

Semantic Sentence Matching with Densely-connected Recurrent and Co-attentive Information

Seonhoon Kim^{1,2}, Inho Kang¹, Nojun Kwak²

¹Naver Search, ²Seoul National University

{seonhoon.kim|once.ihkang}@navercorp.com, nojunk@snu.ac.kr

Abstract

Sentence matching is widely used in various natural language tasks such as natural language inference, paraphrase identification, and question answering. For these tasks, understanding logical and semantic relationship between two sentences is required but it is yet challenging. Although attention mechanism is useful to capture the semantic relationship and to properly align the elements of two sentences, previous methods of attention mechanism simply use a summation operation which does not retain original features enough. Inspired by DenseNet, a densely connected convolutional network, we propose a densely-connected co-attentive recurrent neural network, each layer of which uses concatenated information of attentive features as well as hidden features of all the preceding recurrent layers. It enables preserving the original and the co-attentive feature information from the bottommost word embedding layer to the uppermost recurrent layer. To alleviate the problem of an ever-increasing size of feature vectors due to dense concatenation operations, we also propose to use an autoencoder after dense concatenation. We evaluate our proposed architecture on highly competitive benchmark datasets related to sentence matching. Experimental results show that our architecture, which retains recurrent and attentive features, achieves state-of-the-art performances for most of the tasks.

Introduction

Semantic sentence matching, a fundamental technology in natural language processing, requires lexical and compositional semantics. In paraphrase identification, sentence matching is utilized to identify whether two sentences have identical meaning or not. In natural language inference also known as recognizing textual entailment, it determines whether a hypothesis sentence can reasonably be inferred from a given premise sentence. In question answering, sentence matching is required to determine the degree of matching 1) between a query and a question for question retrieval, and 2) between a question and an answer for answer selection. However identifying logical and semantic relationship between two sentences is not trivial due to the problem of the semantic gap (Liu et al. 2016).

Recent advances of deep neural network enable to learn textual semantics for sentence matching. Large amount of annotated data such as Quora (Csernai 2017), SNLI (Bowman

et al. 2015), and MultiNLI (Williams, Nangia, and Bowman 2017) have contributed significantly to learning semantics as well. In the conventional methods, a matching model can be trained in two different ways (Gong, Luo, and Zhang 2018). The first methods are sentence-encoding-based ones where each sentence is encoded to a fixed-sized vector in a complete isolated manner and the two vectors for the corresponding sentences are used in predicting the degree of matching. The others are joint methods that allow to utilize interactive features like attentive information between the sentences.

In the former paradigm, because two sentences have no interaction, they can not utilize interactive information during the encoding procedure. In our work, we adopted a joint method which enables capturing interactive information for performance improvements. Furthermore, we employ a substantially deeper recurrent network for sentence matching like deep neural machine translator (NMT) (Wu et al. 2016). Deep recurrent models are more advantageous for learning long sequences and outperform the shallower architectures. However, the attention mechanism is unstable in deeper models with the well-known vanishing gradient problem. Though GNMT (Wu et al. 2016) uses residual connection between recurrent layers to allow better information and gradient flow, there are some limitations. The recurrent hidden or attentive features are not preserved intact through residual connection because the summation operation may impede the information flow in deep networks.

Inspired by Densenet (Huang et al. 2017), we propose a densely-connected recurrent network where the recurrent hidden features are retained to the uppermost layer. In addition, instead of the conventional summation operation, the concatenation operation is used in combination with the attention mechanism to preserve co-attentive information better. The proposed architecture shown in Figure 1 is called DRCN which is an abbreviation for *Densely-connected Recurrent and Co-attentive neural Network*. The proposed DRCN can utilize the increased representational power of deeper recurrent networks and attentive information. Furthermore, to alleviate the problem of an ever-increasing feature vector size due to concatenation operations, we adopted an autoencoder and forwarded a fixed length vector to the higher layer recurrent module as shown in the figure. DRCN is, to our best knowledge, the first generalized version of DenseRNN which is expandable to deeper layers with the property of

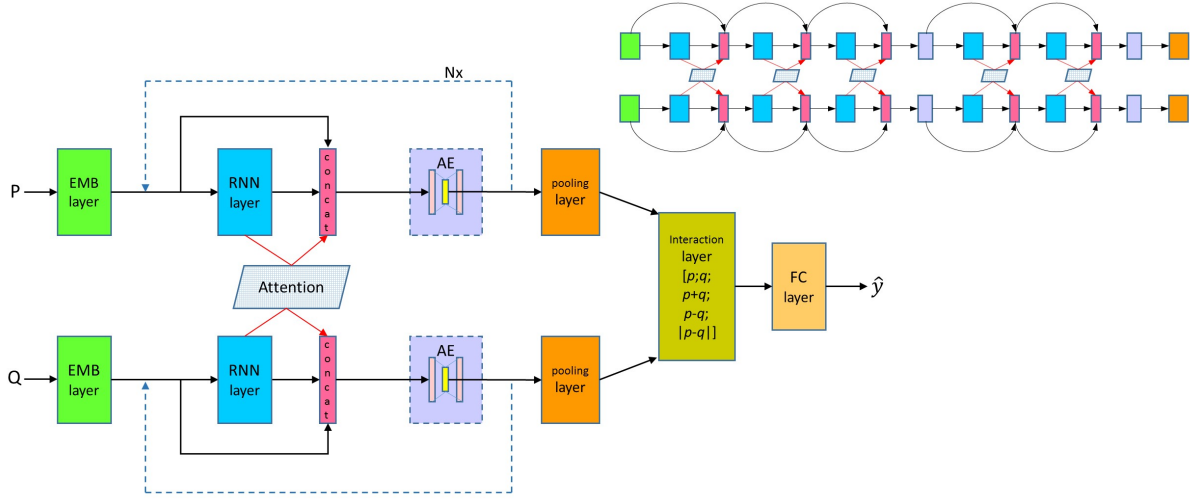


Figure 1: General architecture of our Densely-connected Recurrent and Co-attentive neural Network (DRCN). Dashed arrows indicate that a group of RNN-layer, concatenation and AE can be repeated multiple (N) times (like a repeat mark in a music score). The bottleneck component denoted as AE, inserted to prevent the ever-growing size of a feature vector, is optional for each repetition. The upper right diagram is our specific architecture for experiments with 5 RNN layers ($N = 4$).

controllable feature sizes by the use of an autoencoder.

We evaluate our model on three sentence matching tasks: *natural language inference*, *paraphrase identification* and *answer sentence selection*. Experimental results on five highly competitive benchmark datasets (SNLI, MultiNLI, QUORA, TrecQA and SelQA) show that our model significantly outperforms the current state-of-the-art results on most of the tasks.

Related Work

Earlier approaches of sentence matching mainly relied on conventional methods such as **syntactic features, transformations or relation extraction** (Romano et al. 2006; Wang, Smith, and Mitamura 2007). These are restrictive in that they work only on very specific tasks.

The developments of large-scale annotated datasets (Bowman et al. 2015; Williams, Nangia, and Bowman 2017) and deep learning algorithms have led a big progress on matching natural language sentences. Furthermore, the well-established attention mechanisms endowed richer information for sentence matching by providing alignment and dependency relationship between two sentences. The release of the large-scale datasets also has encouraged the developments of the learning-centered approaches to semantic representation. The first type of these approaches is sentence-encoding-based methods (Conneau et al. 2017; Choi, Yoo, and goo Lee 2017; Nie and Bansal 2017; Shen et al. 2018) where sentences are encoded into their own sentence representation without any cross-interaction. Then, a classifier such as a neural network is applied to decide the relationship based on these independent sentence representations. These sentence-encoding-based methods are simple to extract sentence representation and are able to be used for transfer learning to other natural language tasks (Conneau

et al. 2017). On the other hand, the joint methods, which make up for the lack of interaction in the former methods, use cross-features as an attention mechanism to express the word- or phrase-level alignments for performance improvements (Wang, Hamza, and Florian 2017; Chen et al. 2017b; Gong, Luo, and Zhang 2018; Yang et al. 2016).

Recently, the architectural developments using deeper layers have led more progress in performance. The residual connection is widely and commonly used to increase the depth of a network stably (He et al. 2016; Wu et al. 2016). More recently, Huang *et al.* (Huang et al. 2017) enable the features to be connected from lower to upper layers using the concatenation operation without any loss of information on lower-layer features.

External resources are also used for sentence matching. Chen *et al.* (Chen et al. 2017a; Chen et al. 2017b) used syntactic parse trees or lexical databases like WordNet to measure the semantic relationship among the words and Pavlick *et al.* (Pavlick et al. 2015) added interpretable semantics to the paraphrase database.

Unlike these, in this paper, we do not use any such external resources. Our work belongs to the joint approaches which uses densely-connected recurrent and co-attentive information to enhance representation power for semantic sentence matching.

Methods

In this section, we describe our sentence matching architecture DRCN which is composed of the following three components: (1) word representation layer, (2) attentively connected RNN and (3) interaction and prediction layer. We denote two input sentences as $P = \{p_1, p_2, \dots, p_I\}$ and $Q = \{q_1, q_2, \dots, q_J\}$ where p_i/q_j is the i^{th}/j^{th} word of the sentence P/Q and I/J is the word length of P/Q . The overall

architecture of the proposed DRCN is shown in Fig. 1.

Word Representation Layer

To construct the **word representation layer**, we concatenate **word embedding, character representation and the exact matched flag** which was used in (Gong, Luo, and Zhang 2018).

In word embedding, **each word is represented as a d -dimensional vector** by using a pre-trained word embedding method such as GloVe (Pennington, Socher, and Manning 2014) or Word2vec (Mikolov et al. 2013). In our model, a word embedding vector can be updated or fixed during training. The strategy whether to make the pre-trained word embedding be trainable or not is heavily task-dependent. Trainable word embeddings capture the characteristics of the training data well but can result in overfitting. On the other hand, **fixed (non-trainable) word embeddings lack flexibility** on task-specific data, while it can be robust for overfitting, especially for less frequent words. We use both the **trainable embedding $e_{p_i}^{tr}$ and the fixed (non-trainable) embedding $e_{p_i}^{fix}$** to let them play complementary roles in enhancing the performance of our model. This technique of mixing trainable and non-trainable word embeddings is simple but yet effective.

The character representation c_{p_i} is calculated by feeding randomly initialized character embeddings into a convolutional neural network with the max-pooling operation. **The character embeddings and convolutional weights are jointly learned during training.**

Like (Gong, Luo, and Zhang 2018), the exact match flag f_{p_i} is activated if the same word is found in the other sentence.

Our final word representational feature p_i^w for the word p_i is composed of four components as follows:

$$\begin{aligned} e_{p_i}^{tr} &= E^{tr}(p_i), \quad e_{p_i}^{fix} = E^{fix}(p_i) \\ c_{p_i} &= \text{Char-Conv}(p_i) \\ p_i^w &= [e_{p_i}^{tr}; e_{p_i}^{fix}; c_{p_i}; f_{p_i}]. \end{aligned} \quad (1)$$

Here, E^{tr} and E^{fix} are the trainable and non-trainable (fixed) word embeddings respectively. Char-Conv is the character-level convolutional operation and $[\cdot; \cdot]$ is the concatenation operator. For each word in both sentences, the same above procedure is used to extract word features.

Densely connected Recurrent Networks

The **ordinal** stacked RNNs (Recurrent Neural Networks) are composed of multiple RNN layers on top of each other, with the output sequence of previous layer forming the input sequence for the next. More concretely, let H_l be the l^{th} RNN layer in a stacked RNN. Note that in our implementation, we employ the bidirectional LSTM (BiLSTM) as a base block of H_l . At the time step t , an ordinal stacked RNN is expressed as follows:

$$\begin{aligned} h_t^l &= H_l(x_t^l, h_{t-1}^l) \\ x_t^l &= h_t^{l-1}. \end{aligned} \quad (2)$$

While this architecture enables us to build up higher level representation, deeper networks have difficulties in training due to the exploding or vanishing gradient problem.

To encourage gradient to flow in the backward pass, **residual connection** (He et al. 2016) is introduced which bypasses the non-linear transformations with an identity mapping. Incorporating this into (2), it becomes

$$\begin{aligned} h_t^l &= H_l(x_t^l, h_{t-1}^l) \\ x_t^l &= h_t^{l-1} + x_t^{l-1}. \end{aligned} \quad (3)$$

However, the summation operation in the residual connection may impede the information flow in the network (Huang et al. 2017). Motivated by Densenet (Huang et al. 2017), we employ direct connections using the concatenation operation from any layer to all the subsequent layers so that the features of previous layers are not to be modified but to be retained as they are as depicted in Figure 1. The densely connected recurrent neural networks can be described as

$$\begin{aligned} h_t^l &= H_l(x_t^l, h_{t-1}^l) \\ x_t^l &= [h_t^{l-1}; x_t^{l-1}]. \end{aligned} \quad (4)$$

The concatenation operation enables the hidden features to be preserved until they reach to the uppermost layer and all the previous features work for prediction as collective knowledge (Huang et al. 2017).

Densely-connected Co-attentive networks

Attention mechanism, which has largely succeeded in many domains (Wu et al. 2016; Vaswani et al. 2017), is a technique to learn effectively where a context vector is matched conditioned on a specific sequence.

Given two sentences, a context vector is calculated based on an attention mechanism focusing on the relevant part of the two sentences at each RNN layer. The calculated attentive information represents soft-alignment between two sentences. In this work, we also use an attention mechanism. We incorporate co-attentive information into densely connected recurrent features using the concatenation operation, so as not to lose any information (Fig. 1). This concatenated recurrent and co-attentive features which are obtained by densely connecting the features from the undermost to the uppermost layers, enrich the collective knowledge for lexical and compositional semantics.

The attentive information a_{p_i} of the i^{th} word $p_i \in P$ against the sentence Q is calculated as a weighted sum of h_{q_j} 's which are weighted by the softmax weights as follows :

$$\begin{aligned} a_{p_i} &= \sum_{j=1}^J \alpha_{i,j} h_{q_j} \\ \alpha_{i,j} &= \frac{\exp(e_{i,j})}{\sum_{k=1}^J \exp(e_{i,k})} \\ e_{i,j} &= \cos(h_{p_i}, h_{q_j}) \end{aligned} \quad (5)$$

Similar to the densely connected RNN hidden features, we concatenate the attentive context vector a_{p_i} with triggered vector h_{p_i} so as to retain attentive information as an input to the next layer:

$$\begin{aligned} h_t^l &= H_l(x_t^l, h_{t-1}^l) \\ x_t^l &= [h_t^{l-1}; a_t^{l-1}; x_t^{l-1}]. \end{aligned} \quad (6)$$

Bottleneck component

Our network uses all layers’ outputs as a community of semantic knowledge. However, this network is a structure with increasing input features as layers get deeper, and has a large number of parameters especially in the fully-connected layer. To address this issue, we employ an autoencoder as a bottleneck component. Autoencoder is a compression technique that reduces the number of features while retaining the original information, which can be used as a distilled semantic knowledge in our model. Furthermore, this component increased the test performance by working as a regularizer in our experiments.

Interaction and Prediction Layer

To extract a proper representation for each sentence, we apply the step-wise max-pooling operation over densely connected recurrent and co-attentive features (pooling in Fig. 1). More specifically, if the output of the final RNN layer is a 100d vector for a sentence with 30 words, a 30×100 matrix is obtained which is max-pooled column-wise such that the size of the resultant vector p or q is 100. Then, we aggregate these representations p and q for the two sentences P and Q in various ways in the interaction layer and the final feature vector v for semantic sentence matching is obtained as follows:

$$v = [p; q; p + q; p - q; |p - q|]. \quad (7)$$

Here, the operations $+$, $-$ and $|\cdot|$ are performed element-wise to infer the relationship between two sentences. The element-wise subtraction $p - q$ is an **asymmetric** operator for one-way type tasks such as *natural language inference* or *answer sentence selection*.

Finally, based on previously aggregated features v , we use two fully-connected layers with ReLU activation followed by one fully-connected output layer. Then, the softmax function is applied to obtain a probability distribution of each class. The model is trained end-to-end by minimizing the multi-class cross entropy loss and the reconstruction loss of autoencoders.

Experiments

We evaluate our matching model on five popular and well-studied benchmark datasets for three challenging sentence matching tasks: (i) SNLI and MultiNLI for natural language inference; (ii) Quora Question Pair for paraphrase identification; and (iii) TrecQA and SelQA for answer sentence selection in question answering. Additional details about the above datasets can be found in the supplementary materials.

Implementation Details

We initialized word embedding with 300d GloVe vectors pre-trained from the 840B Common Crawl corpus (Pennington, Socher, and Manning 2014), while the word embeddings for the out-of-vocabulary words were initialized randomly. We also randomly initialized character embedding with a 16d vector and extracted 32d character representation with a convolutional network. For the densely-connected recurrent layers, we stacked 5 layers each of which have 100 hidden

Premise	<i>two bicyclists in spandex and helmets in a race pedaling uphill.</i>
Hypothesis	<i>A pair of humans are riding their bicycle with tight clothing, competing with each other.</i>
Label	{ <i>entailment</i> ; <i>neutral</i> ; <i>contradiction</i> }

Premise	<i>Several men in front of a white building.</i>
Hypothesis	<i>Several people in front of a gray building.</i>
Label	{ <i>entailment</i> ; <i>neutral</i> ; <i>contradiction</i> }

Table 1: Examples of *natural language inference*.

units. We set 1000 hidden units with respect to the fully-connected layers. The dropout was applied after the word and character embedding layers with a keep rate of 0.5. It was also applied before the fully-connected layers with a keep rate of 0.8. For the bottleneck component, we set 200 hidden units as encoded features of the autoencoder with a dropout rate of 0.2. The batch normalization was applied on the fully-connected layers, only for the one-way type datasets. The RMSProp optimizer with an initial learning rate of 0.001 was applied. The learning rate was decreased by a factor of 0.85 when the dev accuracy does not improve. All weights except embedding matrices are constrained by L2 regularization with a regularization constant $\lambda = 10^{-6}$. The sequence lengths of the sentence are all different for each dataset: 35 for SNLI, 55 for MultiNLI, 25 for Quora question pair and 50 for TrecQA. The learning parameters were selected based on the best performance on the dev set. We employed 8 different randomly initialized sets of parameters with the same model for our ensemble approach.

Experimental Results

SNLI and MultiNLI We evaluated our model on the natural language inference task over SNLI and MultiNLI datasets. Table 2 shows the results on SNLI dataset of our model with other published models. Among them, ESIM+ELMo and LM-Transformer are the current state-of-the-art models. However, they use additional contextualized word representations from language models as an external knowledge. The proposed DRCN obtains an accuracy of 88.9% which is a competitive score although we do not use any external knowledge like ESIM+ELMo and LM-Transformer. The ensemble model achieves an accuracy of 90.1%, which sets the new state-of-the-art performance. Our ensemble model with 53m parameters ($6.7m \times 8$) outperforms the LM-Transformer whose the number of parameters is 85m. Furthermore, in case of the encoding-based method, we obtain the best performance of 86.5% without the co-attention and exact match flag.

Table 3 shows the results on MATCHED and MISMATCHED problems of MultiNLI dataset. Our plain DRCN has a competitive performance without any contextualized knowledge. And, by combining DRCN with the ELMo, one of the contextualized embeddings from language models, our model outperforms the LM-Transformer which has 85m parameters with fewer parameters of 61m. From this point of view, the combination of our model with a contextualized knowledge

Models	Acc.	$ \theta $
<i>Sentence encoding-based method</i>		
BiLSTM-Max (Conneau et al. 2017)	84.5	40m
Gumbel TreeLSTM (Choi, Yoo, and goo Lee 2017)	85.6	2.9m
CAFE (Tay, Tuan, and Hui 2017)	85.9	3.7m
Gumbel TreeLSTM (Choi, Yoo, and goo Lee 2017)	86.0	10m
Residual stacked (Nie and Bansal 2017)	86.0	29m
Reinforced SAN (Shen et al. 2018)	86.3	3.1m
Distance SAN (Im and Cho 2017)	86.3	3.1m
DRCN (- Attn, - Flag)	86.5	5.6m
<i>Joint method (cross-features available)</i>		
DIIN (Gong, Luo, and Zhang 2018)	88.0 / 88.9	4.4m
ESIM (Chen et al. 2017b)	88.0 / 88.6	4.3m
BCN+CoVe+Char (McCann et al. 2017)	88.1 / -	22m
DR-BiLSTM (Ghaeini et al. 2018)	88.5 / 89.3	7.5m
CAFE (Tay, Tuan, and Hui 2017)	88.5 / 89.3	4.7m
KIM (Chen et al. 2017a)	88.6 / 89.1	4.3m
ESIM+ELMo (Peters et al. 2018)	88.7 / 89.3	8.0m
LM-Transformer (Radford et al. 2018)	89.9 / -	85m
DRCN (- AE)	88.7 / -	20m
DRCN	88.9 / 90.1	6.7m

Table 2: Classification accuracy (%) for natural language inference on SNLI test set. $|\theta|$ denotes the number of parameters in each model.

Models	Accuracy (%)	
	MATCHED	MISMATCHED
ESIM (Williams, Nangia, and Bowman 2017)	72.3	72.1
DIIN (Gong, Luo, and Zhang 2018)	78.8	77.8
CAFE (Tay, Tuan, and Hui 2017)	78.7	77.9
LM-Transformer (Radford et al. 2018)	82.1	81.4
DRCN	79.1	78.4
DIIN* (Gong, Luo, and Zhang 2018)	80.0	78.7
CAFE* (Tay, Tuan, and Hui 2017)	80.2	79.0
DRCN*	80.6	79.5
DRCN+ELMo*	82.3	81.4

Table 3: Classification accuracy for natural language inference on MultiNLI test set. * denotes ensemble methods.

is a good option to enhance the performance.

Quora Question Pair Table 4 shows our results on the Quora question pair dataset. BiMPM using the multi-perspective matching technique between two sentences reports baseline performance of a L.D.C. network and basic multi-perspective models (Wang, Hamza, and Florian 2017). We obtained accuracies of 90.15% and 91.30% in single and ensemble methods, respectively, surpassing the previous state-of-the-art model of DIIN.

TrecQA and SelQA Table 5 shows the performance of different models on TrecQA and SelQA datasets for answer sentence selection task that aims to select a set of candidate answer sentences given a question. Most competitive models (Shen, Yang, and Deng 2017; Bian et al. 2017; Wang, Hamza, and Florian 2017; Shen et al. 2017) also use attention methods for words alignment between question and candidate answer sentences. However, the proposed DRCN using collective attentions over multiple layers, achieves the new state-of-the-art performance, exceeding the current state-of-the-art performance significantly on both datasets.

Analysis

Ablation study We conducted an ablation study on the SNLI dev set as shown in Table 6, where we aim to exam-

Models	Accuracy (%)
Siamese-LSTM (Wang, Hamza, and Florian 2017)	82.58
MP LSTM (Wang, Hamza, and Florian 2017)	83.21
L.D.C. (Wang, Hamza, and Florian 2017)	85.55
BiMPM (Wang, Hamza, and Florian 2017)	88.17
pt-DecAttchar.c (Tomar et al. 2017)	88.40
DIIN (Gong, Luo, and Zhang 2018)	89.06
DRCN	90.15
DIIN* (Gong, Luo, and Zhang 2018)	89.84
DRCN*	91.30

Table 4: Classification accuracy for paraphrase identification on Quora question pair test set. * denotes ensemble methods.

Models	MAP	MRR
<i>Raw version</i>		
aNMM (Yang et al. 2016)	0.750	0.811
PWIM (He and Lin 2016)	0.758	0.822
MP CNN (He, Gimpel, and Lin 2015)	0.762	0.830
HyperQA (Tay, Luu, and Hui 2017)	0.770	0.825
PR+CNN (Rao, He, and Lin 2016)	0.780	0.834
DRCN	0.804	0.862
<i>clean version</i>		
HyperQA (Tay, Luu, and Hui 2017)	0.801	0.877
PR+CNN (Rao, He, and Lin 2016)	0.801	0.877
BiMPM (Wang, Hamza, and Florian 2017)	0.802	0.875
Comp.-Aggr. (Bian et al. 2017)	0.821	0.899
IWAN (Shen, Yang, and Deng 2017)	0.822	0.889
DRCN	0.830	0.908

(a) TrecQA: raw and clean

Models	MAP	MRR
CNN-DAN (Santos, Wadhawan, and Zhou 2017)	0.866	0.873
CNN-hinge (Santos, Wadhawan, and Zhou 2017)	0.876	0.881
ACNN (Shen et al. 2017)	0.874	0.880
AdaQA (Shen et al. 2017)	0.891	0.898
DRCN	0.925	0.930

(b) SelQA

Table 5: Performance for answer sentence selection on TrecQA and selQA test set.

ine the effectiveness of our word embedding technique as well as the proposed densely-connected recurrent and co-attentive features. Firstly, we verified the effectiveness of the autoencoder as a bottleneck component in (2). Although the number of parameters in the DRCN significantly decreased as shown in Table 2, we could see that the performance was rather higher because of the regularization effect. Secondly, we study how the technique of mixing trainable and fixed word embeddings contributes to the performance in models (3-4). After removing E^{tr} or E^{fix} in eq. (1), the performance degraded, slightly. The trainable embedding E^{tr} seems more effective than the fixed embedding E^{fix} . Next, the effectiveness of dense connections was tested in models (5-9). In (5-6), we removed dense connections only over co-attentive or recurrent features, respectively. The result shows that the dense connections over attentive features are more effective. In (7), we removed dense connections over both co-attentive and recurrent features, and the performance degraded to 88.5%. In (8), we replace dense connection with residual connection

Models	Accuracy (%)
(1) DRCN	89.4
(2) – autoencoder	89.1
(3) – E^{tr}	88.7
(4) – E^{fix}	88.9
(5) – dense(Attn.)	88.7
(6) – dense(Rec.)	88.8
(7) – dense(Rec. & Attn.)	88.5
(8) – dense(Rec. & Attn.)	88.7
+ res(Rec. & Attn.)	
(9) – dense(Rec. & Attn. & Emb)	88.4
+ res(Rec. & Attn.)	
(10) – dense(Rec. & Attn. & Emb)	87.8
(11) – dense(Rec. & Attn. & Emb) - Attn.	85.3

Table 6: Ablation study results on the SNLI dev sets.

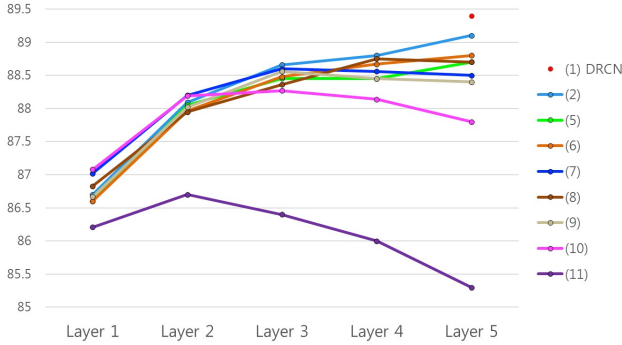


Figure 2: Comparison of models on every layer in ablation study. (best viewed in color)

only over recurrent and co-attentive features. It means that only the word embedding features are densely connected to the uppermost layer while recurrent and attentive features are connected to the upper layer using the residual connection. In (9), we removed additional dense connection over word embedding features from (8). The results of (8-9) demonstrate that the dense connection using concatenation operation over deeper layers, has more powerful capability retaining collective knowledge to learn textual semantics. The model (10) is the basic 5-layer RNN with attention and (11) is the one without attention. The result of (10) shows that the connections among the layers are important to help gradient flow. And, the result of (11) shows that the attentive information functioning as a soft-alignment is significantly effective in semantic sentence matching.

The performances of models having different number of recurrent layers are also reported in Fig. 2. The models (5-9) which have connections between layers, are more robust to the increased depth of network, however, the performances of (10-11) tend to degrade as layers get deeper. In addition, the models with dense connections rather than residual connections, have higher performance in general. Figure 2 shows that the connection between layers is essential, especially in deep models, endowing more representational power, and the dense connection is more effective than the residual connection.

Category	ESIM	DIIN	CAFE	DRCN
Matched				
Conditional	100	57	70	65
Word overlap	50	79	82	89
Negation	76	78	76	80
Antonym	67	82	82	82
Long Sentence	75	81	79	83
Tense Difference	73	84	82	82
Active/Passive	88	93	100	87
Paraphrase	89	88	88	92
Quantity/Time	33	53	53	73
Coreference	83	77	80	80
Quantifier	69	74	75	78
Modal	78	84	81	81
Belief	65	77	77	76
Mean	72.8	77.46	78.9	80.6
Stddev	16.6	10.75	10.2	6.7
Mismatched				
Conditional	60	69	85	89
Word overlap	62	92	87	89
Negation	71	77	80	78
Antonym	58	80	80	80
Long Sentence	69	73	77	84
Tense Difference	79	78	89	83
Active/Passive	91	70	90	100
Paraphrase	84	100	95	90
Quantity/Time	54	69	62	80
Coreference	75	79	83	87
Quantifier	72	78	80	82
Modal	76	75	81	87
Belief	67	81	83	85
Mean	70.6	78.53	82.5	85.7
Stddev	10.2	8.55	7.6	5.5

Table 7: Accuracy (%) of Linguistic correctness on MultiNLI dev sets.

Word Alignment and Importance Our densely-connected recurrent and co-attentive features are connected to the classification layer through the max pooling operation such that all max-valued features of every layer affect the loss function and perform a kind of deep supervision (Huang et al. 2017). Thus, we could cautiously interpret the classification results using our attentive weights and max-pooled positions. The attentive weights contain information on how two sentences are aligned and the numbers of max-pooled positions in each dimension play an important role in classification.

Figure 3 shows the attention map ($\alpha_{i,j}$ in eq. (5)) on each layer of the samples in Table 1. The Avg(Layers) is the average of attentive weights over 5 layers and the gray heatmap right above the Avg(Layers) is the rate of max-pooled positions. The darker indicates the higher importance in classification. In the figure, we can see that *tight*, *competing* and *bicycle* are more important words than others in classifying the label. The word *tight clothing* in the hypothesis can be inferred from *spandex* in the premise. And *competing* is also inferred from *race*. Other than that, the *riding* is matched with *pedaling*, and *pair* is matched with *two*. Judging by the matched terms, the model is undoubtedly able to classify the label as an entailment, correctly.

In Figure 3 (b), most of words in both the premise and the hypothesis coexist except *white* and *gray*. In attention

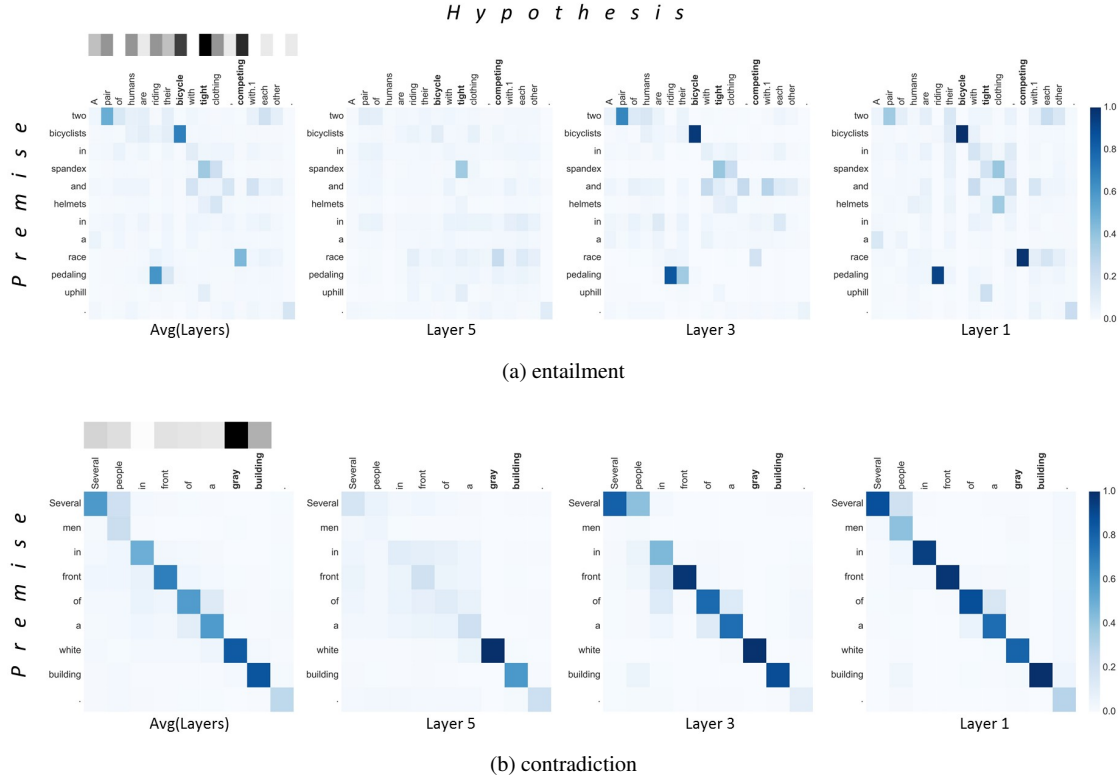


Figure 3: Visualization of attentive weights and the rate of max-pooled position. The darker, the higher. See supplementary materials for a comparison with other models that use the residual connections.

map of layer 1, the same or similar words in each sentence have a high correspondence (*gray* and *white* are not exactly matched but have a linguistic relevance). However, as the layers get deeper, the relevance between *white building* and *gray building* is only maintained as a clue of classification (See layer 5). Because *white* is clearly different from *gray*, our model determines the label as a contradiction.

The densely connected recurrent and co-attentive features are well-semanticized over multiple layers as collective knowledge. And the max pooling operation selects the soft-positions that may extract the clues on inference correctly.

Linguistic Error Analysis We conducted a linguistic error analysis on MultiNLI, and compared DRCN with the ESIM, DIIN and CAFE. We used annotated subset provided by the MultiNLI dataset, and each sample belongs to one of the 13 linguistic categories. The results in table 7 show that our model generally has a good performance than others on most categories. Especially, we can see that ours outperforms much better on the Quantity/Time category which is one of the most difficult problems. Furthermore, our DRCN shows the highest mean and the lowest stddev for both MATCHED and MISMATCHED problems, which indicates that it not only results in a competitive performance but also has a consistent performance.

Conclusion

In this paper, we introduce a densely-connected recurrent and co-attentive network (DRCN) for semantic sentence matching. We connect the recurrent and co-attentive features from the bottom to the top layer without any deformation. These intact features over multiple layers compose a community of semantic knowledge and outperform the previous deep RNN models using residual connections. In doing so, bottleneck components are inserted to reduce the size of the network. Our proposed model is the first generalized version of DenseRNN which can be expanded to deeper layers with the property of controllable feature sizes by the use of an autoencoder. We additionally show the interpretability of our model using the attentive weights and the rate of max-pooled positions. Our model achieves the state-of-the-art performance on most of the datasets of three highly challenging natural language tasks. Our proposed method using the collective semantic knowledge is expected to be applied to the various other natural language tasks.

References

- [Bian et al. 2017] Bian, W.; Li, S.; Yang, Z.; Chen, G.; and Lin, Z. 2017. A compare-aggregate model with dynamic-clip attention for answer selection. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, 1987–1990. ACM.
- [Bowman et al. 2015] Bowman, S. R.; Angeli, G.; Potts, C.; and Manning, C. D. 2015. A large annotated corpus for learning

- natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics.
- [Chen et al. 2017a] Chen, Q.; Zhu, X.; Ling, Z.-H.; and Inkpen, D. 2017a. Natural language inference with external knowledge. *arXiv preprint arXiv:1711.04289*.
- [Chen et al. 2017b] Chen, Q.; Zhu, X.; Ling, Z.-H.; Wei, S.; Jiang, H.; and Inkpen, D. 2017b. Enhanced lstm for natural language inference. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, 1657–1668.
- [Choi, Yoo, and goo Lee 2017] Choi, J.; Yoo, K. M.; and goo Lee, S. 2017. Learning to compose task-specific tree structures. AAAI.
- [Conneau et al. 2017] Conneau, A.; Kiela, D.; Schwenk, H.; Barrault, L.; and Bordes, A. 2017. Supervised learning of universal sentence representations from natural language inference data. *arXiv preprint arXiv:1705.02364*.
- [Csernai 2017] Csernai, K. 2017. Quora question pair dataset.
- [Ghaeini et al. 2018] Ghaeini, R.; Hasan, S. A.; Datla, V.; Liu, J.; Lee, K.; Qadir, A.; Ling, Y.; Prakash, A.; Fern, X. Z.; and Farri, O. 2018. Dr-bilstm: Dependent reading bidirectional lstm for natural language inference. *arXiv preprint arXiv:1802.05577*.
- [Gong, Luo, and Zhang 2018] Gong, Y.; Luo, H.; and Zhang, J. 2018. Natural language inference over interaction space. In *International Conference on Learning Representations*.
- [He and Lin 2016] He, H., and Lin, J. 2016. 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*, 937–948.
- [He et al. 2016] He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778.
- [He, Gimpel, and Lin 2015] He, H.; Gimpel, K.; and Lin, J. 2015. Multi-perspective sentence similarity modeling with convolutional neural networks. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, 1576–1586.
- [Huang et al. 2017] Huang, G.; Liu, Z.; Weinberger, K. Q.; and van der Maaten, L. 2017. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, volume 1, 3.
- [Im and Cho 2017] Im, J., and Cho, S. 2017. Distance-based self-attention network for natural language inference. *arXiv preprint arXiv:1712.02047*.
- [Jurczyk, Zhai, and Choi 2016] Jurczyk, T.; Zhai, M.; and Choi, J. D. 2016. Selqa: A new benchmark for selection-based question answering. In *Tools with Artificial Intelligence (ICTAI), 2016 IEEE 28th International Conference on*, 820–827. IEEE.
- [Liu et al. 2016] Liu, P.; Qiu, X.; Chen, J.; and Huang, X. 2016. Deep fusion lstms for text semantic matching. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, 1034–1043.
- [McCann et al. 2017] McCann, B.; Bradbury, J.; Xiong, C.; and Socher, R. 2017. Learned in translation: Contextualized word vectors. In *Advances in Neural Information Processing Systems*, 6297–6308.
- [Mikolov et al. 2013] Mikolov, T.; Sutskever, I.; Chen, K.; Corrado, G. S.; and Dean, J. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, 3111–3119.
- [Nie and Bansal 2017] Nie, Y., and Bansal, M. 2017. Shortcut-stacked sentence encoders for multi-domain inference. *arXiv preprint arXiv:1708.02312*.
- [Pavlick et al. 2015] Pavlick, E.; Bos, J.; Nissim, M.; Beller, C.; Van Durme, B.; and Callison-Burch, C. 2015. Adding semantics to data-driven paraphrasing. 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)*, volume 1, 1512–1522.
- [Pennington, Socher, and Manning 2014] Pennington, J.; Socher, R.; and Manning, C. D. 2014. Glove: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, 1532–1543.
- [Peters et al. 2018] Peters, M. E.; Neumann, M.; Iyyer, M.; Gardner, M.; Clark, C.; Lee, K.; and Zettlemoyer, L. 2018. Deep contextualized word representations. In *Proc. of NAACL*.
- [Radford et al. 2018] Radford, A.; Narasimhan, K.; Salimans, T.; and Sutskever, I. 2018. Improving language understanding by generative pre-training.
- [Rao, He, and Lin 2016] Rao, J.; He, H.; and Lin, J. 2016. Noise-contrastive estimation for answer selection with deep neural networks. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*, 1913–1916. ACM.
- [Romano et al. 2006] Romano, L.; Kouylekov, M.; Szpektor, I.; Dagan, I.; and Lavelli, A. 2006. Investigating a generic paraphrase-based approach for relation extraction. In *11th Conference of the European Chapter of the Association for Computational Linguistics*.
- [Santos, Wadhawan, and Zhou 2017] Santos, C. N. d.; Wadhawan, K.; and Zhou, B. 2017. Learning loss functions for semi-supervised learning via discriminative adversarial networks. *arXiv preprint arXiv:1707.02198*.
- [Shen et al. 2017] Shen, D.; Min, M. R.; Li, Y.; and Carin, L. 2017. Adaptive convolutional filter generation for natural language understanding. *arXiv preprint arXiv:1709.08294*.
- [Shen et al. 2018] Shen, T.; Zhou, T.; Long, G.; Jiang, J.; Wang, S.; and Zhang, C. 2018. Reinforced self-attention network: a hybrid of hard and soft attention for sequence modeling. *arXiv preprint arXiv:1801.10296*.
- [Shen, Yang, and Deng 2017] Shen, G.; Yang, Y.; and Deng, Z.-H. 2017. Inter-weighted alignment network for sentence pair modeling. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, 1179–1189.
- [Tay, Luu, and Hui 2017] Tay, Y.; Luu, A. T.; and Hui, S. C. 2017. Enabling efficient question answer retrieval via hyperbolic neural networks. *CoRR abs/1707.07847*.
- [Tay, Tuan, and Hui 2017] Tay, Y.; Tuan, L. A.; and Hui, S. C. 2017. A compare-propagate architecture with alignment factorization for natural language inference. *arXiv preprint arXiv:1801.00102*.
- [Tomar et al. 2017] Tomar, G. S.; Duque, T.; Täckström, O.; Uszkoreit, J.; and Das, D. 2017. Neural paraphrase identification of questions with noisy pretraining. *arXiv preprint arXiv:1704.04565*.
- [Vaswani et al. 2017] Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.; and Polosukhin, I. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, 6000–6010.
- [Wang, Hamza, and Florian 2017] Wang, Z.; Hamza, W.; and Florian, R. 2017. Bilateral multi-perspective matching for natural language sentences. *arXiv preprint arXiv:1702.03814*.
- [Wang, Smith, and Mitamura 2007] Wang, M.; Smith, N. A.; and Mitamura, T. 2007. What is the jeopardy model? a quasi-synchronous grammar for qa. In *Proceedings of the 2007 Joint*

Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL).

[Williams, Nangia, and Bowman 2017] Williams, A.; Nangia, N.; and Bowman, S. R. 2017. A broad-coverage challenge corpus for sentence understanding through inference. *arXiv preprint arXiv:1704.05426*.

[Wu et al. 2016] Wu, Y.; Schuster, M.; Chen, Z.; Le, Q. V.; Norouzi, M.; Macherey, W.; Krikun, M.; Cao, Y.; Gao, Q.; Macherey, K.; et al. 2016. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*.

[Yang et al. 2016] Yang, L.; Ai, Q.; Guo, J.; and Croft, W. B. 2016. anmm: Ranking short answer texts with attention-based neural matching model. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*, 287–296. ACM.

Supplementary Material

Datasets

A. SNLI is a collection of 570k human written sentence pairs based on image captioning, supporting the task of natural language inference (Bowman et al. 2015). The labels are composed of entailment, neutral and contradiction. The data splits are provided in (Bowman et al. 2015).

B. MultiNLI, also known as Multi-Genre NLI, has 433k sentence pairs whose size and mode of collection are modeled closely like SNLI. MultiNLI offers ten distinct genres (FACE-TO-FACE, TELEPHONE, 9/11, TRAVEL, LETTERS, OUP, SLATE, VERBATIM, GOVERNMENT and FICTION) of written and spoken English data. Also, there are matched dev/test sets which are derived from the same sources as those in the training set, and mismatched sets which do not closely resemble any seen at training time. The data splits are provided in (Williams, Nangia, and Bowman 2017).

C. Quora Question Pair consists of over 400k question pairs based on actual `quora.com` questions. Each pair contains a binary value indicating whether the two questions are paraphrase or not. The training-dev-test splits for this dataset are provided in (Wang, Hamza, and Florian 2017).

D. TrecQA provided in (Wang, Smith, and Mitamura 2007) was collected from TREC Question Answering tracks 8-13. There are two versions of data due to different pre-processing methods, namely clean and raw (Rao, He, and Lin 2016). We evaluate our model on both data and follow the same data split as provided in (Wang, Smith, and Mitamura 2007). We use official evaluation metrics of MAP (Mean Average Precision) and MRR (Mean Reciprocal Rank), which are standard metrics in information retrieval and question answering tasks.

E. SelQA consists of questions generated through crowdsourcing and the answer sentences are extracted from the ten most prevalent topics (Arts, Country, Food, Historical Events, Movies, Music, Science, Sports, Travel and TV) in the English Wikipedia. We also use MAP and MRR for our evaluation metrics, and the data splits are provided in (Jurczyk, Zhai, and Choi 2016).

Visualization on the comparable models

We study how the attentive weights flow as layers get deeper in each model using the dense or residual connection. We used the samples of the SNLI dev set in Table 1.

Figure 4 and 5 show the attention map on each layer of the models of DRCN, Table 6 (8), and Table 6 (9). In the model of Table 6 (8), we replaced the dense connection with the residual connection only over recurrent and co-attentive features. And, in the model of Table 6 (9), we removed additional dense connection over word embedding features from Table 6 (8). We denote the model of Table 6 (9) as Res1 and the model of Table 6 (8) as Res2 for convenience.

In Figure 4, DRCN does not try to find the right alignments at the upper layer if it already finds the rationale for the prediction at the relatively lower layer. This is expected that the DRCN use the features of all the preceding layers as a collective knowledge. While Res1 and Res2 have to find correct alignments at the top layer, however, there are some misalignments such as *competing* and *bicyclists* rather than *competing* and *race* in Res2 model.

In the second example in Figure 5, although the DRCN couldn’t find the clues at the lower layer, it gradually finds the alignments, which can be a rationale for the prediction. At the 5th layer of DRCN, the attentive weights of *gray building* and *white building* are significantly higher than others. On the other hand, the attentive weights are spread in several positions in both Res1 and Res2 which use residual connection.

Hypothesis

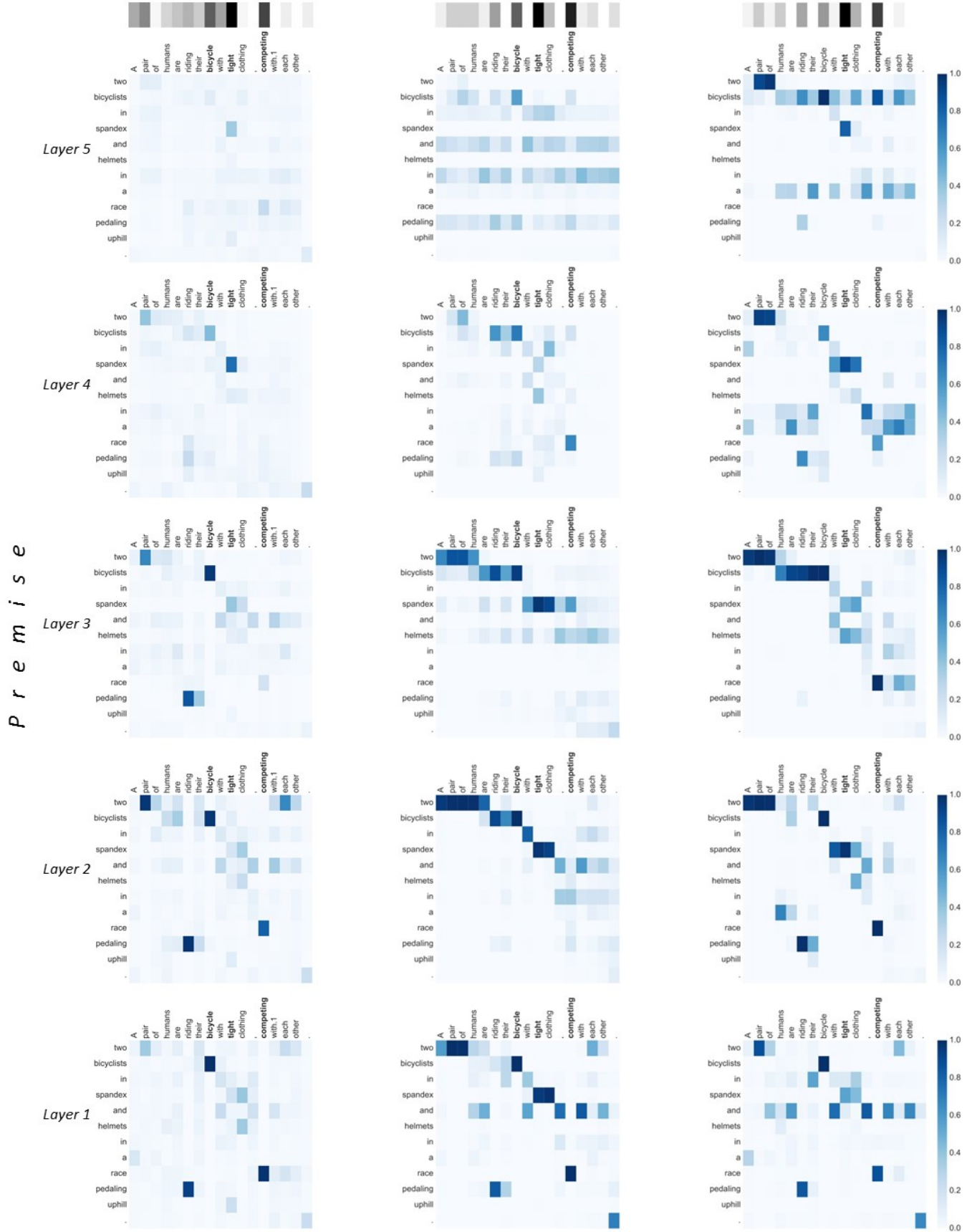


Figure 4: Visualization of attentive weights on the *entailment* example. The premise is “two bicyclists in spandex and helmets in a race pedaling uphill.” and the hypothesis is “A pair of humans are riding their bicycle with tight clothing, competing with each other.”. The attentive weights of DRCN, Res1, and Res2 are presented from left to right.

Hypothesis

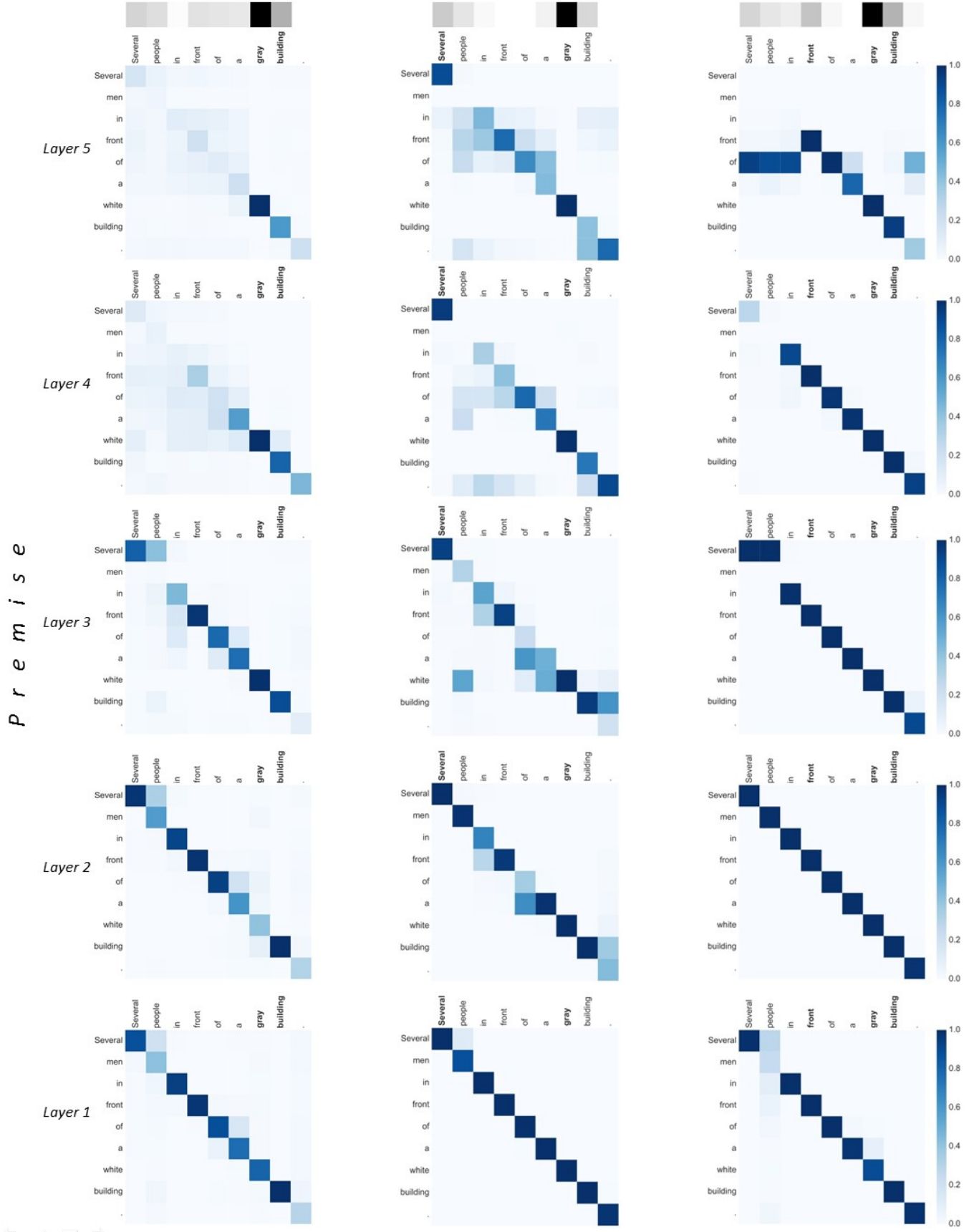


Figure 5: Visualization of attentive weights on the *contradiction* example. The premise is "Several men in front of a white building." and the hypothesis is "Several people in front of a gray building.". The attentive weights of DRCN, Res1, and Res2 are presented from left to right.