) generate word embedding vectors carrying semantic meanings. The embedding vectors of words which share similar meanings are close to each other. To capture the morphology of words, Bojanowski et al. (2017) enrich the word embedding with character  $n$ -grams information. Closest to this approach, Wieting et al. (2016a) also propose to represent a word or sentence using a character  $n$ -gram count vector. However, the objective function for learning these embeddings is based on paraphrase pairs.

For modeling sentences, composition approach attracted many studies. Yessenalina and Cardie (2011) model each word as a matrix and used iterated matrix multiplication to present a phrase. Tai et al. (2015) design a Dependency Tree-Structured LSTM for modeling sentences. This model outperforms the linear chain LSTM in STS tasks. Convolutional neural network (CNN) has recently been applied efficiently for semantic composition (Kalchbrenner et al., 2014; Kim, 2014; Shao, 2017). This technique uses convolutional filters to capture local dependencies in term of context windows and applies a pooling layer to extract global features. He et al. (2015) use CNN to extract features at multiple level of granularity. The authors then compare their sentence representations via multiple similarity metrics at several granularities. Gan et al. (2017) propose a hierarchical CNN-LSTM architecture for modeling sentences. In this approach, CNN is used as an encoder to encode an sentence into a continuous representa-

tion, and LSTM is used as a decoder. [Conneau et al. \(2017\)](#) train a sentence encoder on a textual entailment recognition database using a BiLSTM-Maxpooling network. This encoder achieves competitive results on a wide range of transfer tasks.

At SemEval-2017 STS task, hybrid approaches obtain strong performances. [Wu et al. \(2017\)](#) train a linear regression model with WordNet, alignment features and the word embedding *word2vec*<sup>1</sup>. [Tian et al. \(2017\)](#) develop an ensemble model with multiple boosting techniques (i.e., Random Forest, Gradient Boosting, and XGBoost). This model incorporates traditional features (i.e., n-gram overlaps, syntactic features, alignment features, bag-of-words) and sentence modeling methods (i.e., Averaging Word Vectors, Projecting Averaging Word Vectors, LSTM).

MVCNN model ([Yin and Schütze, 2015](#)) and MGNC-CNN model ([Zhang et al., 2016](#)) are close to our approach. In MVCNN, the authors use variable-size convolution filters on various pre-trained word embeddings for extracting features. However, MVCNN requires word embeddings to have the same size. In MGNC-CNN, the authors apply independently CNN on each pre-trained word embedding for extracting features and then concatenate these features for sentence classification. By contrast, our M-MaxLSTM-CNN model jointly applies CNN on all pre-trained word embeddings to learn a multi-aspect word embedding. From this word representation, we encode sentences via the max-pooling and LSTM. To learn the similarity/relation between two sentences, we employ Multi-level comparison.

### 3 Model description

Our model (shown in Figure 1) consists of three main components: i) learning a multi-aspect word embedding (Section 3.1); ii) modeling sentences from this embedding (Section 3.2); iii) measuring the similarity/relation between two sentences via Multi-level comparison (section 3.3).

#### 3.1 Multi-aspect word embedding

Given a word  $w$ , we transfer it into a word vector  $e_w^{concat}$  via  $K$  pre-trained word embeddings as follows:

$$e_w^{concat} = e_w^1 \oplus e_w^2 \oplus \dots \oplus e_w^K \quad (1)$$

<sup>1</sup><https://code.google.com/p/word2vec/>

where  $\oplus$  is concatenation operator,  $e_w^i$  is the word embedding vector of  $w$  in the  $i$ th pre-trained embedding.

To learn a multi-aspect word embedding  $e_w^{multi}$  from the representation  $e_w^{concat}$ , we design  $H$  convolutional filters. Each filter  $r_i$  is denoted as a weight vector with the same dimension as  $e_w^{concat}$  and a bias value  $b_{r_i}$ . The  $e_w^{multi}$  is obtained by applying these filters on the  $e_w^{concat}$  as follows:

$$e_w^{r_i} = \sigma(e_w^{concat} r_i^T + b_{r_i}) \quad (2)$$

$$e_w^{multi} = [e_w^{r_1}, e_w^{r_2}, \dots, e_w^{r_H}] \quad (3)$$

where  $\sigma$  denotes a logistic sigmoid function.

The next section explains how to model a sentence from its multiple-aspect word embeddings.

#### 3.2 Sentence modeling

Given an input sentence  $s = [w_1, w_2, \dots, w_n]$ , we obtain a sequence of multiple-aspect word embeddings  $s^{multi} = [e_{w_1}^{multi}, e_{w_2}^{multi}, \dots, e_{w_n}^{multi}]$  using Eq. (1-3). For modeling the sentence from the representation  $s^{multi}$ , we use two schemes: max-pooling and LSTM.

**Max-pooling scheme:** To construct a max-pooling sentence embedding  $e_s^{max}$ , the most potential features are extracted from the representation  $s^{multi}$  as follows:

$$e_s^{max}[i] = \max(e_{w_1}^{multi}[i], e_{w_2}^{multi}[i], \dots, e_{w_n}^{multi}[i]) \quad (4)$$

where  $e_{w_k}^{multi}[i]$  is the  $i$ th element of  $e_{w_k}^{multi}$ .

**LSTM scheme:** From Eq. (4), we find that the max-pooling scheme ignores the property of word order. Therefore, we construct a LSTM sentence embedding  $e_s^{lstm}$  to support the sentence embedding  $e_s^{max}$ . The representation  $s^{multi}$  is transformed to a fix-length vector by recursively applying a LSTM unit to each input  $e_{w_t}^{multi}$  and the previous step  $h_{t-1}$ . At each time step  $t$ , the LSTM unit with  $l$ -memory dimension defines six vectors in  $\mathbb{R}^l$ : input gate  $i_t$ , forget gate  $f_t$ , output gate  $o_t$ , tanh layer  $u_t$ , memory cell  $c_t$  and hidden state  $h_t$  as follows (from [Tai et al. \(2015\)](#)):

$$i_t = \sigma(W_i e_{w_t}^{multi} + U_i h_{t-1} + b_i) \quad (5)$$

$$f_t = \sigma(W_f e_{w_t}^{multi} + U_f h_{t-1} + b_f) \quad (6)$$

$$o_t = \sigma(W_o e_{w_t}^{multi} + U_o h_{t-1} + b_o) \quad (7)$$

$$u_t = \tanh(W_u e_{w_t}^{multi} + U_u h_{t-1} + b_u) \quad (8)$$

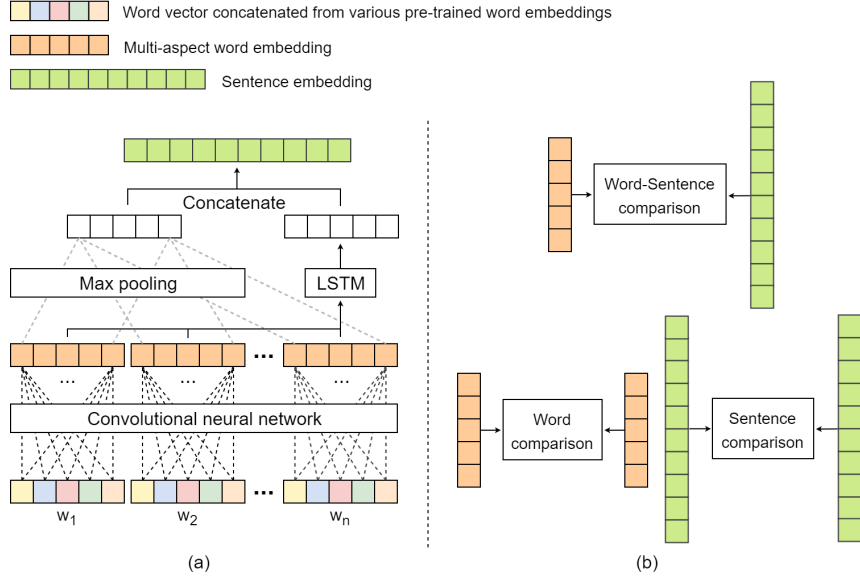


Figure 1: The proposed M-MaxLSTM-CNN model: (a) MaxLSTM-CNN encoder; (b) Multi-level comparison.

$$c_t = f_t \odot c_{t-1} + i_t \odot u_t \quad (9)$$

$$h_t = o_t \odot \tanh(c_t) \quad (10)$$

$$e_s^{lstm} = h_n \quad (11)$$

where  $\sigma, \odot$  respectively denote a logistic sigmoid function and element-wise multiplication;  $W_i, U_i, b_i$  are respectively two weights matrices and a bias vector for input gate  $i$ . The denotation is similar to forget gate  $f$ , output gate  $o$ , tanh layer  $u$ , memory cell  $c$  and hidden state  $h$ .

Finally, the sentence embedding  $e_s$  is obtained by concatenating the two sentence embeddings  $e_s^{max}$  and  $e_s^{lstm}$ :

$$e_s = e_s^{max} \oplus e_s^{lstm} \quad (12)$$

### 3.3 Multi-level comparison

In this section, we describe the process for evaluating the similarity/relation between two sentences. We compare two sentences via three levels: word-word, sentence-sentence and word-sentence.

#### 3.3.1 Word-word comparison

Given two input sentences  $s_1$  and  $s_2$ , we encode them into two sequences of multi-aspect word embeddings  $s_1^{multi}$  and  $s_2^{multi}$  (Section 3.2). We then compute a word-word similarity vector  $sim^{word}$  as follows:

$$A_{ij} = \frac{s_1^{multi}[i] \cdot s_2^{multi}[j]}{\|s_1^{multi}[i]\| \|s_2^{multi}[j]\|} \quad (13)$$

$$sim^{word} = \sigma(W^{word}g(A) + b^{word}) \quad (14)$$

where  $s_t^{multi}[i]$  is the  $i$ th multi-aspect word embedding of sentence  $s_t$ ;  $g()$  is a function to flatten a matrix into a vector; and  $W^{word}$  and  $b^{word}$  are respectively a weight matrix and a bias parameter.

#### 3.3.2 Sentence-sentence comparison

Given two input sentences  $s_1$  and  $s_2$ , we encode them into two sentence embeddings  $e_{s_1}$  and  $e_{s_2}$  (Section 3.1 and 3.2). To compute the similarity/relation between the two embeddings, we introduce four comparison metrics:

**Cosine similarity:**

$$d_{cosine} = \frac{e_{s_1} \cdot e_{s_2}}{\|e_{s_1}\| \|e_{s_2}\|} \quad (15)$$

**Multiplication vector & Absolute difference:**

$$d_{mul} = e_{s_1} \odot e_{s_2} \quad (16)$$

$$d_{abs} = |e_{s_1} - e_{s_2}| \quad (17)$$

where  $\odot$  is element-wise multiplication.

**Neural difference:**

$$x = e_{s_1} \oplus e_{s_2} \quad (18)$$

$$d_{neu} = W^{neu}x + b^{neu} \quad (19)$$

where  $W^{neu}$  and  $b^{neu}$  are respectively a weight matrix and a bias parameter.

As a result, we have a sentence-sentence similarity vector  $sim^{sent}$  as follows:

$$d^{sent} = d_{cosine} \oplus d_{mul} \oplus d_{abs} \oplus d_{neu} \quad (20)$$

$$sim^{sent} = \sigma(W^{sent}d^{sent} + b^{sent}) \quad (21)$$

where  $W^{sent}$  and  $b^{sent}$  are respectively a weight matrix and a bias parameter.

### 3.3.3 Word-sentence comparison

Given a sentence embedding  $e_{s_1}$  and a sequence of multi-aspect word embeddings  $s_2^{multi}$ , we compute a word-sentence similarity matrix  $sim_{s_1}^{ws}$  as follows:

$$e_{s_1}^{ws}[i] = e_{s_1} \oplus s_2^{multi}[i] \quad (22)$$

$$sim_{s_1}^{ws}[i] = \sigma(W^{ws}e_{s_1}^{ws}[i] + b^{ws}) \quad (23)$$

where  $s_2^{multi}[i]$  is the multi-aspect word embedding of the  $i$ th word in sentence  $s_2$ ;  $W^{ws}$  and  $b^{ws}$  are respectively a weight matrix and a bias parameter.

As a result, we have a word-sentence similarity vector  $sim^{ws}$  for the two sentences as follows:

$$sim^{ws} = \sigma(W^{ws'}[g(sim_{s_1}^{ws}) \oplus g(sim_{s_2}^{ws})] + b^{ws'}) \quad (24)$$

where  $g()$  is a function to flatten a matrix into a vector;  $W^{ws'}$  and  $b^{ws'}$  are respectively a weight matrix and a bias parameter.

Finally, we compute a target score/label of a sentence pair as follows:

$$sim = sim^{word} \oplus sim^{sent} \oplus sim^{ws} \quad (25)$$

$$h_s = \sigma(W^{l1}sim + b^{l1}) \quad (26)$$

$$\hat{y} = softmax(W^{l2}h_s + b^{l2}) \quad (27)$$

where  $W^{l1}$ ,  $W^{l2}$ ,  $b^{l1}$  and  $b^{l2}$  are model parameters;  $\hat{y}$  is a predicted target score/label.

## 4 Tasks & Datasets

We evaluate our model on three tasks:

- **Textual entailment recognition:** given a pair of sentences, we predict a directional relation between the sentences (entailment/contradiction/neutral). We evaluate this task on **SICK** dataset. It was collected for the 2014 SemEval competition and includes examples of the lexical, syntactic and semantic phenomena and ignores other aspects of existing sentential datasets (i.e., idiomatic multiword expressions, named entities, telegraphic language).
- **Semantic textual similarity:** given a pair of sentences, we measure a semantic similarity score of this pair. We use two datasets for this task:

- **STSB:** comprises a careful selection of the English data sets used in SemEval and \*SEM STS shared tasks from 2012 to 2017. This dataset cover three genres: image captions, news headlines and user forums. Each sentence pair is annotated with a relatedness score  $\in [0, 5]$ .
- **SICK:** Each sentence pair is annotated with a relatedness score  $\in [1, 5]$ .

- **Paraphrase identification:** given a pair of sentences, we predict a binary label indicating whether the two sentences are paraphrases. Microsoft Research Paraphrase Corpus (MRPC) is used for this task. It includes pairs of sentences extracted from news source on the web.

Table 1 shows the statistic of the three datasets. Because of not dealing with name entities and multi-word idioms, the vocabulary size of SICK is quite small compared to the others.

Dataset	Train	Validation	Test	$l$	$ V $
STSB	5,749	1,500	1,379	11	15,997
SICK	4,500	500	4,927	9	2,312
MRPC	3,576	500	1,725	21	18,003

Table 1: Statistic of datasets.  $|V|$ ,  $l$  denote the vocabulary size, and the average length of sentences respectively.

## 5 Experiment setting

### 5.1 Pre-trained word embeddings

We study five pre-trained word embeddings<sup>2</sup> for our model:

- **word2vec** is trained on Google News dataset (100 billion tokens). The model contains 300-dimensional vectors for 3 million words and phrases.
- **fastText** is learned via skip-gram with sub-word information on Wikipedia text. The embedding representations in fastText are 300-dimensional vectors.
- **GloVe** is a 300-dimensional word embedding model learned on aggregated global word-word co-occurrence statistics from Common Crawl (840 billion tokens).

<sup>2</sup>These embeddings are available at *anonymous*



- **Baroni** uses a context-predict approach to learn a 400-dimensional semantic embedding model. It is trained on 2.8 billion tokens constructed from ukWaC, the English Wikipedia and the British National Corpus.
- **SL999** is trained under the skip-gram objective with negative sampling on word pairs from the paraphrase database **PPDB**. This 300-dimensional embedding model is tuned on SimLex-999 dataset (Hill et al., 2016).

## 5.2 Model configuration

In all of the tasks, we used the same model configuration as follows:

- Convolutional filters: we used 1600 filters. It is also the dimension of the word embedding concatenated from the five pre-trained word embeddings.
- LSTM dimension: we also selected 1600 for LSTM dimension.
- Neural similarity layers: the dimension of  $b^{word}$ ,  $b^{sent}$ ,  $b^{ws}$  and  $b^{ws'}$  are respectively 50, 5, 5 and 100.
- Penultimate fully-connected layer: has the dimension of 250 and is followed by a drop-out layer ( $p = 0.5$ ).

We conducted a grid search on 30% of STSB dataset to select these optimal hyper-parameters.

## 5.3 Training Setting

### 5.3.1 Textual entailment recognition & Paraphrase identification

In these tasks, we use the cross-entropy objective function and employ AdaDelta as the stochastic gradient descent (SGD) update rule with mini-batch size as 30. Details of Adadelta method can be found in Zeiler (2012). During the training phase, the pre-trained word embeddings are fixed.

### 5.3.2 Semantic Textual Similarity

To compute a similarity score of a sentence pair in the range  $[1, K]$ , where  $K$  is an integer, we replace Eq. (27) with the equations in Tai et al. (2015) as follows:

$$\hat{p}_\theta = \text{softmax}(W^{l2}h_s + b^{l2}) \quad (28)$$

$$\hat{y} = r^T \hat{p}_\theta \quad (29)$$

where  $W^{l1}$ ,  $W^{l2}$ ,  $b^{l1}$  and  $b^{l2}$  are parameters;  $r^T = [1, 2, \dots, K]$ ;  $\hat{y}$  is a predicted similarity score.

A sparse target distribution  $p$  which satisfies  $y = r^T p$  is computed as:

$$p_i = \begin{cases} y - \lfloor y \rfloor, & i = \lfloor y \rfloor + 1 \\ \lfloor y \rfloor - y + 1, & i = \lfloor y \rfloor \\ 0 & \text{otherwise} \end{cases} \quad (30)$$

for  $i \in [1, K]$ , and  $y$  is the similarity score.

To train the model, we minimize the regularized KL-divergence between  $p$  and  $\hat{p}_\theta$ :

$$J(\theta) = \frac{1}{m} \sum_{k=1}^m KL(p^{(k)} || \hat{p}_\theta^{(k)}) \quad (31)$$

where  $m$  is the number of training pairs and  $\theta$  denotes the model parameters. The gradient descent optimization Adadelta is used to learn the model parameters. We also use mini-batch size as 30 and keep the pre-trained word embeddings fixed during the training phase. We evaluate our models through Pearson correlation  $r$ .

## 6 Experiments and Discussion

This section describes two experiments: i) compare our model against recent systems; ii) evaluate the efficiency of using multiple pre-trained word embeddings.

### 6.1 Overall evaluation

Besides existing methods, we also compare our model with several sentence modeling approaches using multiple pre-trained word embeddings:

- **Word Average:**

$$e_s = \frac{1}{n} \sum_{i=1}^n e_{w_i}^{concat} \quad (32)$$

where  $e_s$  is the sentence embedding of a  $n$ -words sentence, and  $e_{w_i}^{concat}$  is from Eq. (1)

- **Project Average:**

$$e_s = \sigma(W(\frac{1}{n} \sum_{i=1}^n e_{w_i}^{concat}) + b) \quad (33)$$

where  $W$  is a  $1600 \times 1600$  weight matrix, and  $b$  is a 1600 bias vector.

- **LSTM:** apply Eq. (5-11) on  $e_{w_i}^{concat}$  to construct the 1600-dimension  $e_s$  sentence embedding.

Method	STSB	SICK-R	SICK-E	MRPC
<i>Ensemble models/Feature engineering</i>				
DT_TEAM (Maharjan et al., 2017)	79.2	-	-	-
ECNU (Tian et al., 2017)	81	-	-	-
BIT (Wu et al., 2017)	80.9	-	-	-
TF-KLD (Ji and Eisenstein, 2013)	-	-	-	<b>80.41/85.96</b>
<i>Neural representation models with one embedding</i>				
Multi-Perspective CNN (He et al., 2015)	-	86.86	-	78.6/84.73
InferSent (Conneau et al., 2017)	75.8	88.4	<b>86.1</b>	76.2/83.1
GRAN (Wieting and Gimpel, 2017)	76.4	86	-	-
Paragram-Phrase (Wieting et al., 2016b)	73.2	86.84	85.3	-
HCTI (Shao, 2017)	78.4	-	-	-
<i>Neural representation models with the five embeddings using sentence-sentence comparison (S)</i>				
S-Word Average	71.06	81.18	80.88	71.48/81.1
S-Project Average	75.12	86.53	85.12	75.48/82.47
S-LSTM	77.14	85.15	85.6	70.43/79.71
S-Max-CNN	81.87	88.3	84.33	76.35/83.75
S-MaxLSTM-CNN	82.2	88.47	84.9	77.91/84.31
<i>Neural representation models with the five embeddings using Multi-level comparison (M)</i>				
M-Max-CNN	82.11	88.45	84.7	76.75/83.64
M-MaxLSTM-CNN	<b>82.45</b>	<b>88.76</b>	84.95	78.1/84.5

Table 2: Test set results with Pearson’s  $r$  score  $\times 100$  for STS tasks, and accuracy for other tasks. Bold-face values show the highest scores in each dataset. SICK-R and SICK-E denote the STS task and the entailment task in SICK dataset respectively.

- **Max-CNN**: apply Eq. (2-4) on  $e_{w_i}^{concat}$  to construct the 1600-dimension  $e_s$  sentence embedding.

We report the results of these methods in Table 2. Overall, our M-MaxLSTM-CNN shows competitive performances in these tasks. Especially in the STS task, M-MaxLSTM-CNN outperforms the state-of-the-art methods on the two datasets. Because STSB includes complicated samples compared to SICK, the performances of methods on STSB are quite lower. In STSB, the prior top performance methods use ensemble approaches mixing hand-crafted features (word alignment, syntactic features, N-gram overlaps) and neural sentence representations, while our approach is only based on a neural sentence modeling architecture. In addition, we observed that InferSent shows the strong performance on SICK-R but quite low on STSB while our model consistently obtains the strong performances on both of the datasets. InferSent uses transfer knowledge on textual entailment data, consequently it obtains the strong performance on this entailment task.

According to Wieting et al. (2016b), using Word Average as the compositional architecture outperforms the other architectures (e.g., Project Average, LSTM) for STS tasks. In a multiple word embeddings setting, however, Word Average does not show its efficiency. Each word embedding

model has its own architecture as well as objective function. These factors makes the vector spaces of word embeddings are different. Therefore, we intuitively need a step to learn or refine a representation from a set of pre-trained word embeddings rather than only averaging them. Because Project Average model, LSTM model and Max-CNN model have their parameters for learning sentence embeddings, they significantly outperform Word Average model.

We observed that MaxLSTM-CNN outperforms Max-CNN in both of the settings (i.e., sentence-sentence comparison, Multi-level comparison). As mentioned in Section 1, Max-CNN ignores the property of word order. Therefore, our model achieves improvement compared to Max-CNN by additionally employing LSTM for capturing this property.

We only applied Multi-level comparison on Max-CNN and MaxLSTM-CNN because these encoders generate multi-aspect word embeddings. The experimental results prove the efficiency of using Multi-level comparison. In the textual entailment dataset **SICK-E**, the task mainly focuses on interpreting the meaning of a whole sentence pair rather than comparing word by word. Therefore, the performance of Multi-level comparison is quite similar to sentence-sentence comparison in the SICK-E task. This is also the reason why

Word embedding	STSB		SICK-R & SICK-E			MRPC	
	Pearson	$ V _{avai}(\%)$	Pearson	Acc	$ V _{avai}(\%)$	Acc/F1	$ V _{avai}(\%)$
word2Vec	78.9	75.64	87.27	84.09	98.53	75.42/82.13	67.81
fastText	79.95	84.27	87.59	83.42	99.18	74.31/81.75	79.04
Glove	80.1	91.71	88.21	84.71	99.78	74.9/82.782	89.85
SL999	80.31	94.76	87.26	84.55	99.83	76.46/83.13	94.19
Baroni	79.81	90.54	86.9	83.99	98.83	74.84/82.4	87.92
All	<b>82.45</b>	95.65	<b>88.76</b>	<b>84.95</b>	99.83	<b>78.1/84.5</b>	95.97

Table 3: Evaluation of exploiting multiple pre-trained word embeddings.  $|V|_{avai}$  is the proportion of vocabulary available in a word embedding. In case of using all word embeddings,  $|V|_{avai}$  denotes the proportion of vocabulary where each word is available in at least one word embedding.

LSTM, which captures global relationships in sentences, has the strong performance in this task.

## 6.2 Evaluation of exploiting multiple pre-trained word embeddings

In this section, we evaluate the efficiency of using multiple pre-trained word embeddings. We compare our multiple pre-trained word embeddings model against models using only one pre-trained word embedding. The same objective function and Multi-level comparison are applied for these models. In case of using one pre-trained word embedding, the dimension of LSTM and the number of convolutional filters are set to the length of the corresponding word embedding. Table 3 shows the experimental results of this comparison. Because the approach using five word embeddings outperforms the approaches using two, three, or four word embeddings, we only report the performance of using five word embeddings. We also report  $|V|_{avai}$  which is the proportion of vocabulary available in a pre-trained word embedding. SICK dataset ignores idiomatic multi-word expressions, and named entities, consequently the  $|V|_{avai}$  of SICK is quite high.

We observed that no word embedding has strong results on all the tasks. Although trained on the paraphrase database and having the highest  $|V|_{avai}$ , the SL999 embedding could not outperform the Glove embedding in SICK-R. HCTI (Shao, 2017), which is the current state-of-the-art in the group of neural representation models on STSB, also used the Glove embedding. However, the performance of HCTI in STSB (78.4) is lower than our model using the Glove embedding. In SICK-R, InferSent (Conneau et al., 2017) achieves a strong performance (88.4) using the Glove embedding with transfer knowledge, while our model with only the Glove embedding achieves a performance close to the performance of InferSent.

These results confirm the efficiency of Multi-level comparison.

In STSB and MRPC, as employing the five pre-trained embeddings, the  $|V|_{avai}$  is increased. This factor limits the number of random values when initializing word embedding representations because a word out of a pre-trained word embedding is assigned a random word embedding representation. In other words, a word out of a pre-trained word embedding is assigned a random semantic meaning. Therefore, the increase of the  $|V|_{avai}$  improves the performance of measuring textual similarity. In STSB and MRPC, our multiple pre-trained word embedding achieves a significant improvement in performance compared against using one word embedding. In SICK-R and SICK-E, although the  $|V|_{avai}$  is not increased when employing five pre-trained embeddings, the performance of our model is improved. This fact shows that our model learned an efficient word embedding via these pre-trained word embeddings.

## 7 Conclusion

In this work, we study an approach employing multiple pre-trained word embeddings and Multi-level comparison for measuring semantic textual relation. The proposed M-MaxLSTM-CNN architecture consistently obtains strong performances on several tasks. Compared to the state-of-the-art methods in STS tasks, our model does not require handcrafted features (e.g., word alignment, syntactic features) as well as transfer learning knowledge. In addition, it allows using several pre-trained word embeddings with different dimensions.

Future work could apply our multiple word embeddings approach for transfer learning tasks. This strategy allows making use of pre-trained word embeddings as well as available resources.



## Acknowledgments

This work was done while Nguyen Tien Huy was an intern at Toshiba Research Center.

## References

- Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. 2003. A neural probabilistic language model. *Journal of machine learning research* 3(Feb):1137–1155.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *Transactions of the Association of Computational Linguistics* 5:135–146.
- Alexis Conneau, Douwe Kiela, Holger Schwenk, Loic Barrault, and Antoine Bordes. 2017. Supervised learning of universal sentence representations from natural language inference data. *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*.
- Dipanjan Das and Noah A Smith. 2009. Paraphrase identification as probabilistic quasi-synchronous recognition. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 1-Volume 1*. Association for Computational Linguistics, pages 468–476.
- Zhe Gan, Yunchen Pu, Ricardo Henao, Chunyuan Li, Xiaodong He, and Lawrence Carin. 2017. Learning generic sentence representations using convolutional neural networks. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. pages 2380–2390.
- Hua He, Kevin Gimpel, and Jimmy J Lin. 2015. Multi-perspective sentence similarity modeling with convolutional neural networks. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. pages 1576–1586.
- Felix Hill, Roi Reichart, and Anna Korhonen. 2016. Simlex-999: Evaluating semantic models with (genuine) similarity estimation. *Computational Linguistics*.
- Nguyen Huy Tien and Nguyen Minh Le. 2017. An ensemble method with sentiment features and clustering support. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Asian Federation of Natural Language Processing, Taipei, Taiwan, pages 644–653.
- Yangfeng Ji and Jacob Eisenstein. 2013. Discriminative improvements to distributional sentence similarity. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, pages 891–896.
- Nal Kalchbrenner, Edward Grefenstette, and Phil Blunsom. 2014. A convolutional neural network for modeling sentences. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, pages 655–665. <https://doi.org/10.3115/v1/P14-1062>.
- Yoon Kim. 2014. Convolutional neural networks for sentence classification. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, pages 1746–1751. <https://doi.org/10.3115/v1/D14-1181>.
- Omer Levy and Yoav Goldberg. 2014. Dependency-based word embeddings. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*. pages 302–308.
- Nitin Madnani, Joel Tetreault, and Martin Chodorow. 2012. Re-examining machine translation metrics for paraphrase identification. In *Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, pages 182–190.
- Nabin Maharjan, Rajendra Banjade, Dipesh Gautam, Lasang J Tamang, and Vasile Rus. 2017. Dt.team at semeval-2017 task 1: Semantic similarity using alignments, sentence-level embeddings and gaussian mixture model output. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*. pages 120–124.
- Rada Mihalcea, Courtney Corley, Carlo Strapparava, et al. 2006. Corpus-based and knowledge-based measures of text semantic similarity. In *AAAI*. volume 6, pages 775–780.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013a. Efficient estimation of word representations in vector space. In *Proceedings of Workshop at the International Conference on Learning Representations*.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013b. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*. pages 3111–3119.
- Yang Shao. 2017. Hcti at semeval-2017 task 1: Use convolutional neural network to evaluate semantic textual similarity. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*. pages 130–133.
- Md Arafat Sultan, Steven Bethard, and Tamara Sumner. 2014. Back to basics for monolingual alignment: Exploiting word similarity and contextual evidence. *Transactions of the Association of Computational Linguistics* 2:219–230.

- Kai Sheng Tai, Richard Socher, and Christopher D. Manning. 2015. Improved semantic representations from tree-structured long short-term memory networks. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Association for Computational Linguistics, pages 1556–1566.
- Junfeng Tian, Zhiheng Zhou, Man Lan, and Yuanbin Wu. 2017. Ecnv at semeval-2017 task 1: Leverage kernel-based traditional nlp features and neural networks to build a universal model for multilingual and cross-lingual semantic textual similarity. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*. pages 191–197.
- Stephen Wan, Mark Dras, Robert Dale, and Cécile Paris. 2006. Using dependency-based features to take the para-farce out of paraphrase.
- John Wieting, Mohit Bansal, Kevin Gimpel, and Karen Livescu. 2016a. Charagram: Embedding words and sentences via character n-grams. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. pages 1504–1515.
- John Wieting, Mohit Bansal, Kevin Gimpel, and Karen Livescu. 2016b. Towards universal paraphrastic sentence embeddings. *The International Conference on Learning Representations*.
- John Wieting, Mohit Bansal, Kevin Gimpel, Karen Livescu, and Dan Roth. 2015. From paraphrase database to compositional paraphrase model and back. *Transactions of the ACL (TACL)*.
- John Wieting and Kevin Gimpel. 2017. Revisiting recurrent networks for paraphrastic sentence embeddings. In *ACL (1)*. Association for Computational Linguistics, pages 2078–2088.
- Hao Wu, Heyan Huang, Ping Jian, Yuhang Guo, and Chao Su. 2017. Bit at semeval-2017 task 1: Using semantic information space to evaluate semantic textual similarity. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*. pages 77–84.
- Wei Xu, Alan Ritter, Chris Callison-Burch, William B. Dolan, and Yangfeng Ji. 2014. Extracting lexically divergent paraphrases from twitter. *Transactions of the Association for Computational Linguistics* 2:435–448.
- Ainur Yessenalina and Claire Cardie. 2011. Compositional matrix-space models for sentiment analysis. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, pages 172–182.
- Wenpeng Yin and Hinrich Schütze. 2015. Multichannel variable-size convolution for sentence classification. In *Proceedings of the Nineteenth Conference on Computational Natural Language Learning*. Association for Computational Linguistics, pages 204–214.
- Matthew D. Zeiler. 2012. ADADELTA: an adaptive learning rate method. *CoRR* abs/1212.5701.
- Ye Zhang, Stephen Roller, and Byron C. Wallace. 2016. Mgnc-cnn: A simple approach to exploiting multiple word embeddings for sentence classification. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, pages 1522–1527. <https://doi.org/10.18653/v1/N16-1178>.