

Combining Knowledge And Corpus Based Techniques with Convolutional Neural Nets for Semantic Textual Similarity

Shivam Soni^{a,c}, Kavi Narayana Murthy^{b,c}

^a*sonishivam131@gmail.com*

^b*knmuh@yahoo.com*

^c*University of Hyderabad*

Abstract

Semantic Textual Similarity(STS) is the task of measuring the semantic similarity between two snippets of text. For instance, given two sentences, the idea is to compute upto what degree the two sentences are similar in meaning. We have employed a hybrid architecture which combines various *Knowledge and Corpus Based techniques* [43] with a Convolutional Neural Network followed by a Fully-Connected-Layer inspired from [44]. Moreover, we have enhanced the *Knowledge and Corpus Based techniques* called *WordNet-Augmented Word Overlap* and *Vector Space Sentence* presented in [43]. In *WordNet-Augmented Word Overlap*, we have replaced *pathlen similarity*(used to calculate the similarity between the senses in the WordNet) with: *Jiang- Conrath similarity*(JCS)[23] and *Lin Similarity*(LS)[29] and in *Vector Space Sentence Similarity*, we have replaced LSA word vectors with vectors given by Google’s pre-trained Word2Vec model [35]. We have also tried various modifications in CNN architecture like changing the pre-trained vectors, loss functions and number of hidden layers and hidden unit. After extracting the features for the input pair of sentences using our combined architecture, a regressor model was trained that outputs the final similarity score. With our final hybrid architecture we are able achieve better results than some of the benchmark results as per SemEval 2017 [13].

Keywords: Knowledge and Corpus Based techniques, WordNet-Augmented Word Overlap, Jiang- Conrath similarity, Lin Similarity, Vector Space Sentence Similarity, CNN.

1. Introduction

Semantic Textual Similarity (STS) task follows the idea that, given two sentences, compute upto what degree the two sentences are similar. For instance, given two sentences, the idea is to compute upto what degree the two sentences are similar in meaning. The degree of similarity is measured on a scale of 0 to 5 (real values), according to the shared tasks competition held by SemEval (International Workshop on Semantic Evaluation) [4], where a score of '0' represents complete unrelatedness and '5' represents semantic equivalence. For example :

- A score of '5' can be given if two sentences means exactly the same thing, as shown below:

Bob is diving into the pool.

Bob is jumping into the swimming pool.

- A score of '3' can be given if two sentences roughly means the same thing but some information is missing in one or the other, as shown below:

Akshay said chapter 5 is very important than chapter 3 from the book.

Chapter 3 is not so important. Akshay said.

- A score of '0' can be given if two sentences are completely independent from each other, as shown below:

John is going for a walk with his friends.

There is bird flying in the sky.

The dataset released by SemEval for STS task from the year 2012-2017 contains gold standard labels for every pair of sentences which is annotated by various experienced annotators ([4] [3] [2] [1] [10] [13]). To measure the performance of any model, i.e, how well a particular model has captured the semantic similarity between two sentences, **Pearson Correlation Coefficient** is used as standard across all datasets. SemEval had also released various datasets for measuring similarity between cross-lingual sentences [2] [1] [10] [13]. In this work we have worked on monolingual task of measuring similarity between two English sentences.

Recent researches have shown a need of measuring semantic similarity between the sentences. There are various applications where STS task can be applied such as for Text Summarization, Machine Translation, Question Answering, Dialogue and Conversational Systems and Text Generation. In text mining, it can be used to locate unseen text from the data [6]. STS is related to various other tasks such as Paraphrase Detection [52] [18], Textual Entailment [9] and Semantic relatedness [31]. A lot of work using Knowledge Based and Corpus technique(KACB) has been done in the past for finding the similarity between two sentences. The Knowledge Based techniques commonly use the lexical database WordNet[37] as the source and includes various semantic similarity measures such as *Pairwise Word Similarity* ([7], [33], [43]). [46] made use of word alignment technique which used Paraphrase Database (PPDB), which is a large lexical and phrasal database. The Corpus Based techniques involves measures such as Pointwise Mutual Information (PMI)[33], Vector Space Similarity [43], LSA models [26] etc. Lately, various deep learning techniques like [20], [22], [21], [8], [38],[44] etc. were proposed for the STS task. These Neural based techniques became very competitive with traditional knowledge based and corpus based techniques.

[7] in their observed that the earlier methods like [17] [16] [28] [26] [33] had a limitation due their assumption, that a single measure in itself is enough to capture the semantic relations between sentences and therefore they combined various techniques in their proposed work. In our work, we have followed a similar notion of combining various techniques for the task of measuring similarity between sentences. We have proposed a hybrid architecture which combines various *Knowledge and Corpus Based techniques* (KACB) [43] with

Convolutional Neural Network followed by a Fully-Connected-Layer inspired from [44]. The combined output was then finally passed into a regressor model to produce the final output. Our overall view of the system can be understood from the figure 2. Moreover, we have modified two techniques presented in [43] which are, *WordNet-Augmented Word Overlap* and *Vector Space Sentence Similarity*. In *WordNet-Augmented Word Overlap*, we have replaced *pathlen similarity* (used to calculate the similarity between the senses in the WordNet) with: *Jiang- Conrath similarity* (JCS) [23] and *Lin Similarity* (LS) [29]. The drawback of *pathlen similarity* is that it assigns same similarity score to different pairs of concepts without considering the generality or specificity of the pairs in WordNet heirarchy. Whereas, JCS and LS similarities rely on WordNet structure and adds probabilistic information derived from a corpus. Also, we have made various modifications in the CNN architecture by changing the loss functions and pre-trained word embeddings.

In the next section 2 we will discuss about the background of STS task which is followed by some discussion about various different methods in the section 3. In the section 4 we will discuss about the proposed work and finally we will conclude by discussing about various experiments and about future work in section 5 .

2. Literature Review

A rise in the research on measuring similarity between two sentences mainly started when SemEval (International Workshop on Semantic Evaluation) announced Semantic Textual Similarity shared task, which successfully ran from the year 2012 to 2017 [4] [3] [2] [1] [10] [13]. Before the SemEval STS task was announced, a substantial amount of work had been done for measuring the similarity between the pair of texts but it was mainly focused on comparison between two long texts or documents [5] [19] [25] [32] and very few publications related to measuring similarity between short texts were reported [16] [28] [26] [33]. And the methods used for comparing long texts are generally inefficient when applied on short text because two similar long texts are usually dependent on degree of co-occurring words. Whereas, in short texts co-occurring words maybe less or null sometimes. Therefore, in this field, some of the earliest works can be considered to be given by [16] [28] [26] [33]. Also, [28] [26] were among the first to provide dataset for this task. Where [28] contributed 65 English sentence pairs and [26] contributed 50 document pairs suitable for the STS task. There are various different kinds of measures that has been used till now to compute the similarity between a pair of sentences, where, the most commonly used measures for a long period of time were based on knowledge based and corpus based techniques (Section 2.1). We have also discussed about word alignment technique in the same section. Recently, the impact of Deep Neural Nets based techniques proved quite effective, which is discussed in Section 2.2.

2.1. Knowledge Based, Corpus Based and Word Alignment Techniques

[33] in their work used *Corpus-based Measures* (which included Pointwise Mutual Information (PMI) [49] and LSA [24]) and *Knowledge-based Measures* which makes use word similarity metrics: Leacock and Chodorow [15], Lesk [27], Wu and Palmer [51], Resnik [40],

Lin [29], Jiang and Conrath [23]. These word similarity metrics make use of lexical database (Wordnet database in thier work). The PMI technique and word similarity metrics were used to find similarity between each word from the first text segment and each word from the second text segment which had the same part-of-speech tag. Then using a similarity metric, which they proposed in their work [33] , they determined the similarity score between the two texts. [28] followed a step by step process, where they first created an *order vector* and a *raw semantic vector* by using a lexical database(they used Wordnet) and a joint word set formed by the distinct words from both the sentences. Then the *raw semantic vector* was converted into a *semantic vector* by using Brown corpus. Finally, the sentence similarity was computed by combining semantic similarity(calculated using two semantic vectors) and word order similarity (calculated using two word order). [7] combined various *content based*, *style based* and *structure based* similarity measures which resulted in more than 300 score vectors (out of which 20 features were finally used), which when combined together, served as a feature set for a simple log-linear classifier . Their results gave one of the best performances in the STS task-2012. Similar approach of combining various measures was followed by [43] , and by the combination of various knowledge based, corpus based and syntactical techniques which they proposed, a feature set of scores was constructed. These extracted feature set was used to train a SVR model, which finally predicted sentence similarity score. Some examples of their knowledge based techniques are: *WordNet-Augmented Word Overlap*, *Weighted Word Overlap etc.*, all them used *WordNet* as their lexical database. The corpus based techniques like : *Vector Space Sentence Similarity* and *Weighted Vector Space Similarity* used LSA vectors.

In the year 2014, [46] gave one the best and successfull measure for the STS task using word alignment technique presented by them in [45]. Their work follows an alignment pipeline, where, in the first step, all the Identical Word Sequences are aligned and in the second step, all the Named Entities are aligned. In the third and the fourth step content words are aligned basically using two properties of the words: 1. if they are semantically similar, 2. if they occur in similar contexts in respective sentences. The context of the word pairs from the two respective sentences is defined by using syntactic dependencies in the third step and by taking a surrounding window of size 3 in the step four. The final score is then computed using harmonic mean of $prop_{Al}^{(1)}$ and $prop_{Al}^{(2)}$, where $prop_{Al}^{(1)}$ is the proportion of content words in $S(1)$ that are aligned and vice-versa for $prop_{Al}^{(2)}$.

2.2. Deep Neural Network Based Work

Lately, deep learning related work started giving competition to various syntactical and lexical based techniques. One of the top performances were given by [20]. They developed a *sentence model* using a Convolutional Neural Network(CNN) architecture in which they used two different types of filters. The first one was a holistic filter, which by the help of a particular window size matches the complete word vectors and the second one was a per-dimensional filter which independently match each dimension of word embeddings. The output of the sentence model after passing it through various different types of pooling layer was given into a *similarity measurement layer*, which compared different local region of the sentence representation, and finally a fully-connected layer was employed at the end to get

the final result. In 2016, they released another STS measurement techniques [21], which gave even better results. Initially, a Bidirectional Long Short-Term Memory Network (Bi-LSTMs) was used to model the context of the input sentences and then using a *Pairwise Word Interaction Model* each word from first sentence was compared with each word of the second sentence and 13 different compared values were computed, which resulted in a 3D matrix of size $13 \times |sent1| \times |sent2|$. This was finally passed through a *focus layer* and then into 19 layer Deep ConvNet to produce the final output. It was the first paper ever, which combined the use of both Bi-LSTMs and Deep CNNs for any task. The multi-perspective CNN model (MPCNN) [20] discussed above, was modified in [22] by adding an interaction layer before passing the input sentences in MPCNN model. This modification was mainly done to make the input sentences more inter-related, because earlier, the model lacked contextual interaction information between the sentences. [8] in 2017 presented an end-to-end LSTM system which was a version of Siamese LSTM [38]. In Siamese LSTM [38], an LSTM model is employed for each input sentence, weights are shared between layers and updated parallelly during learning phase. Whereas, [8] allowed two sentences to be given as an input into the same LSTM model and shared the weights. This work by [8] also gave quality results. Another paper in 2017 [44], gave a simple CNN model without adding any dropout or batchnorm layer and passed the output of the CNN model into a Fully-Connected Neural Network (FCNN) to get the final scores.

3. Methods

3.1. Knowledge Based and Corpus Based Techniques

The various knowledge and corpus based techniques which we have used in our work are inspired from [43], which we discuss as follows:

3.1.1. *N - Gram Overlap*

The ngram overlap between two sentences (S_1 and S_2) is expressed as follows:

$$ngo(S_1, S_2) = 2 \cdot \left(\frac{|S_1|}{|S_1 \cap S_2|} + \frac{|S_2|}{|S_1 \cap S_2|} \right)^{-1} \quad (1)$$

where $|S_1|$ and $|S_2|$ is the length of unigrams, bigrams or trigrams in S_1 and S_2 respectively. $|S_1 \cap S_2|$ is the length of common unigrams, bigrams or trigrams between S_1 and S_2 .

The idea behind ngram-overlap is, it computes that upto what extent sentences S_1 and S_2 covers each other and then finally computing a harmonic mean between the two quantities.

3.1.2. *WordNet - Augmented Word Overlap(WAWO)*

Wordnet augmented coverage between sentences (S_1 and S_2) is expressed as follows:

$$wn(S_1, S_2) = \frac{1}{|S_2|} \sum_{w_1 \in S_1} score(w_1, S_2) \quad (2)$$

$$score(w, S) = \begin{cases} 1 & \text{if } w \in S \\ \max_{\bar{w} \in S} sim_{pl}(w, \bar{w}) & \text{otherwise} \end{cases}$$

$wn(S_1, S_2)$ in 2 finds the wordnet coverage between S_1 and S_2 , similarly $wn(S_2, S_1)$ is also calculated between S_2 and S_1 , where, finally a harmonic mean is calculated between $wn(S_1, S_2)$ and $wn(S_2, S_1)$ to compute the final score for this method.

$score(w_1, S_2)$ helps to find the score between each word from S_1 with each word from S_2 . Where $sim_{pl}(w, \bar{w})$ computes the pathlen similarity between w and \bar{w} . And pathlen similarity between two concatenated c_1 and c_2 is calculated as:

$$sim_{pl}(c_1, c_2) = \frac{1}{Pathlen(c_1, c_2)}$$

$Pathlen(c_1, c_2)$ = No. of edges in the shortest path (in hypernym graph) between 2 concepts.

3.1.3. Weighted Word Overlap

The following equation computes weighted word coverage of the first sentence S_1 by the second sentence S_2 :

$$wvc(S_2, S_1) = \frac{\sum_{w \in S_1 \cap S_2} IC(w)}{\sum_{\bar{w} \in S_2} IC(\bar{w})} \quad (3)$$

Similarly $wvc(S_1, S_2)$ is also calculated and a harmonic mean between $wvc(S_2, S_1)$ and $wvc(S_1, S_2)$ is computed to get the final score for *weighted Word Overlap* measure.

$IC(w)$ in the above equation 3 is the information content of a word w from a sentence, which is expressed as follows:

$$IC(w) = \ln \frac{\sum_{\bar{w} \in C} freq(\bar{w})}{freq(w)}$$

where $freq(w)$ refers to the frequency of a word in a corpus and C is the a of words from the corpus.

3.1.4. Vector Space Sentence Similarity

In this method LSA vectors are used to represent a word in a sentence and the whole sentence is represented by taking the summation of all the word vectors present in the sentence, expressed as following:

$$\vec{v}(S) = \sum_{w \in S} u_w$$

where $\vec{v}(S)$ is the vector representation of the sentence and u_w is the LSA vector representation of a word w . An additional representation v_W is also used, where each word vector

is assigned some weight, which is basically the information content of the word w and as expressed as following :

$$v_W(S) = \sum_{w \in S} IC(w)u_w$$

Finally, vector space sentence similarity features between the two sentences is computed as $|\cos(v(S_1), v(S_2))|$ and $|\cos(v_W(S_1), v_W(S_2))|$.

3.1.5. Normalized Difference and Number Overlap

The normalized differences between the two sentences is calculated using three different ways: 1) using the chunks of noun, verb, predicate counts from sentences, 2) using lengths of the sentences and 3) by aggregating the information content of the words in a sentence.

Numbers Overlap includes three different features given as follows:

$$(A \subseteq B \vee B \subseteq A), (\log(1 + |A| + |B|)) \text{ and } (|A \cap B|)/(|A| + |B|)$$

where A and B are the sets of numbers in S1 and S2.

3.1.6. Named Entity Feature

In this feature, they considered capitalized words as named entities and found the overlap between them as done equation 1. They also found the overlap between the stock index symbols as well. Stock index symbols are the tokens which are in all caps and begin with a period.

3.2. Convolutional Neural Network

Convolution in most simple terms is a mathematical operation on two functions which computes a third a function, expressing how a particular function can bring changes to the other function. It has been used in various fields like, *probability, statistics, signal processing, computer vision, natural language processing and differential equations*. In our work we have used CNNs for NLP task and we will be discussing about it in section 3.3.

CNNs are a variety of different convolution layers with *non-linear activations* like ReLU, tanh or sigmoid computed on to the output of each of them. Earlier, when CNNs were not so popular, a lot of research in *computer vision* was focused on what kind of different filters can be used to extract features out of an image or any type of input given. But with CNNs, by the help of various loss functions, it learned the weights for the filters by its own and was able to perform much better.

CNNs are able to deal with various shortcomings of Deep Neural Nets (DNNs) with all fully-connected layers. In DNNs we have a multiple fully-connected layers which results in too many parameters, because of which they are often prone to overfitting. So, the most important questions that arose was, can we have DNNs which can have many non-linearity, but have fewer parameters and hence less prone to overfitting? The answer is yes, and all these problems are resolved by CNNs. Formally, two main aspects of CNNs are **Sparse Connectivity** and **Weight Sharing**, which makes them perform better than DNNs. We discuss the aspects of CNNs below:

Sparse Connectivity: In traditional feedforward neural network, each neuron of the output layer is connected to each neuron of the previous layer. Which finally results in more number

of parameters and hence overfitting. Whereas, in CNNs one particular neuron in the output is a result of convolution of filters with a specific region of the input. Therefore, it is connected only to some neighbouring neurons of the input, which provides sparse connectivity and results in *fewer number of parameters*.

Weight Sharing: In CNNs, **weights are nothing but the kernel**, and the **kernel remains constant as we pass over the entire input matrix**, creating different output values based on the neighborhood of inputs. One complete pass of the kernel over the input matrix constitutes one Convolutional output. As we know, the kernel is constant, thereby **the weights get shared for all the neurons in that output area**. Only the inputs vary.

Thereby, we are effectively reducing the number of parameters yet still retaining model complexity (many non-linearity) and overcoming one of the shortcomings of DNNs.

3.3. Why and How are CNNs related to NLP tasks

CNNs are usually used for images as input. But when used for texts, the document or sentence is represented in a matrix form. As an example, for a sentence of length 15 words and each word represented as a vector of dimension 300, the input would be a 15×300 matrix.

In NLP, generally, filters slide over full dimensions of words, i.e., they move only in 1-dimension, taking two words or three words at a time depending on the size of window. A full picture of how it works has been shown in figure 1, which has been taken from [53]. We can see that, 3 different choices of window sizes for filter has been used, 1st with a window size of two, 2nd with a window size of three and 4th with a window size of four, with 2 filters for each. Therefore we get 6 outputs, 1 from each of the filters, which is then max-pooled and concatenated to give final classification among two categories. For intuition it can be looked upon as extracting n-grams features out of a sentence.

Yes, we can say that, there are not very clear intuitions of how CNNs provide benefits for local invariance and what exactly higher-level representation mean in case of NLP as in computer vision. So, it may seem that Recurrent Neural Networks (RNNs) works better as it is more intuitive in case of texts. But this does not mean that CNNs do not work at all. Many models might not work well but some have proved to be very useful in recent studies. A big defence for CNNs is, they are very fast and are very comfortably used at hardware level on GPUs. If compared to n-grams, CNNs are more efficient in terms of representation, because a large size of vocabulary is computationally expensive, even Google has not provided anything greater than 5-grams, but with CNNs we can even work with 6-grams or greater than that. The filters are automatically able to learn good representations, without a need of any large vocabulary. We have discussed below about different applications for CNNs on NLP tasks.

3.3.1. Loss function

The loss functions which we have used in CNNs for computing the errors between the *actual value* and the *predicted value* are given as follows:

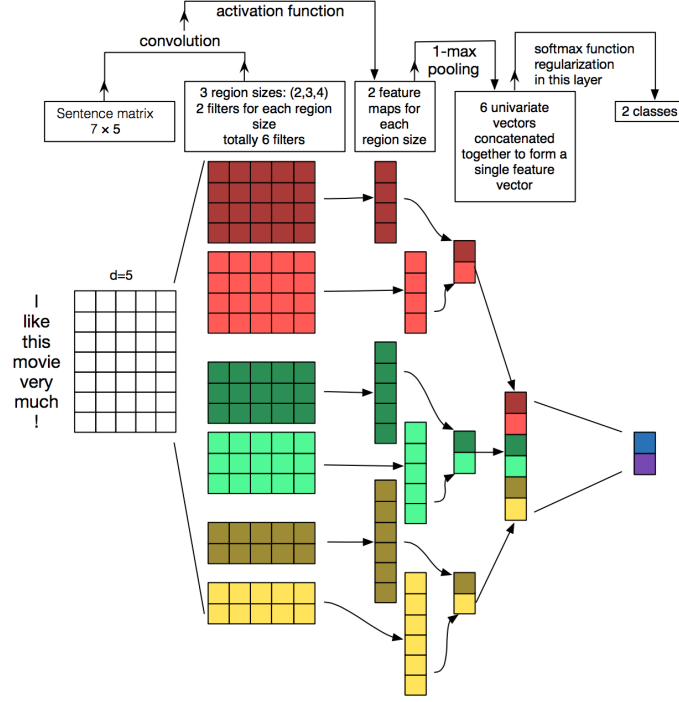


Figure 1: CNNs for Sentence classification [53]

1) Pearson Correlation Coefficient

Pearson Correlation Coefficient (ρ) is a statistical measure that captures linear correlation between two quantitative, continuous variables. Its range is between -1 to +1, where, -1 represents a negative linear correlation and +1 represents a positive linear correlation. ρ between two random variable (A and B) is given as follows:

$$\rho_{A,B} = \frac{\text{cov}(A,B)}{\sigma_A \sigma_B} \quad (4)$$

where, σ_A and σ_B represents the standard deviation of A and B respectively and cov is the *covariance* between A and B .

Pearson's correlation coefficient for sample data is represented as r_{xy} and is expressed as follows:

$$r_{ab} = \frac{\sum_{i=1}^m (a_i - \bar{a})(b_i - \bar{b})}{\sqrt{\sum_{i=1}^m (a_i - \bar{a})^2} \sqrt{\sum_{i=1}^m (b_i - \bar{b})^2}} \quad (5)$$

where, m is the total number of samples, a_i , b_i are the individual points indexed with i and \bar{a} , \bar{b} are sample mean.

2) Kullback-Leibler Divergence(KL-Divergence)

KL-Divergence capture differences between two probability distribution, i.e., it measures how on probability distribution is different from another. If there are no differences,i.e., if probability distribution perfectly matches the KL-Divergence gives a score of zero , otherwise it can take any value 0 and ∞ .

KL-Divergence between two probability distribution $A(i)$ and $B(i)$ (where $B(i)$ is the approximation and $A(i)$ is the true distribution is expressed as follows:

$$D_{KL}(A||B) = \sum_{i \in I} A(i) \log \left(\frac{A(i)}{B(i)} \right) \quad (6)$$

4. Proposed Work

4.1. Architecture

We have employed a hybrid architecture by combining various knowledge based and corpus techniques (KACB) with a CNN architecture, as shown in figure 2. Each technique present in KAB computes a similarity score between sentences in range 0 to 1 (real values) when a sentence pair is given as an input to them. Also, CNN followed by an FCNN produces a similarity score between two sentences. All the results are then combined to train a regression model(SVR or MLP regressor), which produces the final score.

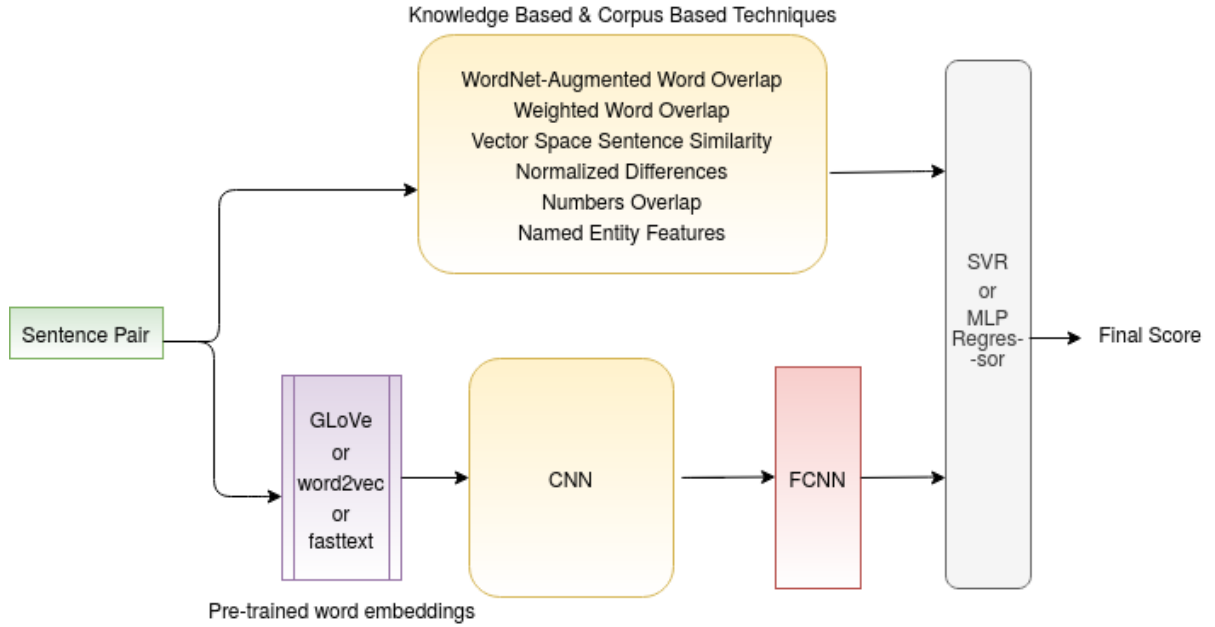


Figure 2: Overview of the sytem

4.2. Knowledge Based and Corpus Techniques

The knowledge based and corpus based techniques have been taken from [43] (as discussed in section (3.1)) with some modifications in two techniques called *WordNet-Augmented Word Overlap* which is knowledge based technique and *Vector Space Sentence Similarity* which is a corpus based technique. We will discuss the drawbacks and the modifications which have been employed to overcome the drawbacks in following sections.

4.2.1. WordNet - Augmented Word Overlap (WAWO)

The name of the technique in itself clarifies the knowledge based source used in this technique is the lexical database *WordNet* [37]. As discussed in section 3.1.2, WAWO makes use of *pathlen similarity* to measure the similarity between two concepts present in the wordnet hierarchy. But a drawback with pathlen similarity is, it assigns equal similarity scores to different pairs of concepts without considering the generality or specificity of the pairs in WordNet hierarchy. To understand it better, we can consider an example from the figure 3.

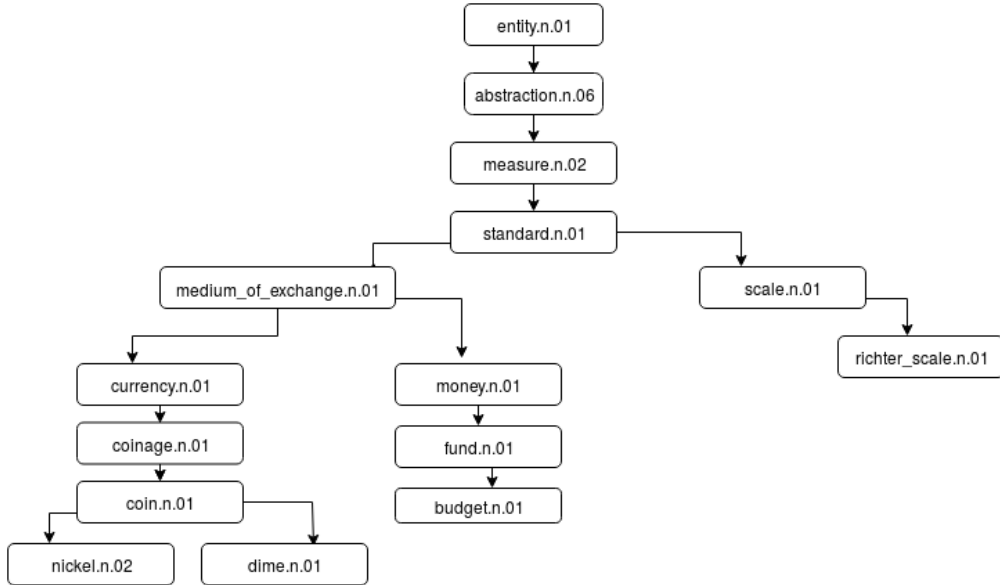


Figure 3: WordNet Hierarchy Example

Let us consider two pairs of concepts (*nickel.n.02*, *dime.n.01*) and (*medium_of_exchange.n.01*, *scale.n.01*) from the figure 3. Both the pairs have a pathlen of 2 (i.e, the count of the edges in the shortest path between them), therefore *pathlen similarity* between them is 0.5 given by the equation 7.

$$sim_{pl} = \frac{1}{pathlen(C1, C2)} \quad (7)$$

But in actual, if we look intuitively, (*nickel.n.02*, *dime.n.01*) are more similar to each other than (*medium_of_exchange.n.01*, *scale.n.01*), also (*nickel.n.02*, *dime.n.01*) are more specific in hierarchy than (*medium_of_exchange.n.01*, *scale.n.01*), therefore the similarity score for (*nickel.n.02*, *dime.n.01*) should have been greater. So, we can observe that the similarities given by *pathlen similarity* are not accurate. To overcome these drawbacks we have used *Jiang-Conrath Similarity (JC similarity)* and *Lin Similarity* proposed by [23] and [29] respectively, replacing *pathlen similarity*, which we discuss in the next sections.

1) Jiang-Conrath Similarity (JC Similarity)

. Jiang- Conrath similarity(JC Similarity) [23] makes use of lengths of graph edges to overcome the drawbacks discussed above and is expressed as:

$$sim_{JC}(c_1, c_2) = \frac{1}{IC(c_1) + IC(c_2) - 2 * IC(LCS(c_1, c_2))} \quad (8)$$

where $LCS(c_1, c_2)$ is the least common subsumer between the concepts c_1 and c_2 , $IC(c)$ is the *Information Content* of the concept c . In order to calculate $IC(c)$ we have used the same formula as given by *Resnik(1998)* in [41], which is expressed as follows:

$$IC(c) = -\log P(c) \quad (9)$$

we have used *Brown Corpus*, *BNC Corpus* and *treebank Corpus* from NLTK to calculate IC of each concept in WordNet. $P(c)$ is the probability that any random word *say* x from the corpus is a *hyponym/instance or a part* of concept c in WordNet hierarchy. This random word x ranges over whole corpus. Therefore, we can say that $P(\text{root}) = 1$ because any word x from the corpus will be an instance of root concept. This probability is given as:

$$P(c) = \frac{\sum_{w \in S_w(c)} count(w)}{T} \quad (10)$$

where S_w represents the set of words that are an instance of concept c , and T is the total number of words which are present in corpus as well as in the thesaurus.

2) Lin- Similarity

. Lin - similarity [29] is the information content common to c_1 and c_2 , normalized by their average information content. The similarity is given as follows:

$$sim_{Lin}(c_1, c_2) = \frac{2 * IC(LCS(c_1, c_2))}{IC(c_1) + IC(c_2)} \quad (11)$$

It basically tells the more information they don't share , the less similar they are.

4.2.2. Vector Space Sentence Similarity

This feature as we discussed in section 3.1.4 makes use of *LSA vectors* to represent each word in a sentence. Now, in this work, rather than using *LSA vectors*, vectors provided by Google's pre-trained *word2vec model*[36] has been used. This word2vec model makes use of *skip-gram architecture* and is trained on various news dataset, which contributed approximately *100 billion* word . Also, same modification has been done with the weighted word vectors, where word vectors are weighted by their Information Content. These modifications have been done because *word2vec algorithm* has shown to capture better similarity between the words from a given corpus ([34]). But it gives better performance when the training data is in large amount, that is the reason a pre-trained model is preferred here.

4.2.3. Experiments

We have done five different types of experiments in our work on *knowledge based and corpus techniques*, the abbreviations and interpretation of them are as follows:

jc: It contains results for the changes made in *WordNet-Augmented Word Overlap*, where *pathlen similarity* is replaced by *JC similarity* keeping rest of the techniques same.

lin: It contains results for the changes made in *WordNet-Augmented Word Overlap*, where *pathlen similarity* is replaced by *Lin similarity* keeping rest of the techniques same.

jc.lin: It contains results for the changes made in *WordNet-Augmented Word Overlap*, where *pathlen similarity* is replaced by *Lin similarity* in one feature and by *JC similarity* in another, keeping rest of features same.

ww_wwv: It contains results for the changes made in *Vector Space Sentence Similarity*, where LSA vectors are replaced by Word2Vec vectors and also an additional feature is used where LSA vectors is replaced by weighted Word2Vec vectors as discussed in Section(3.1.4), rest all the features are the same.

jc.lin_ww_wwv : In this technique all the proposed changes are merged together keeping rest of the features same.

These experiments were carried on two different sets of dataset. First, we see how the five experiments performed on the dataset provided by SemEval in the year 2012 [4] to compare our results with [43]. Secondly, we see how our experiments performed on the latest monolingual dataset on English provided by SemEval in 2017 [13], to compare it with baseline models.

1) Results on 2012 SemEval Dataset

. The dataset is taken from the SemEval STS task (2012)[4] which contains in total 5342 sentence pairs in English language and a corresponding score in range 0 to 5 (all real values). This data is collected from various different sources as shown in the 3.

The MSRvid and MSRpar dataset is collected by Microsoft Research(MSR). SMTeuroparl dataset was created from ACL Workshop based on Machine Translation in 2007 and

Dataset	Pairs	Source
MSRpar	1500	News
MSRvid	1500	Video
SMTeuroparl	750	MT eval
SMTnews	399	MT eval
OnWn	1193	Glosses

Table 1: Dataset Description

Technique	MSRvid	MSRpar	SMTeuroparl	SMTnews	OnWn
(saric-etal)[43]	0.8803	0.7343	0.4771	0.3989	0.6797
jc	0.8747	0.7381	0.5364	0.4237	0.6914
lin	0.8705	0.7371	0.5362	0.4319	0.6904
jc_lin	0.8713	0.7377	0.5354	0.4068	0.6899
wv_wwv	0.8837	0.7378	0.5376	0.5014	0.7067
jc_lin_wv_wwv	0.8791	0.7367	0.5323	0.4720	0.7027

Table 2: Pearson Correlation Values for each dataset (from 2012) corresponding to each experiments .

2008 . [11] [12]. In addition, two surprise datasets SMTnews and OnWn were also released which are used here.

The results are computed by finding the *Pearson Correlation* 3.3.1 between the scores given by the model on the test dataset and the gold-standard human annotated scores.

In table 2, we have compared our experiments with [43]. We have observed that except for *MSRpar dataset*, *wv_wwv* has performed better for all the datasets. For *MSRpar dataset* *jc* gave slightly better result.

2) Results on Monolingual SemEval Dataset(2017)

In table 3 we have compared our results with averaged word embedding baseline [42] on monolingual(English) test dataset [13]. We notice that all our results are above the baseline scores.

Experiments	Scores
word_vec_baseline [42]	55.8
jc	73.0
lin	72.4
jc_lin	72.7
wv_wwv	71.4
jc_lin_wv_wwv	72.7

Table 3: Pearson Correlation $\times 100$ Scores On Test Dataset 2017

In table 4 we have shown benchmark results of STS 2017[13] on test dataset (with the techniques which they have adopted) and compared with our top performing experiment.

We notice that except [14] any of our results are not above rest of the benchmark results.

Experiments	Scores
jc	73.0
Ensemble [48]	81.0
Wordnet + Embeddinglin [50]	80.9
Ensemble[30]	79.2
CNN[44]	78.4
Doc2Vec[14]	59.2

Table 4: STS 2017 benchmark

4.3. CNNs And FCNN

The CNN layer which is then followed by an FCNN to calculate the scores for similarity between a pair of sentences is shown in the figure 4. The architecture has been taken from [44]. Each sentence is first represented in a matrix form, by passing it into word embedding layer. Then each sentence matrix is given as an input to the convolution layer and we get two vectors v_1 and v_2 , after the *max-pooling operation*, for both the sentences respectively. At last, by the concatenation of element-wise absolute difference and the element-wise multiplication between the two representations of input sentences v_1 and v_2 , a *Semantic Difference Vector* ($S\vec{D}V$) is created. This $S\vec{D}V$ is finally passed through the Full-Connected Layer (FCNN) to get the final output.

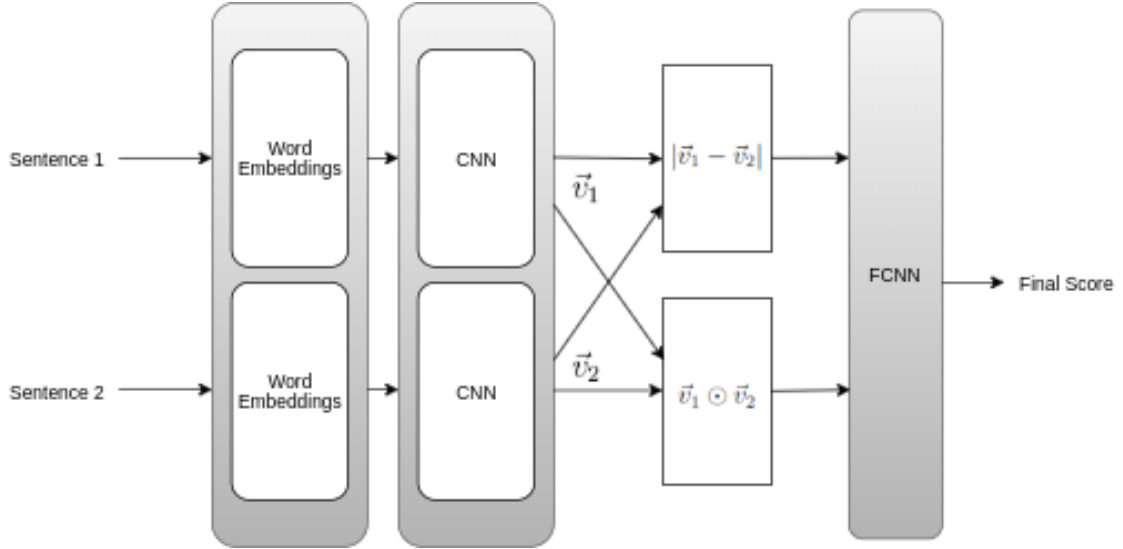


Figure 4: CNN Architecture [44]. Note that there is only one CNN in the shaded region onto which both sentences are passed.

During training, we considered the task of semantic relatedness as a classification problem instead of a regression problem, by considering the outputs from the last layer as

probability distribution over 6 classes. And the gold labels (referred as a in 12) which were real values in range 0-5 were converted into a probability distribution as follows [47]:

$$p_i = \begin{cases} a - \lfloor a \rfloor, & i = \lfloor a \rfloor + 1 \\ \lfloor a \rfloor - a + 1, & i = \lfloor a \rfloor \\ 0 & otherwise \end{cases} \quad (12)$$

And then used the loss functions (section 3.3.1) to compute the errors, which were then backpropogated for training the weights. **During evaluation** a weighted sum or expectation of the scores was calculated, which was considered as the final score.

4.3.1. Reason for preferring CNNs over RNNS or LSTMs

The model proposed by [44] (which we have used) is an extension of Siamese LSTM [38] with two major changes. First, rather than using an LSTM, a CNN is used in their work. Second, in [38] an LSTM model is employed for each input sentence and weights are shared between layers and updated parallelly during learning phase, whereas, in [44] two sentences are given as an input into the same CNN model and then the weights are shared. The reason for the first, i.e, using CNNs for training our model is that though RNNs or LSTMs works better on dynamic length of inputs sentences but if the input length is fixed, CNNs generally performs better. And in the STS 2017 dataset(for English sentence pairs) which we have considered in our work, the lenghts of the sentences among each other does not vary much, i.e., the standard deviation of the lengths of sentences is less so CNNs will work better here. To compare both LSTMs and CNNs on this dataset, we can compare [44] work with [8]. The whole architecture is exactly the same except that, [44] used CNNs in their work whereas [8] used LSTMs. And we can clearly conclude from their work that CNNs are working much better than LSTMS, where CNNs and LSTMs got a *Pearson Correlation* of **0.78** and **0.47** on English test dataset respectively. That is the reason we have also used CNNs here.

4.3.2. Experiments

[44] have used 300 dimensional GloVe word embeddings for the representation of each word in a sentence and Pearson Correlation Coefficient as the loss function. We have tried five different modifications in their work by using different loss functions (section 3.3.1) and pre-trained word embeddings. The abbreviations and the interpretation of them are as follows:

glove_kl: In this experiment we have replaced loss function correlation coefficient with KL-Divergence loss function.

word2vec_coef: In this experiment we have replaced pretrained GloVe word embeddings with pretrained word2vec embeddings.

word2vec_kl: Here both embeddings and loss functions have been replaced with *word2vec* and *KL-Divergence* respectively.

fasttext_kl: Here also, both embeddings and loss functions have been replaced with *word2vec* and *KL-Divergence* respectively.

fasttext_coef: In this experiment embeddings are replaced with fasttext embeddings, keeping the loss function same.

Each experiment discussed above produces a score between the two sentences, expressing the semantic relatedness between them. These scores, corresponding to every English sentence pair from the SemEval 2017 test dataset [13] are compared with Gold Standard scores using person correlation coefficient 3.3.1, which are shown in table 5.

Experiments	Scores
(shao) [44]	78.4
glove_kl	70.26
word2vec_coef	76.9
word2vec_kl	70.10
fasttext_kl	72.5
fasttext_coef	73.4

Table 5: (Pearson Correlation Scores $\times 100$) On Test Dataset(2017) on English Pairs

We observe that none of our systems were able to beat the scores of [44], which proves that GloVe word embeddings with correlation coefficient as loss function was certainly a good combination for this architecture.

4.4. Whole architecture

The results in this section are computed by finding *Pearson Correlation Coefficient* between the output scores from our proposed architecture (figure 2) and the gold standard labels. The table 6 shows the results when the combined outputs of *CNN* 4.3.2 and *Knowledge Based and Corpus Based (KACB)* 4.2.3 is passed into an MLP regressor and a *Pearson Correlation Coefficient* is computed between output of the scores from MLP regressor and the gold standard labels.

MLP Regressor	jc	lin	jc_lin	wv_wwv	jc_lin_wv_wwv	[43]
(glove_coef)[43]	0.787	0.742	0.731	0.742	0.752	0.786
glove_kl	0.761	0.752	0.732	0.733	0.745	0.734
word2vec_coef	0.731	0.725	0.713	0.713	0.716	0.715
word2vec_kl	0.721	0.724	0.712	0.715	0.70	0.730
fasttext_kl	0.732	0.721	0.716	0.712	0.708	0.730
fasttext_coef	0.736	0.727	0.713	0.712	0.715	0.732

Table 6: Pearson Correlation Values b/w MLP regressor Scores and Gold labels

The table 7 shows the results when the combined outputs of *CNN* and *Knowledge Based and Corpus Based (KACB)* is passed into an MLP regressor and a *Pearson Correlation*

SVR	jc	lin	jc_lin	wv_wwv	jc_lin_wv_wwv	[43]
(shao-2017)[44]	0.771	0.762	0.761	0.752	0.753	0.770
glove_kl	0.731	0.721	0.745	0.752	0.742	0.751
word2vec_coef	0.733	0.723	0.751	0.743	0.742	0.751
word2vec_kl	0.761	0.752	0.761	0.752	0.754	0.753
fasttext_kl	0.751	0.752	0.730	0.731	0.732	0.752
fasttext_coef	0.763	0.752	0.752	0.741	0.731	0.762

Table 7: Pearson Correlation Values b/w SVR Scores and Gold labels

Coefficient is computed between output of the scores from SVR and the gold standard labels.

We have observed from table 6 that the best performance is given when CNN architecture as proposed by [44] is combined with our proposed KACB technique (*jc*). In table 8 we compared our top performing architecture with the benchmark results as per SemEval 2017 test dataset on English sentence pairs [13] and noticed that we are getting slightly better results than [44].

Experiments	Scores
Ensemble [48]	81.0
Wordnet + Embeddinglin [50]	80.9
Ensemble[30]	79.2
glove_coef & <i>jc</i>	78.7
CNN[44]	78.4
Doc2Vec[14]	59.2

Table 8: STS 2017 benchmark. The bold score is our score on the benchmark database.

5. Conclusion

With some modifications in *Knowledge and Corpus Based* techniques (as seen in section 4.2) we were able to improve scores for measuring semantic similarity between sentences (as per STS 2012 test dataset for English pairs [4]), as seen in table 2. But for the 2017 STS dataset [13], except [14] our scores were not able to beat other benchmark results (table 4).

Some modifications in CNN architecture given by [44] (as seen in section 4.3.2), like changing the word embeddings, loss functions and hidden layers or unit sizes did not really bring any improvement in measuring similarity between sentences (table 5). Which therefore proved that the combination of word embeddings and the loss function as suggested in [44] (i.e. *pretrained GLoVe vectors* as word embedding [39] and *Pearson Correlation Coefficient* as loss function 3.3.1) was actually the best among all the experiments we carried out (table 5).

Our final hybrid architecture (figure 2) was able to perform much better than [14] and

slightly better than [44] which are part of the benchmark results 2017 (table 8). So, this shows that by combining the traditional *Knowledge and Corpus Based* techniques with the recent techniques with DNNs, we can extract more relevant feature set out of the sentence pairs, to compute the similarity. Therefore, in the future work, we will include some more traditional techniques, like using word alignment [45] to improve the capability of our system to measure similarity between the two sentences.

References

- [1] Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Weiwei Guo, Inigo Lopez-Gazpio, Montse Maritxalar, Rada Mihalcea, et al. Semeval-2015 task 2: Semantic textual similarity, english, spanish and pilot on interpretability. In *Proceedings of the 9th international workshop on semantic evaluation (SemEval 2015)*, pages 252–263, 2015.
- [2] Eneko Agirre, Carmen Banea, Claire Cardie, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, Weiwei Guo, Rada Mihalcea, German Rigau, and Janyce Wiebe. Semeval-2014 task 10: Multilingual semantic textual similarity. In *Proceedings of the 8th international workshop on semantic evaluation (SemEval 2014)*, pages 81–91, 2014.
- [3] Eneko Agirre, Daniel Cer, Mona Diab, Aitor Gonzalez-Agirre, and Weiwei Guo. * sem 2013 shared task: Semantic textual similarity. In *Second Joint Conference on Lexical and Computational Semantics (*SEM), Volume 1: Proceedings of the Main Conference and the Shared Task: Semantic Textual Similarity*, pages 32–43, 2013.
- [4] Eneko Agirre, Mona Diab, Daniel Cer, and Aitor Gonzalez-Agirre. Semeval-2012 task 6: A pilot on semantic textual similarity. In *Proceedings of the First Joint Conference on Lexical and Computational Semantics, SemEval ’12*, pages 385–393, Stroudsburg, PA, USA, 2012. Association for Computational Linguistics.
- [5] James Allen. *Natural Language Understanding*. Benjamin-Cummings Publishing Co., Inc., USA, 1988.
- [6] J. Atkinson-Abutridy, C. Mellish, and S. Aitken. Combining information extraction with genetic algorithms for text mining. *IEEE Intelligent Systems*, 19(3):22–30, 2004.
- [7] Daniel Bär, Chris Biemann, Iryna Gurevych, and Torsten Zesch. Ukp: Computing semantic textual similarity by combining multiple content similarity measures. In *Proceedings of the First Joint Conference on Lexical and Computational Semantics, SemEval ’12*, pages 435–440, Stroudsburg, PA, USA, 2012. Association for Computational Linguistics.
- [8] Joe Barrow and Denis Peskov. UMDeep at SemEval-2017 task 1: End-to-end shared weight LSTM model for semantic textual similarity. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 180–184, Vancouver, Canada, August 2017. Association for Computational Linguistics.
- [9] Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, Lisbon, Portugal, September 2015. Association for Computational Linguistics.
- [10] Tomáš Brychcín and Lukáš Svoboda. Uwb at semeval-2016 task 1: Semantic textual similarity using lexical, syntactic, and semantic information. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, pages 588–594, 2016.
- [11] Chris Callison-Burch, Cameron Fordyce, Philipp Koehn, Christof Monz, and Josh Schroeder. (meta-) evaluation of machine translation. In *Proceedings of the Second Workshop on Statistical Machine Translation*, pages 136–158. Association for Computational Linguistics, 2007.
- [12] Chris Callison-Burch, Cameron Fordyce, Philipp Koehn, Christof Monz, and Josh Schroeder. Further meta-evaluation of machine translation. In *Proceedings of the third workshop on statistical machine translation*, pages 70–106. Association for Computational Linguistics, 2008.

- [13] Daniel Cer, Mona Diab, Eneko Agirre, Iñigo Lopez-Gazpio, and Lucia Specia. SemEval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 1–14, Vancouver, Canada, August 2017. Association for Computational Linguistics.
- [14] Mirela-Stefania Duma and Wolfgang Menzel. SEF@UHH at SemEval-2017 task 1: Unsupervised knowledge-free semantic textual similarity via paragraph vector. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 170–174, Vancouver, Canada, August 2017. Association for Computational Linguistics.
- [15] C. Fellbaum and G. Miller. *Combining Local Context and Wordnet Similarity for Word Sense Identification*, pages 265–283. 1998.
- [16] Peter W. Foltz, Walter Kintsch, and Thomas K Landauer. The measurement of textual coherence with latent semantic analysis. *Discourse Processes*, 25(2-3):285–307, 1998.
- [17] Evgeniy Gabrilovich and Shaul Markovitch. Computing semantic relatedness using wikipedia-based explicit semantic analysis. In *Proceedings of the 20th International Joint Conference on Artificial Intelligence, IJCAI07*, page 16061611, San Francisco, CA, USA, 2007. Morgan Kaufmann Publishers Inc.
- [18] Juri Ganitkevitch, Benjamin Van Durme, and Chris Callison-Burch. PPDB: The paraphrase database. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 758–764, Atlanta, Georgia, June 2013. Association for Computational Linguistics.
- [19] Vasileios Hatzivassiloglou, Judith L. Klavans, and Eleazar Eskin. Detecting text similarity over short passages: Exploring linguistic feature combinations via machine learning. In *1999 Joint SIGDAT Conference on Empirical Methods in Natural Language Processing and Very Large Corpora*, 1999.
- [20] Hua He, Kevin Gimpel, and Jimmy Lin. Multi-perspective sentence similarity modeling with convolutional neural networks. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1576–1586, Lisbon, Portugal, September 2015. Association for Computational Linguistics.
- [21] Hua He and Jimmy Lin. Pairwise word interaction modeling with deep neural networks for semantic similarity measurement. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 937–948, San Diego, California, June 2016. Association for Computational Linguistics.
- [22] Hua He, John Wieting, Kevin Gimpel, Jinfeng Rao, and Jimmy Lin. UMD-TTIC-UW at SemEval-2016 task 1: Attention-based multi-perspective convolutional neural networks for textual similarity measurement. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, pages 1103–1108, San Diego, California, June 2016. Association for Computational Linguistics.
- [23] Jay J. Jiang and David W. Conrath. Semantic similarity based on corpus statistics and lexical taxonomy. In *Proceedings of the 10th Research on Computational Linguistics International Conference*, pages 19–33, Taipei, Taiwan, August 1997. The Association for Computational Linguistics and Chinese Language Processing (ACLCLP).
- [24] Thomas K Landauer, Peter W. Foltz, and Darrell Laham. An introduction to latent semantic analysis. *Discourse Processes*, 25(2-3):259–284, 1998.
- [25] Thomas K. Landauer, Darrell Laham, Bob Rehder, and M. E. Schreiner. How well can passage meaning be derived without using word order? a comparison of latent semantic analysis and humans. 1997.
- [26] Michael D. Lee, Brandon Pincombe, and Matthew Welsh. An empirical evaluation of models of text document similarity. 2005.
- [27] Michael Lesk. Automatic sense disambiguation using machine readable dictionaries: How to tell a pine cone from an ice cream cone. In *Proceedings of the 5th Annual International Conference on Systems Documentation, SIGDOC 86*, page 2426, New York, NY, USA, 1986. Association for Computing Machinery.
- [28] Y. Li, D. McLean, Z. A. Bandar, J. D. O’Shea, and K. Crockett. Sentence similarity based on semantic nets and corpus statistics. *IEEE Transactions on Knowledge and Data Engineering*, 18(8):1138–1150,

- 2006.
- [29] Dekang Lin. An information-theoretic definition of similarity. In *Proceedings of the Fifteenth International Conference on Machine Learning, ICML '98*, pages 296–304, San Francisco, CA, USA, 1998. Morgan Kaufmann Publishers Inc.
 - [30] Nabin Maharjan, Rajendra Banjade, Dipesh Gautam, Lasang J. Tamang, and Vasile Rus. 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, Vancouver, Canada, August 2017. Association for Computational Linguistics.
 - [31] Marco Marelli, Stefano Menini, Marco Baroni, Luisa Bentivogli, Raffaella Bernardi, and Roberto Zamparelli. A SICK cure for the evaluation of compositional distributional semantic models. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 216–223, Reykjavik, Iceland, May 2014. European Language Resources Association (ELRA).
 - [32] Charles T Meadow, Donald H Kraft, and Bert R Boyce. *Text information retrieval systems*. Academic Press, Inc., 1999.
 - [33] Rada Mihalcea, Courtney Corley, and Carlo Strapparava. Corpus-based and knowledge-based measures of text semantic similarity. In *Proceedings of the 21st National Conference on Artificial Intelligence - Volume 1, AAAI06*, page 775780. AAAI Press, 2006.
 - [34] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
 - [35] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. Distributed representations of words and phrases and their compositionality. In *Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 2, NIPS'13*, pages 3111–3119, USA, 2013. Curran Associates Inc.
 - [36] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 3111–3119. Curran Associates, Inc., 2013.
 - [37] George A. Miller. Wordnet: A lexical database for english. *Commun. ACM*, 38(11):3941, November 1995.
 - [38] Jonas Mueller and Aditya Thyagarajan. Siamese recurrent architectures for learning sentence similarity. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, AAAI16*, page 27862792. AAAI Press, 2016.
 - [39] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. Glove: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, 2014.
 - [40] Philip Resnik. Using information content to evaluate semantic similarity in a taxonomy. In *Proceedings of the 14th International Joint Conference on Artificial Intelligence - Volume 1, IJCAI95*, page 448453, San Francisco, CA, USA, 1995. Morgan Kaufmann Publishers Inc.
 - [41] Philip Resnik. Using information content to evaluate semantic similarity in a taxonomy. In *Proceedings of the 14th International Joint Conference on Artificial Intelligence - Volume 1, IJCAI'95*, pages 448–453, San Francisco, CA, USA, 1995. Morgan Kaufmann Publishers Inc.
 - [42] Alexandre Salle, Aline Villavicencio, and Marco Idiart. Matrix factorization using window sampling and negative sampling for improved word representations. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 419–424, Berlin, Germany, August 2016. Association for Computational Linguistics.
 - [43] Frane Šarić, Goran Glavaš, Mladen Karan, Jan Šnajder, and Bojana Dalbelo Bašić. TakeLab: Systems for measuring semantic text similarity. In **SEM 2012: The First Joint Conference on Lexical and Computational Semantics – (SemEval 2012)*, pages 441–448, Montréal, Canada, 7-8 June 2012. Association for Computational Linguistics.
 - [44] Yang Shao. 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, Vancouver, Canada, August 2017. Association for Computational Linguistics.
- [45] Md Arafat Sultan, Steven Bethard, and Tamara Sumner. Back to basics for monolingual alignment: Exploiting word similarity and contextual evidence. *Transactions of the Association for Computational Linguistics*, 2:219–230, 2014.
 - [46] Md Arafat Sultan, Steven Bethard, and Tamara Sumner. DLS@CU: Sentence similarity from word alignment and semantic vector composition. In *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*, pages 148–153, Denver, Colorado, June 2015. Association for Computational Linguistics.
 - [47] Kai Sheng Tai, Richard Socher, and Christopher D. Manning. 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)*, pages 1556–1566, Beijing, China, July 2015. Association for Computational Linguistics.
 - [48] Junfeng Tian, Zhiheng Zhou, Man Lan, and Yuanbin Wu. ECNU 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, Vancouver, Canada, August 2017. Association for Computational Linguistics.
 - [49] Peter D. Turney. Mining the web for synonyms: Pmi-ir versus lsa on toefl. In Luc De Raedt and Peter Flach, editors, *Machine Learning: ECML 2001*, pages 491–502, Berlin, Heidelberg, 2001. Springer Berlin Heidelberg.
 - [50] Hao Wu, Heyan Huang, Ping Jian, Yuhang Guo, and Chao Su. 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, Vancouver, Canada, August 2017. Association for Computational Linguistics.
 - [51] Zhibiao Wu and Martha Palmer. Verbs semantics and lexical selection. In *Proceedings of the 32nd Annual Meeting on Association for Computational Linguistics*, ACL 94, page 133138, USA, 1994. Association for Computational Linguistics.
 - [52] Wei Xu, Chris Callison-Burch, and Bill Dolan. SemEval-2015 task 1: Paraphrase and semantic similarity in twitter (PIT). In *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*, pages 1–11, Denver, Colorado, June 2015. Association for Computational Linguistics.
 - [53] Ye Zhang and Byron Wallace. A sensitivity analysis of (and practitioners’ guide to) convolutional neural networks for sentence classification. *arXiv preprint arXiv:1510.03820*, 2015.