

# ABCNN: Attention-Based Convolutional Neural Network for Modeling Sentence Pairs

Wenpeng Yin, Hinrich Schütze  
Center for Information and Language Processing  
LMU Munich, Germany  
wenpeng@cis.lmu.de

Bing Xiang, Bowen Zhou  
IBM Watson  
Yorktown Heights, NY, USA  
bingxia, zhou@us.ibm.com

## Abstract

How to model a pair of sentences is a critical issue in many NLP tasks such as answer selection (AS), paraphrase identification (PI) and textual entailment (TE). Most prior work (i) deals with one individual task by fine-tuning a specific system; (ii) models each sentence’s representation separately, rarely considering the impact of the other sentence; or (iii) relies fully on manually designed, task-specific linguistic features. This work presents a general Attention Based Convolutional Neural Network (ABCNN) for modeling a pair of sentences. We make three contributions. (i) The ABCNN can be applied to a wide variety of tasks that require modeling of sentence pairs. (ii) We propose three attention schemes that integrate mutual influence between sentences into CNNs; thus, the representation of each sentence takes into consideration its counterpart. These interdependent sentence pair representations are more powerful than isolated sentence representations. (iii) ABCNNs achieve state-of-the-art performance on AS, PI and TE tasks. We release code at: [https://github.com/yinwenpeng/Answer\\_Selection](https://github.com/yinwenpeng/Answer_Selection).

## 1 Introduction

How to model a pair of sentences is a critical issue in many NLP tasks such as answer selection (AS) (Yu et al., 2014; Feng et al., 2015), paraphrase identification (PI) (Madnani et al., 2012; Yin and Schütze, 2015a), textual entailment (TE) (Marelli et al., 2014a; Bowman et al., 2015a) etc.

AS	$s_0$	how much did Waterboy gross?
	$s_1^+$	the movie <i>earned</i> \$161.5 million
	$s_1^-$	this was Jerry Reed’s final film appearance
PI	$s_0$	she struck a deal with RH to pen a book <i>today</i>
	$s_1^+$	she signed a contract with RH to write a book
	$s_1^-$	she denied <i>today</i> that she struck a deal with RH
TE	$s_0$	an ice skating rink placed <i>outdoors</i> is <i>full of people</i>
	$s_1^+$	a <i>lot of people</i> are in an ice skating park
	$s_1^-$	an ice skating rink placed <i>indoors</i> is <i>full of people</i>

Figure 1: Positive ( $\langle s_0, s_1^+ \rangle$ ) and negative ( $\langle s_0, s_1^- \rangle$ ) examples for AS, PI and TE tasks. RH = Random House

Most prior work derives each sentence’s representation separately, rarely considering the impact of the other sentence. This neglects the mutual influence of the two sentences in the context of the task. It also contradicts what humans do when comparing two sentences. We usually focus on key parts of one sentence by extracting parts from the other sentence that are related by identity, synonymy, antonymy and other relations. Thus, human beings model the two sentences together, using the content of one sentence to guide the representation of the other.

Figure 1 demonstrates that each sentence of a pair partially determines which parts of the other sentence we must focus on. For AS, correctly answering  $s_0$  requires attention on “gross”:  $s_1^+$  contains a corresponding unit (“earned”) while  $s_1^-$  does not. For PI, focus should be removed from “today” to correctly recognize  $\langle s_0, s_1^+ \rangle$  as paraphrases and  $\langle s_0, s_1^- \rangle$  as non-paraphrases. For TE, we need to focus on “full of people” (to recognize TE for  $\langle s_0, s_1^+ \rangle$ ) and on “outdoors” / “indoors” (to recognize non-TE for  $\langle s_0, s_1^- \rangle$ ). These examples show the need for an architecture that computes different representations of  $s_i$  for different  $s_{1-i}$  ( $i \in \{0, 1\}$ ).

Convolutional Neural Networks (CNNs) (LeCun et al., 1998) are widely used to model sentences (Kalchbrenner et al., 2014; Kim, 2014) and sentence pairs (Socher et al., 2011; Yin and Schütze, 2015a), especially in classification tasks. CNNs are supposed to be good at extracting robust and abstract features of input. This work presents the ABCNN, an attention-based convolutional neural network, that has a powerful mechanism for modeling a sentence pair by taking into account the interdependence between the two sentences. The ABCNN is a general architecture that can handle a wide variety of sentence pair modeling tasks.

Some prior work proposes simple mechanisms that can be interpreted as controlling varying attention; e.g., Yih et al. (2013) employ word alignment to match related parts of the two sentences. In contrast, our attention scheme based on CNNs models relatedness between two parts fully automatically. Moreover, attention at multiple levels of granularity, not only at word level, is achieved as we stack multiple convolution layers that increase abstraction.

Prior work on attention in deep learning (DL) mostly addresses long short-term memory networks (LSTMs) (Hochreiter and Schmidhuber, 1997). LSTMs achieve attention usually in a word-to-word scheme, and word representations mostly encode the *whole context* within the sentence (Bahdanau et al., 2015; Rocktäschel et al., 2016). It is not clear whether this is the best strategy; e.g., in the AS example in Figure 1, it is possible to determine that “how much” in  $s_0$  matches “\$161.5 million” in  $s_1$  without taking the entire sentence contexts into account. This observation was also investigated by Yao et al. (2013b) where an information retrieval system retrieves sentences with tokens labeled as DATE by named entity recognition or as CD by POS tagging if there is a “when” question. However, labels or POS tags require extra tools. CNNs benefit from incorporating attention into representations of *local phrases* detected by filters; in contrast, LSTMs encode the *whole context* to form attention-based word representations – a strategy that is more complex than the CNN strategy and (as our experiments suggest) performs less well for some tasks.

Apart from these differences, it is clear that attention has as much potential for CNNs as it does for LSTMs. As far as we know, this is the first NLP

paper that incorporates attention into CNNs. Our ABCNNs get state-of-the-art in AS and TE tasks, and competitive performance in PI, then obtains further improvements over all three tasks when linguistic features are used.

## 2 Related Work

**Non-DL on Sentence Pair Modeling.** Sentence pair modeling has attracted lots of attention in the past decades. Many tasks can be reduced to a semantic text matching problem. Due to the variety of word choices and inherent ambiguities in natural language, bag-of-word approaches with simple surface-form word matching tend to produce brittle results with poor prediction accuracy (Bilotti et al., 2007). As a result, researchers put more emphasis on exploiting syntactic and semantic structure. Representative examples include methods based on deeper semantic analysis (Shen and Lapata, 2007; Moldovan et al., 2007), tree edit-distance (Punyakanok et al., 2004; Heilman and Smith, 2010) and quasi-synchronous grammars (Wang et al., 2007) that match the dependency parse trees of the two sentences. Instead of focusing on the high-level semantic representation, Yih et al. (2013) turn their attention to improving the shallow semantic component, lexical semantics, by performing semantic matching based on a latent word-alignment structure (cf. Chang et al. (2010)). Lai and Hockenmaier (2014) explore finer-grained word overlap and alignment between two sentences using negation, hypernym, synonym and antonym relations. Yao et al. (2013a) extend word-to-word alignment to phrase-to-phrase alignment by a semi-Markov CRF. However, such approaches often require more computational resources. In addition, employing syntactic or semantic parsers – which produce errors on many sentences – to find the best match between the structured representations of two sentences is not trivial.

**DL on Sentence Pair Modeling.** To address some of the challenges of non-DL work, much recent work uses neural networks to model sentence pairs for AS, PI and TE.

For AS, Yu et al. (2014) present a bigram CNN to model question and answer candidates. Yang et al. (2015) extend this method and get state-of-the-art performance on the WikiQA dataset (Section 5.1).

Feng et al. (2015) test various setups of a bi-CNN architecture on an insurance domain QA dataset. Tan et al. (2016) explore bidirectional LSTMs on the same dataset. Our approach is different because we do not model the sentences by two independent neural networks in parallel, but instead as an interdependent sentence pair, using attention.

For PI, Blacoe and Lapata (2012) form sentence representations by summing up word embeddings. Socher et al. (2011) use recursive autoencoders (RAEs) to model representations of local phrases in sentences, then pool similarity values of phrases from the two sentences as features for binary classification. Yin and Schütze (2015a) similarly replace an RAE with a CNN. In all three papers, the representation of one sentence is not influenced by the other – in contrast to our attention-based model.

For TE, Bowman et al. (2015b) use recursive neural networks to encode entailment on SICK (Marelli et al., 2014b). Rocktäschel et al. (2016) present an attention-based LSTM for the Stanford natural language inference corpus (Bowman et al., 2015a). Our system is the first CNN-based work on TE.

Some prior work aims to solve a general sentence matching problem. Hu et al. (2014) present two CNN architectures, ARC-I and ARC-II, for sentence matching. ARC-I focuses on sentence representation learning while ARC-II focuses on matching features on phrase level. Both systems were tested on PI, sentence completion (SC) and tweet-response matching. Yin and Schütze (2015b) propose the MultiGranCNN architecture to model general sentence matching based on phrase matching on multiple levels of granularity and get promising results for PI and SC. Wan et al. (2016) try to match two sentences in AS and SC by multiple sentence representations, each coming from the local representations of two LSTMs. Our work is the first one to investigate attention for the general sentence matching task.

**Attention-Based DL in Non-NLP Domains.** Even though there is little if any work on attention mechanisms in CNNs for NLP, attention-based CNNs have been used in computer vision for visual question answering (Chen et al., 2015), image classification (Xiao et al., 2015), caption generation (Xu et al., 2015), image segmentation (Hong et al., 2016) and object localization (Cao et al., 2015).

symbol	description
$s, s_0, s_1$	sentence or sentence length
$v$	word
$w$	filter width
$d_i$	dimensionality of input to layer $i + 1$
$\mathbf{W}$	weight matrix

Table 1: Notation

Mnih et al. (2014) apply attention in recurrent neural networks (RNNs) to extract information from an image or video by adaptively selecting a sequence of regions or locations and only processing the selected regions at high resolution. Gregor et al. (2015) combine a spatial attention mechanism with RNNs for image generation. Ba et al. (2015) investigate attention-based RNNs for recognizing multiple objects in images. Chorowski et al. (2014) and Chorowski et al. (2015) use attention in RNNs for speech recognition.

**Attention-Based DL in NLP.** Attention-based DL systems have been applied to NLP after their success in computer vision and speech recognition. They mainly rely on RNNs and end-to-end encoder-decoders for tasks such as machine translation (Bahdanau et al., 2015; Luong et al., 2015) and text reconstruction (Li et al., 2015; Rush et al., 2015). Our work takes the lead in exploring attention mechanisms in CNNs for NLP tasks.

### 3 BCNN: Basic Bi-CNN

We now introduce our basic (non-attention) CNN that is based on the Siamese architecture (Bromley et al., 1993), i.e., it consists of two weight-sharing CNNs, each processing one of the two sentences, and a final layer that solves the sentence pair task. See Figure 2. We refer to this architecture as the *BCNN*. The next section will then introduce the ABCNN, an attention architecture that extends the BCNN. Table 1 gives our notational conventions.

In our implementation and also in the mathematical formalization of the model given below, we pad the two sentences to have the same length  $s = \max(s_0, s_1)$ . However, in the figures we show different lengths because this gives a better intuition of how the model works.

We now describe the BCNN’s four types of layers: input, convolution, average pooling and output.

**Input layer.** In the example in the figure, the two input sentences have 5 and 7 words, respectively.

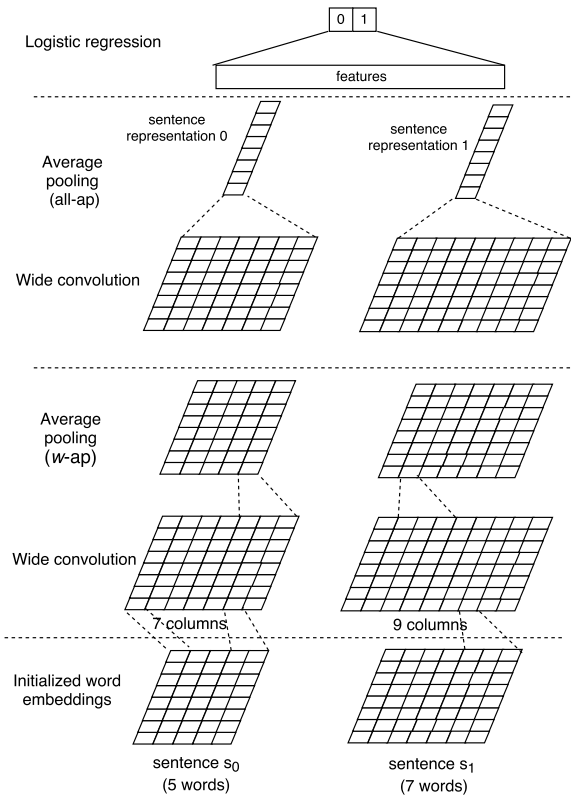


Figure 2: BCNN: ABCNN without Attention

Each word is represented as a  $d_0$ -dimensional pre-computed word2vec (Mikolov et al., 2013) embedding,  $d_0 = 300$ . As a result, each sentence is represented as a feature map of dimension  $d_0 \times s$ .

**Convolution layer.** Let  $v_1, v_2, \dots, v_s$  be the words of a sentence and  $\mathbf{c}_i \in \mathbb{R}^{w \cdot d_0}$ ,  $0 < i < s + w$ , the concatenated embeddings of  $v_{i-w+1}, \dots, v_i$  where embeddings for  $v_j$  are set to zero when  $j < 1$  or  $j > s$ . We then generate the representation  $\mathbf{p}_i \in \mathbb{R}^{d_1}$  for the phrase  $v_{i-w+1}, \dots, v_i$  using the convolution weights  $\mathbf{W} \in \mathbb{R}^{d_1 \times w d_0}$  as follows:

$$\mathbf{p}_i = \tanh(\mathbf{W} \cdot \mathbf{c}_i + \mathbf{b})$$

where  $\mathbf{b} \in \mathbb{R}^{d_1}$  is the bias.

**Average pooling layer.** Pooling (including min, max, average pooling) is commonly used to extract robust features from convolution. In this paper, we introduce attention weighting as an alternative, but use average pooling as a baseline as follows.

For the output feature map of the last convolution layer, we do column-wise averaging over *all columns*, denoted as *all-ap*. This generates a representation vector for each of the two sentences, shown as the top “Average pooling (*all-ap*)” layer

below “Logistic regression” in Figure 2. These two vectors are the basis for the sentence pair decision.

For the output feature map of non-final convolution layers, we do column-wise averaging over *windows of w consecutive columns*, denoted as *w-ap*; shown as the lower “Average pooling (*w-ap*)” layer in Figure 2. For filter width  $w$ , a convolution layer transforms an input feature map of  $s$  columns into a new feature map of  $s + w - 1$  columns; average pooling transforms this back to  $s$  columns. This architecture supports stacking an arbitrary number of convolution-pooling blocks to extract increasingly abstract features. Input features to the bottom layer are words, input features to the next layer are short phrases and so on. Each level generates more abstract features of higher granularity.

The last layer is an **output layer**, chosen according to the task; e.g., for binary classification tasks, this layer is logistic regression (see Figure 2). Other types of output layers are introduced below.

We found that in most cases, performance is boosted if we provide the output of *all pooling layers* as input to the output layer. For each non-final average pooling layer, we perform *w-ap* (pooling over windows of  $w$  columns) as described above, but we also perform *all-ap* (pooling over all columns) and forward the result to the output layer. This improves performance because representations from different layers cover the properties of the sentences at different levels of abstraction and all of these levels can be important for a particular sentence pair.

#### 4 ABCNN: Attention-Based BCNN

We now describe three architectures based on the BCNN, the ABCNN-1, the ABCNN-2 and the ABCNN-3, that each introduces an attention mechanism for modeling sentence pairs; see Figure 3.

**ABCNN-1.** The ABCNN-1 (Figure 3(a)) employs an attention feature matrix  $\mathbf{A}$  to influence convolution. Attention features are intended to weight those units of  $s_i$  more highly in convolution that are relevant to a unit of  $s_{1-i}$  ( $i \in \{0, 1\}$ ); we use the term “unit” here to refer to words on the lowest level and to phrases on higher levels of the network. Figure 3(a) shows two *unit representation feature maps* in red: this part of the ABCNN-1 is the same as in the BCNN (see Figure 2). Each column is the

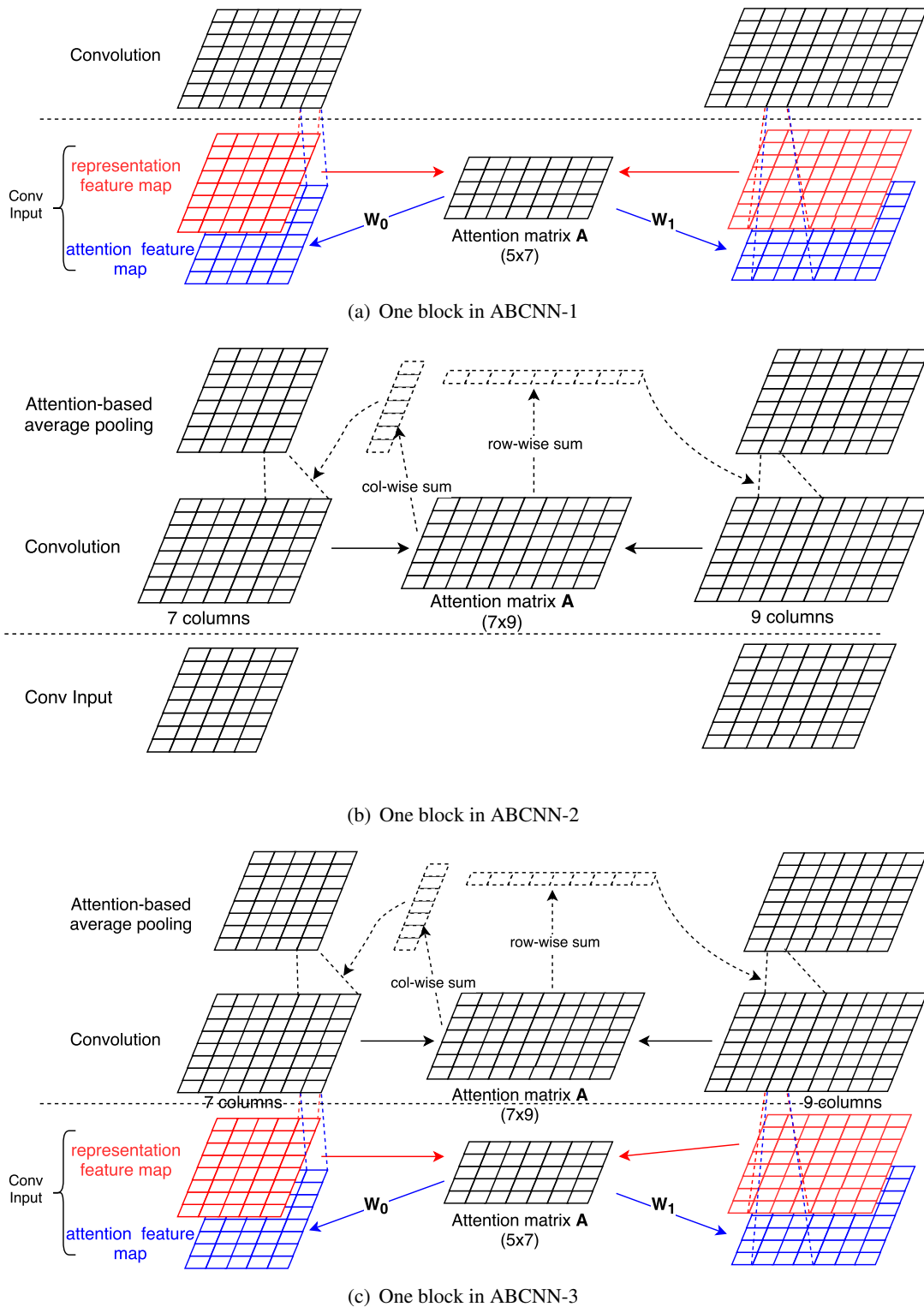


Figure 3: Three ABCNN architectures

representation of a unit, a word on the lowest level and a phrase on higher levels. We first describe the attention feature matrix  $\mathbf{A}$  informally (layer “Conv input”, middle column, in Figure 3(a)).  $\mathbf{A}$  is generated by matching units of the left representation feature map with units of the right representation feature map such that the attention values of row  $i$  in  $\mathbf{A}$  denote the attention distribution of the  $i$ -th unit of  $s_0$  with respect to  $s_1$ , and the attention values of column  $j$  in  $\mathbf{A}$  denote the attention distribution of the  $j$ -th unit of  $s_1$  with respect to  $s_0$ .  $\mathbf{A}$  can be viewed as a new feature map of  $s_0$  (resp.  $s_1$ ) in row (resp. column) direction because each row (resp. column) is a new feature vector of a unit in  $s_0$  (resp.  $s_1$ ). Thus, it makes sense to combine this new feature map with the representation feature maps and use both as input to the convolution operation. We achieve this by transforming  $\mathbf{A}$  into the two blue matrices in Figure 3(a) that have the same format as the representation feature maps. As a result, the new input of convolution has two feature maps for each sentence (shown in red and blue). Our motivation is that the attention feature map will guide the convolution to learn “counterpart-biased” sentence representations.

More formally, let  $\mathbf{F}_{i,r} \in \mathbf{R}^{d \times s}$  be the *representation feature map* of sentence  $i$  ( $i \in \{0, 1\}$ ). Then we define the attention matrix  $\mathbf{A} \in \mathbf{R}^{s \times s}$  as follows:

$$\mathbf{A}_{i,j} = \text{match-score}(\mathbf{F}_{0,r}[:, i], \mathbf{F}_{1,r}[:, j]) \quad (1)$$

The function match-score can be defined in a variety of ways. We found that  $1/(1 + |x - y|)$  works well where  $|\cdot|$  is Euclidean distance.

Given attention matrix  $\mathbf{A}$ , we generate the *attention feature map*  $\mathbf{F}_{i,a}$  for  $s_i$  as follows:

$$\mathbf{F}_{0,a} = \mathbf{W}_0 \cdot \mathbf{A}^\top, \quad \mathbf{F}_{1,a} = \mathbf{W}_1 \cdot \mathbf{A}$$

The weight matrices  $\mathbf{W}_0 \in \mathbf{R}^{d \times s}$ ,  $\mathbf{W}_1 \in \mathbf{R}^{d \times s}$  are parameters of the model to be learned in training.<sup>1</sup>

We stack the representation feature map  $\mathbf{F}_{i,r}$  and the attention feature map  $\mathbf{F}_{i,a}$  as an order 3 tensor and feed it into convolution to generate a higher-level representation feature map for  $s_i$  ( $i \in \{0, 1\}$ ). In Figure 3(a),  $s_0$  has 5 units,  $s_1$  has 7. The output of convolution (shown in the top layer, filter width

<sup>1</sup>The weights of the two matrices are shared in our implementation to reduce the number of parameters of the model.

$w = 3$ ) is a higher-level representation feature map with 7 columns for  $s_0$  and 9 columns for  $s_1$ .

**ABCNN-2.** The ABCNN-1 computes attention weights *directly on the input representation* with the aim of *improving the features computed by convolution*. The ABCNN-2 (Figure 3(b)) instead computes attention weights *on the output of convolution* with the aim of *reweighting this convolution output*. In the example shown in Figure 3(b), the feature maps output by convolution for  $s_0$  and  $s_1$  (layer marked “Convolution” in Figure 3(b)) have 7 and 9 columns, respectively; each column is the representation of a unit. The attention matrix  $\mathbf{A}$  compares all units in  $s_0$  with all units of  $s_1$ . We sum all attention values for a unit to derive a single attention weight for that unit. This corresponds to summing all values in a row of  $\mathbf{A}$  for  $s_0$  (“col-wise sum”, resulting in the column vector of size 7 shown) and summing all values in a column for  $s_1$  (“row-wise sum”, resulting in the row vector of size 9 shown).

More formally, let  $\mathbf{A} \in \mathbf{R}^{s \times s}$  be the attention matrix,  $a_{0,j} = \sum \mathbf{A}[j, :]$  the attention weight of unit  $j$  in  $s_0$ ,  $a_{1,j} = \sum \mathbf{A}[:, j]$  the attention weight of unit  $j$  in  $s_1$  and  $\mathbf{F}_{i,r}^c \in \mathbf{R}^{d \times (s_i + w - 1)}$  the output of convolution for  $s_i$ . Then the  $j$ -th column of the new feature map  $\mathbf{F}_{i,r}^p$  generated by  $w$ -ap is derived by:

$$\mathbf{F}_{i,r}^p[:, j] = \sum_{k=j:w} a_{i,k} \mathbf{F}_{i,r}^c[:, k], \quad j = 1 \dots s_i$$

Note that  $\mathbf{F}_{i,r}^p \in \mathbf{R}^{d \times s_i}$ , i.e., ABCNN-2 pooling generates an output feature map of the same size as the input feature map of convolution. This allows us to stack multiple convolution-pooling blocks to extract features of increasing abstraction.

There are three main differences between the ABCNN-1 and the ABCNN-2. (i) Attention in the ABCNN-1 impacts *convolution indirectly* while attention in the ABCNN-2 influences *pooling through direct* attention weighting. (ii) The ABCNN-1 requires the two matrices  $\mathbf{W}_i$  to convert the attention matrix into attention feature maps; and the input to convolution has two times as many feature maps. Thus, the ABCNN-1 has more parameters than the ABCNN-2 and is more vulnerable to overfitting. (iii) As pooling is performed after convolution, pooling handles larger-granularity units than convolution; e.g., if the input to convolution has word level

granularity, then the input to pooling has phrase level granularity, the phrase size being equal to filter size  $w$ . Thus, the ABCNN-1 and the ABCNN-2 implement attention mechanisms for linguistic units of different granularity. The complementarity of the ABCNN-1 and the ABCNN-2 motivates us to propose the ABCNN-3, a third architecture that combines elements of the two.

**ABCNN-3** (Figure 3(c)) combines the ABCNN-1 and the ABCNN-2 by stacking them; it combines the strengths of the ABCNN-1 and -2 by allowing the attention mechanism to operate (i) both on the convolution and on the pooling parts of a convolution-pooling block and (ii) both on the input granularity and on the more abstract output granularity.

## 5 Experiments

We test the proposed architectures on three tasks: answer selection (AS), paraphrase identification (PI) and textual entailment (TE).

**Common Training Setup.** Words are initialized by 300-dimensional word2vec embeddings and not changed during training. A single randomly initialized embedding is created for all unknown words by uniform sampling from  $[-.01, .01]$ . We employ AdaGrad (Duchi et al., 2011) and  $L_2$  regularization.

**Network Configuration.** Each network in the experiments below consists of (i) an initialization block  $b_1$  that initializes words by word2vec embeddings, (ii) a stack of  $k - 1$  convolution-pooling blocks  $b_2, \dots, b_k$ , computing increasingly abstract features, and (iii) one final *LR layer* (logistic regression layer) as shown in Figure 2.

The input to the LR layer consists of  $kn$  features – each block provides  $n$  similarity scores, e.g.,  $n$  cosine similarity scores. Figure 2 shows the two sentence vectors output by the final block  $b_k$  of the stack (“sentence representation 0”, “sentence representation 1”); this is the basis of the last  $n$  similarity scores. As we explained in the final paragraph of Section 3, we perform *all-ap* pooling for *all blocks*, not just for  $b_k$ . Thus we get one sentence representation each for  $s_0$  and  $s_1$  for each block  $b_1, \dots, b_k$ . We compute  $n$  similarity scores for each block (based on the block’s two sentence representations). Thus, we compute a total of  $kn$  similarity scores and these scores are input to the LR layer.

	#CL	AS		PI		TE	
		lr	$w L_2$	lr	$w L_2$	lr	$w L_2$
ABCNN-1	1	.08	4 .0004	.08	3 .0002	.08	3 .0006
ABCNN-1	2	.085	4 .0006	.085	3 .0003	.085	3 .0006
ABCNN-2	1	.05	4 .0003	.085	3 .0001	.09	3 .00065
ABCNN-2	2	.06	4 .0006	.085	3 .0001	.085	3 .0007
ABCNN-3	1	.05	4 .0003	.05	3 .0003	.09	3 .0007
ABCNN-3	2	.06	4 .0006	.055	3 .0005	.09	3 .0007

Table 2: Hyperparameters. lr: learning rate. #CL: number convolution layers.  $w$ : filter width. The number of convolution kernels  $d_i$  ( $i > 0$ ) is 50 throughout.

Depending on the task, we use different methods for computing the similarity score: see below.

**Layerwise Training.** In our training regime, we first train a network consisting of just one convolution-pooling block  $b_2$ . We then create a new network by adding a block  $b_3$ , initialize its  $b_2$  block with the previously learned weights for  $b_2$  and train  $b_3$  keeping the previously learned weights for  $b_2$  fixed. We repeat this procedure until all  $k - 1$  convolution-pooling blocks are trained. We found that this training regime gives us good performance and shortens training times considerably. Since similarity scores of lower blocks are kept unchanged once they have been learned, this also has the nice effect that “simple” similarity scores (those based on surface features) are learned first and subsequent training phases can focus on complementary scores derived from more complex abstract features.

**Classifier.** We found that performance increases if we do not use the output of the LR layer as the final decision, but instead train a linear SVM or a logistic regression with default parameters<sup>2</sup> directly on the input to the LR layer (i.e., on the  $kn$  similarity scores that are generated by the  $k$ -block stack after network training is completed). Direct training of SVMs/LR seems to get closer to the global optimum than gradient descent training of CNNs.

Table 2 shows hyperparameters, tuned on dev.

We use addition and LSTMs as two **shared baselines** for all three tasks, i.e., for AS, PI and TE. We now describe these two shared baselines.

(i) **Addition.** We sum up word embeddings element-wise to form each sentence representation. The classifier input is then the concatenation of the two sentence representations. (ii) **A-LSTM.** Before this work, most attention mechanisms in NLP

<sup>2</sup> <http://scikit-learn.org/stable/> for both.

were implemented in recurrent neural networks for text generation tasks such as machine translation (e.g., Bahdanau et al. (2015), Luong et al. (2015)). Rocktäschel et al. (2016) present an attention-LSTM for natural language inference. Since this model is the pioneering attention based RNN system for sentence pair classification, we consider it as a baseline system (“A-LSTM”) for all our three tasks. The A-LSTM has the same configuration as our ABCNNs in terms of word initialization (300-dimensional word2vec embeddings) and the dimensionality of all hidden layers (50).

## 5.1 Answer Selection

We use WikiQA,<sup>3</sup> an open domain question-answer dataset. We use the subtask that assumes that there is at least one correct answer for a question. The corresponding dataset consists of 20,360 question-candidate pairs in train, 1,130 pairs in dev and 2,352 pairs in test where we adopt the standard setup of only considering questions with correct answers in test. Following Yang et al. (2015), we truncate answers to 40 tokens.

The task is to rank the candidate answers based on their relatedness to the question. Evaluation measures are mean average precision (MAP) and mean reciprocal rank (MRR).

**Task-Specific Setup.** We use cosine similarity as the similarity score for AS. In addition, we use sentence lengths, *WordCnt* (count of the number of non-stopwords in the question that also occur in the answer) and *WgtWordCnt* (reweight the counts by the IDF values of the question words). Thus, the final input to the LR layer has size  $k + 4$ : one cosine for each of the  $k$  blocks and the four additional features.

We compare with seven **baselines**. The first three are considered by Yang et al. (2015): (i) *WordCnt*; (ii) *WgtWordCnt*; (iii) *CNN-Cnt* (the state-of-the-art system): combine CNN with (i) and (ii). Apart from the baselines considered by Yang et al. (2015), we compare with two *Addition* baselines and two *LSTM* baselines. *Addition* and *A-LSTM* are the *shared baselines* described before. We also combine both with the four extra features; this gives us two additional baselines that we refer to as *Addition(+)* and *A-LSTM(+)*.

<sup>3</sup><http://aka.ms/WikiQA> (Yang et al., 2015)

	method	MAP	MRR
Baselines	WordCnt	0.4891	0.4924
	WgtWordCnt	0.5099	0.5132
	CNN-Cnt	<u>0.6520</u>	<u>0.6652</u>
	Addition	0.5021	0.5069
	Addition(+)	0.5888	0.5929
	A-LSTM	0.5347	0.5483
	A-LSTM(+)	0.6381	0.6537
BCNN	one-conv	0.6629	0.6813
	two-conv	0.6593	0.6738
ABCNN-1	one-conv	0.6810*	0.6979*
	two-conv	0.6855*	0.7023*
ABCNN-2	one-conv	0.6885*	0.7054*
	two-conv	0.6879*	0.7068*
ABCNN-3	one-conv	0.6914*	<b>0.7127*</b>
	two-conv	<b>0.6921*</b>	0.7108*

Table 3: Results on WikiQA. Best result per column is bold. Significant improvements over state-of-the-art baselines (underlined) are marked with \* ( $t$ -test,  $p < .05$ ).

**Results.** Table 3 shows performance of the baselines, of the BCNN and of the three ABCNNs. For CNNs, we test one (one-conv) and two (two-conv) convolution-pooling blocks.

The non-attention network BCNN already performs better than the baselines. If we add attention mechanisms, then the performance further improves by several points. Comparing the ABCNN-2 with the ABCNN-1, we find the ABCNN-2 is slightly better even though the ABCNN-2 is the simpler architecture. If we combine the ABCNN-1 and the ABCNN-2 to form the ABCNN-3, we get further improvement.<sup>4</sup>

This can be explained by the ABCNN-3’s ability to take attention of finer-grained granularity into consideration in each convolution-pooling block while the ABCNN-1 and the ABCNN-2 consider attention only at convolution input or only at pooling input, respectively. We also find that stacking two convolution-pooling blocks does not bring consistent improvement and therefore do not test deeper architectures.

## 5.2 Paraphrase Identification

We use the Microsoft Research Paraphrase (MSRP) corpus (Dolan et al., 2004). The training set contains 2753 true / 1323 false and the test set 1147 true / 578 false paraphrase pairs. We randomly select 400

<sup>4</sup>If we limit the input to the LR layer to the  $k$  similarity scores in the ABCNN-3 (two-conv), results are .660 (MAP) / .677 (MRR).



pairs from train and use them as dev; but we still report results for training on the entire training set. For each triple (label,  $s_0$ ,  $s_1$ ) in the training set, we also add (label,  $s_1$ ,  $s_0$ ) to the training set to make best use of the training data. Systems are evaluated by accuracy and  $F_1$ .

**Task-Specific Setup.** In this task, we add the 15 MT features from (Madnani et al., 2012) and the lengths of the two sentences. In addition, we compute ROUGE-1, ROUGE-2 and ROUGE-SU4 (Lin, 2004), which are scores measuring the match between the two sentences on (i) unigrams, (ii) bigrams and (iii) unigrams and skip-bigrams (maximum skip distance of four), respectively. In this task, we found transforming Euclidean distance into similarity score by  $1/(1 + |x - y|)$  performs better than cosine similarity. Additionally, we use dynamic pooling (Yin and Schütze, 2015a) of the attention matrix  $\mathbf{A}$  in Equation (1) and forward pooled values of all blocks to the classifier. This gives us better performance than only forwarding sentence-level matching features.

We compare our system with representative DL approaches: (i) A-LSTM; (ii) A-LSTM(+): A-LSTM plus handcrafted features; (iii) RAE (Socher et al., 2011), recursive autoencoder; (iv) Bi-CNN-MI (Yin and Schütze, 2015a), a bi-CNN architecture; and (v) MPSSM-CNN (He et al., 2015), the state-of-the-art NN system for PI, and the following four non-DL systems: (vi) Addition; (vii) Addition(+): Addition plus handcrafted features; (viii) MT (Madnani et al., 2012), a system that combines machine translation metrics;<sup>5</sup> (ix) MF-TF-KLD (Ji and Eisenstein, 2013), the state-of-the-art non-NN system.

**Results.** Table 4 shows that the BCNN is slightly worse than the state-of-the-art whereas the ABCNN-1 roughly matches it. The ABCNN-2 is slightly above the state-of-the-art. The ABCNN-3 outperforms the state-of-the-art in accuracy and  $F_1$ .<sup>6</sup> Two convolution layers only bring small improvements over one.

<sup>5</sup>For better comparability of approaches in our experiments, we use a simple SVM classifier, which performs slightly worse than Madnani et al. (2012)’s more complex meta-classifier.

<sup>6</sup>Improvement of .3 (acc) and .1 ( $F_1$ ) over state-of-the-art is not significant. The ABCNN-3 (two-conv) without “linguistic” features (i.e., MT and ROUGE) achieves 75.1/82.7.

	method	acc	$F_1$
Baselines	majority voting	66.5	79.9
	RAE	76.8	83.6
	Bi-CNN-MI	78.4	84.6
	MPSSM-CNN	<u>78.6</u>	<u>84.7</u>
	MT	76.8	83.8
	MF-TF-KLD	<u>78.6</u>	84.6
	Addition	70.8	80.9
	Addition (+)	77.3	84.1
	A-LSTM	69.5	80.1
	A-LSTM (+)	77.1	84.0
BCNN	one-conv	78.1	84.1
	two-conv	78.3	84.3
ABCNN-1	one-conv	78.5	84.5
	two-conv	78.5	84.6
ABCNN-2	one-conv	78.6	84.7
	two-conv	78.8	84.7
ABCNN-3	one-conv	78.8	<b>84.8</b>
	two-conv	<b>78.9</b>	<b>84.8</b>

Table 4: Results for PI on MSRP

### 5.3 Textual Entailment

SemEval 2014 Task 1 (Marelli et al., 2014a) evaluates system predictions of textual entailment (TE) relations on sentence pairs from the SICK dataset (Marelli et al., 2014b). The three classes are entailment, contradiction and neutral. The sizes of SICK train, dev and test sets are 4439, 495 and 4906 pairs, respectively. We call this dataset ORIG.

We also create NONOVER, a copy of ORIG in which *words occurring in both sentences are removed*. A sentence in NONOVER is denoted by the special token <empty> if all words are removed. Table 5 shows three pairs from ORIG and their transformation in NONOVER. We observe that focusing on the non-overlapping parts provides clearer hints for TE than ORIG. In this task, we run two copies of each network, one for ORIG, one for NONOVER; these two networks have a single common LR layer.

Like Lai and Hockenmaier (2014), we train our final system (after fixing hyperparameters) on train and dev (4934 pairs). Eval measure is accuracy.

**Task-Specific Setup.** We found that for this task forwarding two similarity scores from each block (instead of just one) is helpful. We use cosine similarity and Euclidean distance. As we did for paraphrase identification, we add the 15 MT features for each sentence pair for this task as well; our motivation is that entailed sentences resemble paraphrases more than contradictory sentences do.

	ORIG	NONOVER
0	children in red shirts are	children red shirts
	playing in the leaves	playing
	three kids are sitting in the leaves	three kids sitting
1	three boys are jumping in the leaves	boys
	three kids are jumping in the leaves	kids
2	a man is jumping into an empty pool	an empty
	a man is jumping into a full pool	a full

Table 5: SICK data: Converting the original sentences (ORIG) into the NONOVER format

We use the following linguistic features. **Negation** is important for detecting contradiction. Feature NEG is set to 1 if either sentence contains “no”, “not”, “nobody”, “isn’t” and to 0 otherwise. Following Lai and Hockenmaier (2014), we use WordNet (Miller, 1995) to detect **nyms**: synonyms, hypernyms and antonyms in the pairs. But we do this on NONOVER (not on ORIG) to focus on what is critical for TE. Specifically, feature SYN is the number of word pairs in  $s_0$  and  $s_1$  that are synonyms. HYP0 (resp. HYP1) is the number of words in  $s_0$  (resp.  $s_1$ ) that have a hypernym in  $s_1$  (resp.  $s_0$ ). In addition, we collect all *potential antonym pairs* (PAP) in NONOVER. We identify the matched chunks that occur in *contradictory* and *neutral*, but not in *entailed* pairs. We exclude synonyms and hypernyms and apply a frequency filter of  $n = 2$ . In contrast to Lai and Hockenmaier (2014), we constrain the PAP pairs to cosine similarity above 0.4 in word2vec embedding space as this discards many noise pairs. Feature ANT is the number of matched PAP antonyms in a sentence pair. As before we use sentence **lengths**, both for ORIG (LEN00: length  $s_0$ , LEN10: length  $s_1$ ) and for NONOVER (LEN0N: length  $s_0$ , LEN1N: length  $s_1$ ).

On the whole, we have 24 extra features: 15 MT metrics, NEG, SYN, HYP0, HYP1, ANT, LEN00, LEN10, LEN0N and LEN1N.

Apart from the Addition and LSTM baselines, we further compare with the top-3 systems in SemEval and TrRNTN (Bowman et al., 2015b), a recursive neural network developed for this SICK task.

**Results.** Table 6 shows that our CNNs outperform A-LSTM (with or without linguistic features added) and the top three SemEval systems. Comparing ABCNNs with the BCNN, attention mechanisms consistently improve performance. The ABCNN-1 has performance comparable to the ABCNN-2 while

	method	acc
Sem-Eval Top3	(Jimenez et al., 2014)	83.1
	(Zhao et al., 2014)	83.6
	(Lai and Hockenmaier, 2014)	84.6
TrRNTN	(Bowman et al., 2015b)	76.9
Addition	no features	73.1
	plus features	79.4
A-LSTM	no features	78.0
	plus features	81.7
BCNN	one-conv	84.8
	two-conv	85.0
ABCNN-1	one-conv	85.6
	two-conv	85.8
ABCNN-2	one-conv	85.7
	two-conv	85.8
ABCNN-3	one-conv	86.0*
	two-conv	86.2*

Table 6: Results on SICK. Significant improvements over (Lai and Hockenmaier, 2014) are marked with \* (test of equal proportions,  $p < .05$ ).

the ABCNN-3 is better still: a boost of 1.6 points compared to the previous state of the art.<sup>7</sup>

**Visual Analysis.** Figure 4 visualizes the attention matrices for one TE sentence pair in the ABCNN-2 for blocks  $b_1$  (unigrams),  $b_2$  (first convolutional layer) and  $b_3$  (second convolutional layer). Darker shades of blue indicate stronger attention values.

In Figure 4 (top), each word corresponds to exactly one row or column. We can see that words in  $s_i$  with semantic equivalents in  $s_{1-i}$  get high attention while words without semantic equivalents get low attention, e.g., “walking” and “murals” in  $s_0$  and “front” and “colorful” in  $s_1$ . This behavior seems reasonable for the unigram level.

Rows/columns of the attention matrix in Figure 4 (middle) correspond to phrases of length three since filter width  $w = 3$ . High attention values generally correlate with close semantic correspondence: the phrase “people are” in  $s_0$  matches “several people are” in  $s_1$ ; both “are walking outside” and “walking outside the” in  $s_0$  match “are in front” in  $s_1$ ; “the building that” in  $s_0$  matches “a colorful building” in  $s_1$ . More interestingly, looking at the bottom right corner, both “on it” and “it” in  $s_0$  match “building” in  $s_1$ ; this indicates that ABCNNs are able to detect some coreference across sentences. “building” in  $s_1$  has two places in which higher attentions appear, one is with “it” in  $s_0$ , the other is with “the building

<sup>7</sup>If we run the ABCNN-3 (two-conv) without the 24 linguistic features, performance is 84.6.

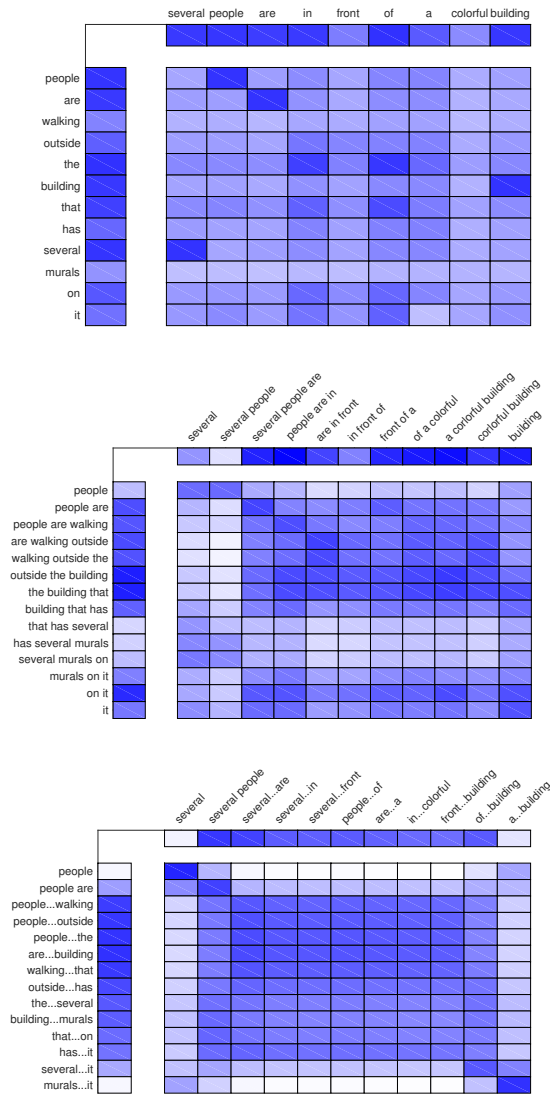


Figure 4: Attention visualization for TE. Top: unigrams,  $b_1$ . Middle: conv1,  $b_2$ . Bottom: conv2,  $b_3$ .

that” in  $s_0$ . This may indicate that ABCNNs recognize that “building” in  $s_1$  and “the building that” / “it” in  $s_0$  refer to the same object. Hence, coreference resolution across sentences as well as within a sentence both are detected. For the attention vectors on the left and the top, we can see that attention has focused on the key parts: “people are walking outside the building that” in  $s_0$ , “several people are in” and “of a colorful building” in  $s_1$ .

Rows/columns of the attention matrix in Figure 4 (bottom, second layer of convolution) correspond to phrases of length 5 since filter width  $w = 3$  in both convolution layers ( $5 = 1 + 2 * (3 - 1)$ ). We use “...” to denote words in the middle if a phrase

like “several...front” has more than two words. We can see that attention distribution in the matrix has focused on some local regions. As granularity of phrases is larger, it makes sense that the attention values are smoother. But we still can find some interesting clues: at the two ends of the main diagonal, higher attentions hint that the first part of  $s_0$  matches well with the first part of  $s_1$ ; “several murals on it” in  $s_0$  matches well with “of a colorful building” in  $s_1$ , which satisfies the intuition that these two phrases are crucial for making a decision on TE in this case. This again shows the potential strength of our system in figuring out which parts of the two sentences refer to the same object. In addition, in the central part of the matrix, we can see that the long phrase “people are walking outside the building” in  $s_0$  matches well with the long phrase “are in front of a colorful building” in  $s_1$ .

## 6 Summary

We presented three mechanisms to integrate attention into CNNs for general sentence pair modeling tasks.

Our experiments on AS, PI and TE show that attention-based CNNs perform better than CNNs without attention mechanisms. The ABCNN-2 generally outperforms the ABCNN-1 and the ABCNN-3 surpasses both.

In all tasks, we did not find any big improvement of two layers of convolution over one layer. This is probably due to the limited size of training data. We expect that, as larger training sets become available, deep ABCNNs will show even better performance.

In addition, linguistic features contribute in all three tasks: improvements by 0.0321 (MAP) and 0.0338 (MRR) for AS, improvements by 3.8 (acc) and 2.1 ( $F_1$ ) for PI and an improvement by 1.6 (acc) for TE. But our ABCNNs can still reach or surpass state-of-the-art even without those features in AS and TE tasks. This indicates that ABCNNs are generally strong NN systems.

Attention-based LSTMs are especially successful in tasks with a strong *generation component* like machine translation (discussed in Sec. 2). CNNs have not been used for this type of task. This is an interesting area of future work for attention-based CNNs.

## Acknowledgments

We gratefully acknowledge the support of Deutsche Forschungsgemeinschaft (DFG): grant SCHU 2246/8-2.

We would like to thank the anonymous reviewers for their helpful comments.

## References

- Jimmy Ba, Volodymyr Mnih, and Koray Kavukcuoglu. 2015. Multiple object recognition with visual attention. In *Proceedings of ICLR*.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *Proceedings of ICLR*.
- Matthew W. Bilotti, Paul Ogilvie, Jamie Callan, and Eric Nyberg. 2007. Structured retrieval for question answering. In *Proceedings of SIGIR*, pages 351–358.
- William Blacoe and Mirella Lapata. 2012. A comparison of vector-based representations for semantic composition. In *Proceedings of EMNLP-CoNLL*, pages 546–556.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015a. A large annotated corpus for learning natural language inference. In *Proceedings of EMNLP*, pages 632–642.
- Samuel R. Bowman, Christopher Potts, and Christopher D. Manning. 2015b. Recursive neural networks can learn logical semantics. In *Proceedings of CVSC Workshop*, pages 12–21.
- Jane Bromley, James W. Bentz, Léon Bottou, Isabelle Guyon, Yann LeCun, Cliff Moore, Eduard Säckinger, and Roopak Shah. 1993. Signature verification using a “siamese” time delay neural network. *IJPRAI*, 7(4):669–688.
- Chunshui Cao, Xianming Liu, Yi Yang, Yinan Yu, Jiang Wang, Zilei Wang, Yongzhen Huang, Liang Wang, Chang Huang, Wei Xu, Deva Ramanan, and Thomas S. Huang. 2015. Look and think twice: Capturing top-down visual attention with feedback convolutional neural networks. In *Proceedings of ICCV*, pages 2956–2964.
- Ming-Wei Chang, Dan Goldwasser, Dan Roth, and Vivek Srikumar. 2010. Discriminative learning over constrained latent representations. In *Proceedings of NAACL-HLT*, pages 429–437.
- Kan Chen, Jiang Wang, Liang-Chieh Chen, Haoyuan Gao, Wei Xu, and Ram Nevatia. 2015. ABC-CNN: An attention based convolutional neural network for visual question answering. *CoRR*, abs/1511.05960.
- Jan Chorowski, Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. End-to-end continuous speech recognition using attention-based recurrent NN: First results. In *Proceedings of Deep Learning and Representation Learning Workshop, NIPS*.
- Jan Chorowski, Dzmitry Bahdanau, Dmitriy Serdyuk, Kyunghyun Cho, and Yoshua Bengio. 2015. Attention-based models for speech recognition. In *Proceedings of NIPS*, pages 577–585.
- Bill Dolan, Chris Quirk, and Chris Brockett. 2004. Un-supervised construction of large paraphrase corpora: Exploiting massively parallel news sources. In *Proceedings of COLING*, pages 350–356.
- John Duchi, Elad Hazan, and Yoram Singer. 2011. Adaptive subgradient methods for online learning and stochastic optimization. *JMLR*, 12:2121–2159.
- Minwei Feng, Bing Xiang, Michael R. Glass, Lidan Wang, and Bowen Zhou. 2015. Applying deep learning to answer selection: A study and an open task. In *Proceedings of IEEE ASRU Workshop*, pages 813–820.
- Karol Gregor, Ivo Danihelka, Alex Graves, Danilo Jimenez Rezende, and Daan Wierstra. 2015. DRAW: A recurrent neural network for image generation. In *Proceedings of ICML*, pages 1462–1471.
- Hua He, Kevin Gimpel, and Jimmy Lin. 2015. Multi-perspective sentence similarity modeling with convolutional neural networks. In *Proceedings of EMNLP*, pages 1576–1586.
- Michael Heilman and Noah A. Smith. 2010. Tree edit models for recognizing textual entailments, paraphrases, and answers to questions. In *Proceedings of NAACL-HLT*, pages 1011–1019.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural Computation*, 9(8):1735–1780.
- Seunghoon Hong, Junhyuk Oh, Honglak Lee, and Bo-hyung Han. 2016. Learning transferrable knowledge for semantic segmentation with deep convolutional neural network. In *Proceedings of CVPR*.
- Baotian Hu, Zhengdong Lu, Hang Li, and Qingcai Chen. 2014. Convolutional neural network architectures for matching natural language sentences. In *Proceedings of NIPS*, pages 2042–2050.
- Yangfeng Ji and Jacob Eisenstein. 2013. Discriminative improvements to distributional sentence similarity. In *Proceedings of EMNLP*, pages 891–896.
- Sergio Jimenez, George Dueñas, Julia Baquero, and Alexander Gelbukh. 2014. UNAL-NLP: Combining soft cardinality features for semantic textual similarity, relatedness and entailment. In *Proceedings of SemEval*, pages 732–742.

- Nal Kalchbrenner, Edward Grefenstette, and Phil Blunsom. 2014. A convolutional neural network for modelling sentences. In *Proceedings of ACL*, pages 655–665.
- Yoon Kim. 2014. Convolutional neural networks for sentence classification. In *Proceedings of EMNLP*, pages 1746–1751.
- Alice Lai and Julia Hockenmaier. 2014. Illinois-LH: A denotational and distributional approach to semantics. In *Proceedings of SemEval*, pages 329–334.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. 1998. Gradient-based learning applied to document recognition. In *Proceedings of the IEEE*, pages 2278–2324.
- Jiwei Li, Minh-Thang Luong, and Dan Jurafsky. 2015. A hierarchical neural autoencoder for paragraphs and documents. In *Proceedings of ACL*, pages 1106–1115.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Proceedings of the ACL Text Summarization Workshop*.
- Thang Luong, Hieu Pham, and Christopher D. Manning. 2015. Effective approaches to attention-based neural machine translation. In *Proceedings of EMNLP*, pages 1412–1421.
- Nitin Madnani, Joel Tetreault, and Martin Chodorow. 2012. Re-examining machine translation metrics for paraphrase identification. In *Proceedings of NAACL-HLT*, pages 182–190.
- Marco Marelli, Luisa Bentivogli, Marco Baroni, Raffaella Bernardi, Stefano Menini, and Roberto Zamparelli. 2014a. Semeval-2014 task 1: Evaluation of compositional distributional semantic models on full sentences through semantic relatedness and textual entailment. In *Proceedings of SemEval*, pages 1–8.
- Marco Marelli, Stefano Menini, Marco Baroni, Luisa Bentivogli, Raffaella Bernardi, and Roberto Zamparelli. 2014b. A SICK cure for the evaluation of compositional distributional semantic models. In *Proceedings of LREC*, pages 216–223.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Proceedings of NIPS*, pages 3111–3119.
- George A. Miller. 1995. WordNet: A lexical database for english. *Commun. ACM*, 38(11):39–41.
- Volodymyr Mnih, Nicolas Heess, Alex Graves, and Koray Kavukcuoglu. 2014. Recurrent models of visual attention. In *Proceedings of NIPS*, pages 2204–2212.
- Dan Moldovan, Christine Clark, Sanda Harabagiu, and Daniel Hodges. 2007. Cogex: A semantically and contextually enriched logic prover for question answering. *Journal of Applied Logic*, 5(1):49–69.
- Vasin Punyakanok, Dan Roth, and Wen-tau Yih. 2004. Mapping dependencies trees: An application to question answering. In *Proceedings of AI&Math 2004 (Special session: Intelligent Text Processing)*.
- Tim Rocktäschel, Edward Grefenstette, Karl Moritz Hermann, Tomáš Kočiský, and Phil Blunsom. 2016. Reasoning about entailment with neural attention. In *Proceedings of ICLR*.
- Alexander M. Rush, Sumit Chopra, and Jason Weston. 2015. A neural attention model for abstractive sentence summarization. In *Proceedings of EMNLP*, pages 379–389.
- Dan Shen and Mirella Lapata. 2007. Using semantic roles to improve question answering. In *Proceedings of EMNLP-CoNLL*, pages 12–21.
- Richard Socher, Eric H. Huang, Jeffrey Pennington, Andrew Y. Ng, and Christopher D. Manning. 2011. Dynamic pooling and unfolding recursive autoencoders for paraphrase detection. In *Proceedings of NIPS*, pages 801–809.
- Ming Tan, Bing Xiang, and Bowen Zhou. 2016. LSTM-based deep learning models for non-factoid answer selection. In *Proceedings of ICLR Workshop*.
- Shengxian Wan, Yanyan Lan, Jiafeng Guo, Jun Xu, Liang Pang, and Xueqi Cheng. 2016. A deep architecture for semantic matching with multiple positional sentence representations. In *Proceedings of AAAI*, pages 2835–2841.
- Mengqiu Wang, Noah A. Smith, and Teruko Mitamura. 2007. What is the jeopardy model? A quasi-synchronous grammar for QA. In *Proceedings of EMNLP-CoNLL*, pages 22–32.
- Tianjun Xiao, Yichong Xu, Kuiyuan Yang, Jiaxing Zhang, Yuxin Peng, and Zheng Zhang. 2015. The application of two-level attention models in deep convolutional neural network for fine-grained image classification. In *Proceedings of CVPR*, pages 842–850.
- Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron C. Courville, Ruslan Salakhutdinov, Richard S. Zemel, and Yoshua Bengio. 2015. Show, attend and tell: Neural image caption generation with visual attention. In *Proceedings of ICML*, pages 2048–2057.
- Yi Yang, Wen-tau Yih, and Christopher Meek. 2015. WikiQA: A challenge dataset for open-domain question answering. In *Proceedings of EMNLP*, pages 2013–2018.
- Xuchen Yao, Benjamin Van Durme, Chris Callison-Burch, and Peter Clark. 2013a. Semi-markov phrase-based monolingual alignment. In *Proceedings of EMNLP*, pages 590–600.
- Xuchen Yao, Benjamin Van Durme, and Peter Clark. 2013b. Automatic coupling of answer extraction and information retrieval. In *Proceedings of ACL*, pages 159–165.

- Wen-tau Yih, Ming-Wei Chang, Christopher Meek, and Andrzej Pastusiak. 2013. Question answering using enhanced lexical semantic models. In *Proceedings of ACL*, pages 1744–1753.
- Wenpeng Yin and Hinrich Schütze. 2015a. Convolutional neural network for paraphrase identification. In *Proceedings of NAACL-HLT*, pages 901–911.
- Wenpeng Yin and Hinrich Schütze. 2015b. Multi-GranCNN: An architecture for general matching of text chunks on multiple levels of granularity. In *Proceedings of ACL-IJCNLP*, pages 63–73.
- Lei Yu, Karl Moritz Hermann, Phil Blunsom, and Stephen Pulman. 2014. Deep learning for answer sentence selection. In *Proceedings of Deep Learning and Representation Learning Workshop, NIPS*.
- Jiang Zhao, Tiantian Zhu, and Man Lan. 2014. ECNU: One stone two birds: Ensemble of heterogenous measures for semantic relatedness and textual entailment. In *Proceedings of SemEval*, pages 271–277.