

# Learning Natural Language Inference with LSTM

**Shuohang Wang**

School of Information Systems  
Singapore Management University  
shwang.2014@phdis.smu.edu.sg

**Jing Jiang**

School of Information Systems  
Singapore Management University  
jingjiang@smu.edu.sg

## Abstract

Natural language inference (NLI) is a fundamentally important task in natural language processing that has many applications. The recently released Stanford Natural Language Inference (SNLI) corpus has made it possible to develop and evaluate learning-centered methods such as deep neural networks for the NLI task. In this paper, we propose a special long short-term memory (LSTM) architecture for NLI. Our model builds on top of a recently proposed neural attention model for NLI but is based on a significantly different idea. Instead of deriving sentence embeddings for the premise and the hypothesis to be used for classification, our solution uses a matching-LSTM that performs word-by-word matching of the hypothesis with the premise. This LSTM is able to place more emphasis on important word-level matching results. In particular, we observe that this LSTM remembers important mismatches that are critical for predicting the contradiction or the neutral relationship label. Our experiments on the SNLI corpus show that our model outperforms the state of the art, achieving an accuracy of 86.1% on the test data.

## 1 Introduction

Natural language inference (NLI) is the problem of determining whether from a sentence  $P$  one can infer another sentence  $H$  (MacCartney, 2009). Here  $P$  is called the *premise* and  $H$  the *hypothesis*. NLI is a fundamentally important problem that has applications in many tasks including question answering, semantic search and automatic text summarization. There has been much interest in NLI in the past decade, especially sur-

rounding the PASCAL Recognizing Textual Entailment (RTE) Challenge (Dagan et al., 2005). Existing solutions to NLI range from shallow approaches based on lexical similarities (Glickman et al., 2005) to advanced methods that consider syntax (Mehdad et al., 2009), perform explicit sentence alignment (MacCartney et al., 2008) or use formal logic (Clark and Harrison, 2009).

Recently, Bowman et al. (2015) released the Stanford Natural Language Inference (SNLI) corpus for the purpose of encouraging more learning-centered approaches to NLI. This corpus contains around 570K sentence pairs with three labels: *entailment*, *contradiction* and *neutral*. The size of the corpus makes it now feasible to train deep neural network models, which typically require a large amount of training data. Bowman et al. (2015) tested a straightforward architecture of deep neural networks for NLI. In their architecture, the premise and the hypothesis are each represented by a sentence embedding vector. The two vectors are then fed into a multi-layer neural network to train a classifier. Bowman et al. (2015) achieved an accuracy of 77.6% on the SNLI corpus when long short-term memory (LSTM) networks were used to obtain the sentence embeddings.

A more recent work by Rocktäschel et al. (2015) improved the performance by applying a neural attention model. While their basic architecture is still based on sentence embeddings for the premise and the hypothesis, a key difference is that the embedding of the premise takes into consideration the alignment between the premise and the hypothesis. This so-called *attention-weighted* representation of the premise was shown to push the accuracy to 83.5% on the SNLI corpus.

A limitation of the aforementioned two models is that they reduce both the premise and the hypothesis to a single embedding vector before matching them, i.e., in the end, they use two embedding vectors to perform sentence-level match-

ing. However, not all word or phrase-level matching results are equally important. For example, the matching between stop words in the two sentences is not likely to contribute much to the final prediction. Another example is that for a hypothesis to *contradict* a premise, a single word or phrase-level mismatch (e.g., a mismatch of the subjects of the two sentences) may be sufficient and other word or phrase-level matching results are less important, but this intuition is hard to be captured if we directly match two sentence embeddings.

In this paper, we propose a new LSTM-based architecture for learning natural language inference. Different from the models in (Bowman et al., 2015) and (Rocktäschel et al., 2015), our prediction is not based on whole sentence embeddings of the premise and the hypothesis. Instead, we use an LSTM to perform *word-by-word* matching of the hypothesis with the premise. Our LSTM sequentially processes the hypothesis, and at each position, it tries to match the current word in the hypothesis with an attention-weighted representation of the premise. Matching results that are critical for the final prediction will be “remembered” by the LSTM while less important matching results will be “forgotten.” We refer to this architecture for natural language inference a match-LSTM, or *mLSTM* for short.

Using the SNLI corpus, we show that our *mLSTM* model improves the state-of-the-art performance on this data set by achieving a classification accuracy of 86.1%. Through qualitative analyses, we also show that the *mLSTM* architecture can indeed pick up the more important word-level matching results that need to be remembered for the final prediction. In particular, we observe that “good” word-level matching results are generally forgotten but important mismatches, which often indicate a *contradiction* or a *neutral* relationship, tend to be remembered.

## 2 Model

In this section, we present our *mLSTM* architecture for natural language inference.

### 2.1 Background

We first present the model by Rocktäschel et al. (2015) because our model builds on top of theirs. Their model uses LSTMs to process the premise and the hypothesis separately, with a neural attention component to find a soft alignment between

the two sentences.

### LSTM

Let us first briefly review LSTM. LSTM is a special form of recurrent neural networks (RNNs), which process sequence data. LSTM uses a few gate vectors at each position to control the passing of information along the sequence and thus improves the modeling of long-range dependencies. While there are different variations of LSTMs, here we present the one adopted in (Rocktäschel et al., 2015). Specifically, let us use  $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$  to denote an input sequence, where  $\mathbf{x}_k \in \mathbb{R}^l$  ( $1 \leq k \leq N$ ). At each position  $k$ , there is a set of internal vectors, including an input gate  $\mathbf{i}_k$ , a forget gate  $\mathbf{f}_k$ , an output gate  $\mathbf{o}_k$  and a memory cell  $\mathbf{c}_k$ . All these vectors will be used together to generate a hidden state  $\mathbf{h}_k$ . The following transition equations define the LSTM architecture:

$$\begin{aligned} \mathbf{i}_k &= \sigma(\mathbf{W}^i \mathbf{x}_k + \mathbf{V}^i \mathbf{h}_{k-1} + \mathbf{b}^i), \\ \mathbf{f}_k &= \sigma(\mathbf{W}^f \mathbf{x}_k + \mathbf{V}^f \mathbf{h}_{k-1} + \mathbf{b}^f), \\ \mathbf{o}_k &= \sigma(\mathbf{W}^o \mathbf{x}_k + \mathbf{V}^o \mathbf{h}_{k-1} + \mathbf{b}^o), \\ \mathbf{c}_k &= \mathbf{f}_k \odot \mathbf{c}_{k-1} + \mathbf{i}_k \odot \tanh(\mathbf{W}^c \mathbf{x}_k + \mathbf{V}^c \mathbf{h}_{k-1} + \mathbf{b}^c), \\ \mathbf{h}_k &= \mathbf{o}_k \odot \tanh(\mathbf{c}_k), \end{aligned} \quad (1)$$

where  $\sigma$  is the sigmoid function,  $\odot$  is the element-wise multiplication of two vectors, and all  $\mathbf{W} \in \mathbb{R}^{d \times l}$ ,  $\mathbf{V} \in \mathbb{R}^{d \times d}$  and  $\mathbf{b} \in \mathbb{R}^d$  are weight matrices and weight vectors to be learned.

### Neural Attention Model

For the natural language inference task, we have two sentences  $\mathbf{X}^s = (\mathbf{x}_1^s, \mathbf{x}_2^s, \dots, \mathbf{x}_M^s)$  and  $\mathbf{X}^t = (\mathbf{x}_1^t, \mathbf{x}_2^t, \dots, \mathbf{x}_N^t)$ , where  $\mathbf{X}^s$  is the premise and  $\mathbf{X}^t$  is the hypothesis. Here each  $\mathbf{x}$  is an embedding vector of the corresponding word and can be initialized using some pre-trained word embedding vectors. The goal is to predict a label  $y$  that indicates the relationship between  $\mathbf{X}^s$  and  $\mathbf{X}^t$ . In this paper, we assume  $y$  is one of *entailment*, *contradiction* and *neutral*.

Rocktäschel et al. (2015) first used two LSTMs to process the premise and the hypothesis, respectively, but initialized the second LSTM (for the hypothesis) with the last cell state of the first LSTM (for the premise). Let us use  $\mathbf{h}_j^s$  ( $1 \leq j \leq M$ ) and  $\mathbf{h}_k^t$  ( $1 \leq k \leq N$ ) to denote the resulting hidden states corresponding to  $\mathbf{x}_j^s$  and  $\mathbf{x}_k^t$ , respectively. The main idea of the word-by-word attention model proposed by Rocktäschel et al. (2015)

is to introduce a series of attention-weighted combinations of the hidden states of the premise, where each weighted version of the premise is for a particular word in the hypothesis. Let us use  $\mathbf{a}_k$  to denote such a vector for word  $\mathbf{x}_k^t$  in the hypothesis. We call these vectors  $\{\mathbf{a}_k\}_{k=1}^N$  the *attention vectors*. Specifically,  $\mathbf{a}_k$  is defined as follows<sup>1</sup>:

$$\mathbf{a}_k = \sum_{j=1}^M \alpha_{kj} \mathbf{h}_j^s, \quad (2)$$

where  $\alpha_{kj}$  is an attention weight that encodes the degree to which  $\mathbf{x}_k^t$  in the hypothesis is aligned with  $\mathbf{h}_j^s$  in the premise. The attention weight  $\alpha_{kj}$  is generated in the following way:

$$\alpha_{kj} = \frac{\exp(e_{kj})}{\sum_{j'} \exp(e_{kj'})}, \quad (3)$$

where

$$e_{kj} = \mathbf{w}^e \cdot \tanh(\mathbf{W}^s \mathbf{h}_j^s + \mathbf{W}^t \mathbf{h}_k^t + \mathbf{W}^a \mathbf{h}_{k-1}^a). \quad (4)$$

Here  $\cdot$  is the dot-product between two vectors, the vector  $\mathbf{w}^e \in \mathbb{R}^d$  and all matrices  $\mathbf{W} \in \mathbb{R}^{d \times d}$  contain weights to be learned, and  $\mathbf{h}_{k-1}^a$  is another hidden state which we will explain below.

The attention-weighted premise  $\mathbf{a}_k$  essentially tries to model the relevant parts in the premise with respect to  $\mathbf{x}_k^t$ , i.e., the  $k^{\text{th}}$  word in the hypothesis. Rocktäschel et al. (2015) further built an RNN model over  $\{\mathbf{a}_k\}_{k=1}^N$  by defining the following hidden states:

$$\mathbf{h}_k^a = \mathbf{a}_k + \tanh(\mathbf{V}^a \mathbf{h}_{k-1}^a), \quad (5)$$

where  $\mathbf{V}^a \in \mathbb{R}^{d \times d}$  is a weight matrix to be learned. We can see that the last  $\mathbf{h}_N^a$  aggregates all the previous  $\mathbf{a}_k$  and can be seen as an attention-weighted representation of the whole premise. Rocktäschel et al. (2015) then used this  $\mathbf{h}_N^a$ , which represents the whole premise, together with  $\mathbf{h}_N^t$ , which is an aggregated representation of the whole hypothesis<sup>2</sup>, to predict the label  $y$ .

<sup>1</sup>We present the word-by-word attention model by Rocktäschel et al. (2015) in a different way but the underlying model is the same. Our presentation is close to the one by Bahdanau et al. (2015), with our attention vectors  $\mathbf{a}$  corresponding to the context vectors  $\mathbf{c}$  in their paper.

<sup>2</sup>Strictly speaking, in (Rocktäschel et al., 2015),  $\mathbf{h}_N^t$  encodes both the premise and the hypothesis because the two sentences are chained. But  $\mathbf{h}_N^t$  places a higher emphasis on the hypothesis given the nature of RNNs.

## 2.2 Our Model

Although the neural attention model by Rocktäschel et al. (2015) achieved better results than Bowman et al. (2015), we see two limitations. First, the model still uses a single vector representation of the premise, namely  $\mathbf{h}_N^a$ , to match the entire hypothesis. We speculate that if we instead use each of the attention-weighted representations of the premise for matching, i.e. use  $\mathbf{a}_k$  at position  $k$  to match the hidden state  $\mathbf{h}_k^t$  of the hypothesis while we go through the hypothesis, we could achieve better matching quality. This can be done using an RNN which at each position takes in both  $\mathbf{a}_k$  and  $\mathbf{h}_k^t$  as its input and determines how well the overall matching of the two sentences is up to the current position. In the end the RNN will produce a single vector representing the matching of the two entire sentences.

The second limitation is that the model by Rocktäschel et al. (2015) does not explicitly allow us to place more emphasis on the more important matches between the premise and the hypothesis and down-weight the less critical matches. For example, matching of stop words is presumably less important than matching of content words. Also, some matching results may be particularly critical for making the final prediction and thus should be remembered. For example, consider the premise “A dog jumping for a Frisbee in the snow.” and the hypothesis “A cat washes his face and whiskers with his front paw.” When we sequentially process the hypothesis, once we see that the subject of the hypothesis *cat* does not match the subject of the premise *dog*, we have a high probability to believe that there is a contradiction. So this mismatch should be remembered.

Based on the two observations above, we propose to use an LSTM to sequentially match the two sentences. At each position the LSTM takes in both  $\mathbf{a}_k$  and  $\mathbf{h}_k^t$  as its input. The LSTM is expected to remember the more important matches between the two sentences for predicting the final label  $y$  and forget the less important ones. Figure 1 gives an overview of our model in contrast to the model by Rocktäschel et al. (2015).

Specifically, our model works as follows. First, similar to Rocktäschel et al. (2015), we process the premise and the hypothesis using two LSTMs, but we do not feed the last cell state of the premise to the LSTM of the hypothesis. This is because

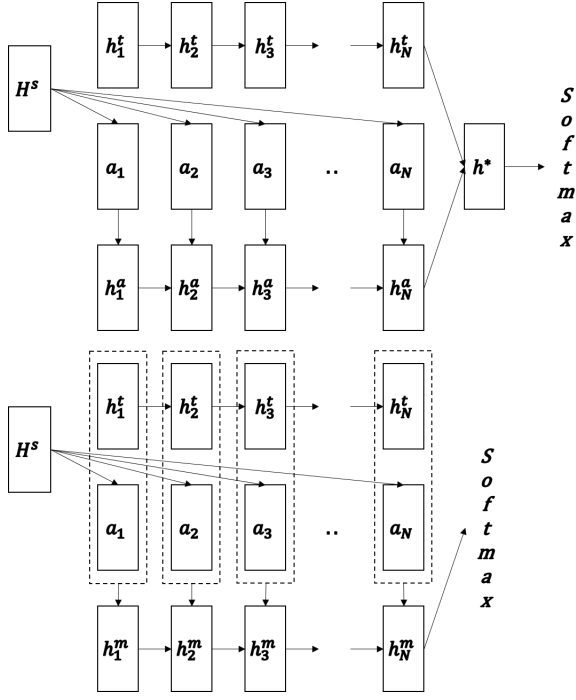


Figure 1: The top figure depicts the model by Rocktäschel et al. (2015) and the bottom figure depicts our model. Here  $\mathbf{H}^s$  represents all the hidden states  $\mathbf{h}_j^s$ . We can see that in the model by Rocktäschel et al. (2015), each  $\mathbf{h}_k^a$  represents a weighted version of the premise only, while in our model, each  $\mathbf{h}_k^m$  represents the matching between the premise and the hypothesis up to position  $k$ .

we do not need the LSTM for the hypothesis to encode any knowledge about the premise but we will match the premise with the hypothesis using the hidden states of the two LSTMs. Again, we use  $\mathbf{h}_j^s$  and  $\mathbf{h}_k^t$  to represent these hidden states.

Next, we generate the attention vectors  $\mathbf{a}_k$  similarly to Eqn (2). However, Eqn (4) will be replaced by the following equation:

$$e_{kj} = \mathbf{w}^e \cdot \tanh(\mathbf{W}^s \mathbf{h}_j^s + \mathbf{W}^t \mathbf{h}_k^t + \mathbf{W}^m \mathbf{h}_{k-1}^m). \quad (6)$$

The only difference here is that we use a hidden state  $\mathbf{h}^m$  instead of  $\mathbf{h}^a$ , and the way we define  $\mathbf{h}^m$  is very different from the definition of  $\mathbf{h}^a$ .

Our  $\mathbf{h}_k^m$  is the hidden state at position  $k$  generated from our *mLSTM*. This LSTM models the *matching* between the premise and the hypothesis. Important matches will be “remembered” by the LSTM while non-essential ones will be “forgotten.” We use the concatenation of  $\mathbf{a}_k$ , which is the attention-weighted version of the premise for the  $k^{\text{th}}$  word in the hypothesis, and  $\mathbf{h}_k^t$ , the hid-

den state for the  $k^{\text{th}}$  word itself, as input to the *mLSTM*.

Specifically, let us define

$$\mathbf{m}_k = \begin{bmatrix} \mathbf{a}_k \\ \mathbf{h}_k^t \end{bmatrix}. \quad (7)$$

We then build the *mLSTM* as follows:

$$\begin{aligned} \mathbf{i}_k^m &= \sigma(\mathbf{W}^{mi} \mathbf{m}_k + \mathbf{V}^{mi} \mathbf{h}_{k-1}^m + \mathbf{b}^{mi}), \\ \mathbf{f}_k^m &= \sigma(\mathbf{W}^{mf} \mathbf{m}_k + \mathbf{V}^{mf} \mathbf{h}_{k-1}^m + \mathbf{b}^{mf}), \\ \mathbf{o}_k^m &= \sigma(\mathbf{W}^{mo} \mathbf{m}_k + \mathbf{V}^{mo} \mathbf{h}_{k-1}^m + \mathbf{b}^{mo}), \\ \mathbf{c}_k^m &= \mathbf{f}_k^m \odot \mathbf{c}_{k-1}^m + \mathbf{i}_k^m \odot \tanh(\mathbf{W}^{mc} \mathbf{m}_k + \mathbf{V}^{mc} \mathbf{h}_{k-1}^m + \mathbf{b}^{mc}), \\ \mathbf{h}_k^m &= \mathbf{o}_k^m \odot \tanh(\mathbf{c}_k^m). \end{aligned} \quad (8)$$

With this *mLSTM*, finally we use only  $\mathbf{h}_N^m$  to predict the label  $y$ .

### 2.3 Implementation Details

Besides the *mLSTM* architecture, which is the main difference of our model from the model by Rocktäschel et al. (2015), we also introduce a few other changes.

First, we insert a special word *NULL* to the premise, and we allow words in the hypothesis to be aligned with this *NULL*. This is inspired by common practice in machine translation. Specifically, we introduce a vector  $\mathbf{h}_0^s$ , which is fixed to be a vector of 0s of dimension  $d$ . This  $\mathbf{h}_0^s$  represents *NULL* and is used together with other  $\mathbf{h}_j^s$  to derive the attention vectors  $\{\mathbf{a}_k\}_{k=1}^N$ .

Second, we use word embeddings trained from GloVe (Pennington et al., 2014) instead of word2vec vectors. The main reason is that GloVe word embeddings cover more words in the SNLI corpus than word2vec<sup>3</sup>.

Third, for words which do not have pre-trained word embeddings, we take the average of the embeddings of all the words (in GloVe) surrounding the unseen word within a window size of 9 (4 on the left and 4 on the right) as an approximation of the embedding of this unseen word. Then we do not update any word embedding when learning our model. Although this is a very crude approximation, it reduces the number of parameters we need to update, and as it turns out, we can still achieve better performance than Rocktäschel et al. (2015).

<sup>3</sup>The SNLI corpus contains 37K unique tokens. Around 12.1K of them cannot be found in word2vec but only around 4.1K of them cannot be found in GloVe.

| Model  | $d$ | $ \theta _{W+M}$ | $ \theta _M$ | Train | Dev         | Test        |
|--|-----|------------------|--------------|-------|-------------|-------------|
| LSTM [Bowman et al. (2015)]                        | 100 | 10M              | 221K         | 84.4  | -           | 77.6        |
| Classifier [Bowman et al. (2015)]                  | -   | -                | -            | 99.7  | -           | 78.2        |
| LSTM shared [Rocktäschel et al. (2015)]            | 159 | 3.9M             | 252K         | 84.4  | 83.0        | 81.4        |
| Word-by-word attention [Rocktäschel et al. (2015)] | 100 | 3.9M             | 252K         | 85.3  | 83.7        | 83.5        |
| Word-by-word attention (our implementation)        | 150 | 340K             | 340K         | 85.5  | 83.3        | 82.6        |
| <i>m</i> LSTM                                      | 150 | 544K             | 544K         | 91.0  | 86.2        | 85.7        |
| <i>m</i> LSTM with bi-LSTM sentence modeling       | 150 | 1.4M             | 1.4M         | 91.3  | 86.6        | 86.0        |
| <i>m</i> LSTM                                      | 300 | 1.9M             | 1.9M         | 92.0  | <b>86.9</b> | <b>86.1</b> |
| <i>m</i> LSTM with word embedding                  | 300 | 1.3M             | 1.3M         | 88.6  | 85.4        | 85.3        |

Table 1: Experiment results in terms of accuracy.  $d$  is the dimension of the hidden states.  $|\theta|_{W+M}$  is the total number of parameters and  $|\theta|_M$  is the number of parameters excluding the word embeddings. Note that the five models in the last section do not update word embeddings. The last three columns are the accuracies of the trained models on the training data, the development data and the test data, respectively.

### 3 Experiments

In this section, we present the evaluation of our model. We first perform quantitative evaluation, comparing our model with the model by Rocktäschel et al. (2015). We then conduct some qualitative analyses to understand how our *m*LSTM model works in matching the premise and the hypothesis.

#### 3.1 Experiment Settings

**Data:** We use the SNLI corpus to test the effectiveness of our model. The original data set contains 570,152 sentence pairs, each labeled with one of the following relationships: *entailment*, *contradiction*, *neutral* and  $-$ , where  $-$  indicates a lack of consensus from the human annotators. We discard the sentence pairs labeled with  $-$  and keep the remaining ones for our experiments. In the end, we have 549,367 pairs for training, 9,842 pairs for development and 9,824 pairs for testing. This follows the same data partition used by Bowman et al. (2015) in their experiments. We perform three-class classification and use accuracy as our evaluation metric.

**Parameters:** We use the Adam method (Kingma and Ba, 2014) with hyperparameters  $\beta_1$  set to 0.9 and  $\beta_2$  set to 0.999 for optimization. The initial learning rate is set to be 0.001 with a decay ratio of 0.95 for each iteration. The batch size is set to be 30. We experiment with  $d = 150$  and  $d = 300$  where  $d$  is the dimension of all the hidden states of the LSTMs.

**Methods for comparison:** We mainly want to compare our model with the word-by-word attention model by Rocktäschel et al. (2015) because this model achieved the state-of-the-art performance on the SNLI corpus. To ensure fair com-

parison, besides comparing with the accuracy reported by Rocktäschel et al. (2015), we also re-implemented their word-by-word attention model ourselves and report the performance of our implementation. We also consider a few variations of our model. Specifically, the following models are implemented and tested in our experiments:

- Word-by-word attention ( $d = 150$ ): This is our implementation of the word-by-word attention model by Rocktäschel et al. (2015), where we set the dimension of the hidden states to 150. The differences between our implementation and the original implementation by Rocktäschel et al. (2015) are the following: (1) We also add a *NULL* token to the premise for matching. (2) We do not feed the last cell state of the LSTM for the premise to the LSTM for the hypothesis, to keep it consistent with the implementation of our model. (3) For word representation, we also use the GloVe word embeddings as in the implementation of our model, and we do not update the word embeddings. For unseen words, we adopt the same strategy as described in Section 2.3.
- *m*LSTM ( $d = 150$ ): This is our *m*LSTM model with  $d$  set to 150.
- *m*LSTM with bi-LSTM sentence modeling ( $d = 150$ ): This is the same as the model above except that when we derive the hidden states  $\mathbf{h}_j^s$  and  $\mathbf{h}_k^t$  of the two sentences, we use bi-LSTMs instead of LSTMs. We implement this model to see whether bi-LSTMs allow us to better align the sentences.
- *m*LSTM ( $d = 300$ ): This is our *m*LSTM model with  $d$  set to 300.

- *mLSTM* with word embedding ( $d = 300$ ): This is the same as the model above except that we directly use the word embedding vectors  $\mathbf{x}_j^s$  and  $\mathbf{x}_k^t$  instead of the hidden states  $\mathbf{h}_j^s$  and  $\mathbf{h}_k^t$  in our model. In this case, each attention vector  $\mathbf{a}_k$  is a weighted sum of  $\{\mathbf{x}_j^s\}_{j=1}^M$ . We experiment with this setting because we hypothesize that the effectiveness of our model is largely related to the *mLSTM* that models the matching rather than the use of LSTMs to process the original sentences.

### 3.2 Quantitative Results

Table 1 compares the performance of the various models we tested together with some previously reported results. We have the following observations: (1) First of all, we can see that when we set  $d$  to 300, our model achieves an accuracy of 86.1% on the test data, which to the best of our knowledge is the highest on this data set. (2) If we compare our *mLSTM* model with our implementation of the word-by-word attention model by Rocktäschel et al. (2015) under the same setting with  $d = 150$ , we can see that our performance on the test data (85.7%) is still higher than that of the other model (82.6%). We also tested statistical significance and found the improvement to be statistically significant at the 0.001 level. (3) The performance of *mLSTM* with bi-LSTM sentence modeling compared with the model with standard LSTM sentence modeling when  $d$  is set to 150 shows that using bi-LSTM to process the original sentences helps (86.0% vs. 85.7% on the test data), but the difference is small. Therefore when we increased  $d$  to 300 we did not experiment with bi-LSTM sentence modeling. (4) Interestingly, when we experimented with the *mLSTM* model using the pre-trained word embeddings instead of LSTM-generated hidden states as initial representations of the premise and the hypothesis, we were able to achieve an accuracy of 85.3% on the test data, which is still better than previously reported state of the art. This suggests that the *mLSTM* architecture coupled with the attention model works well, regardless of whether or not we use LSTM to process the original sentences.

### 3.3 Qualitative Analyses

To obtain a better understanding of how our proposed model actually performs the matching between a premise and a hypothesis, we further conduct the following analyses. First, we look at the

learned word-by-word alignment weights  $\alpha_{kj}$  to check whether the soft alignment makes sense. This is the same as what was done in (Rocktäschel et al., 2015). We then look at the values of the various gate vectors of the *mLSTM*. By looking at these values, we are able to check (1) whether the model is able to differentiate between more important and less important word-level matching results, and (2) whether the model forgets certain matches and remembers certain other matches.

While we have looked at a random sample of sentence pairs, here we show only three examples. These three sentence pairs share the same premise but have different hypotheses and different relationship labels. They are given in Table 2. We can see that the first hypothesis is an entailment. The second hypothesis is a contradiction because it mentions a completely different event. The third hypothesis is neutral to the premise because the phrase “with his owner” cannot be inferred from the premise.

#### Word Alignment

First, let us look at the top-most plots of Figure 2, Figure 3 and Figure 4. These plots show the alignment weights  $\alpha_{kj}$  between the hypothesis and the premise, where a darker color corresponds to a larger value of  $\alpha_{kj}$ . Recall that  $\alpha_{kj}$  is the degree to which the word  $\mathbf{x}_k^t$  in the hypothesis is aligned with the word  $\mathbf{x}_j^s$  in the premise. Also recall that the weights  $\alpha_{kj}$  are configured such that for the same  $k$  all the  $\alpha_{kj}$  add up to 1. This means the weights in the same row in these plots add up to 1.

From the three plots we can see that the alignment weights generally make sense. For example, in Example 1, “animal” is strongly aligned with “dog” while “toy” aligned with “Frisbee.” The phrase “cold weather” is aligned with “snow.” In Example 3, we also see that “pet” is strongly aligned with “dog” while “game” aligned with “Frisbee.”

In Example 2, “cat” is strongly aligned with “dog” and “washes” is aligned with “jumping.” It may appear that these matches are wrong. However, “dog” is likely the best match for “cat” among all the words in the premise, and as we will show later, this match between “cat” and “dog” is actually a strong indication of a contradiction between the two sentences. The same explanation applies to the match between “washes” and “jumping.”

We also observe that some words are aligned

|            | ID        | sentence  | label         |
|------------|-----------|---|---------------|
| Premise    |           | A dog jumping for a Frisbee in the snow.                              |               |
| Hypothesis | Example 1 | An animal is outside in the cold weather, playing with a plastic toy. | entailment    |
|            | Example 2 | A cat washed his face and whiskers with his front paw.                | contradiction |
|            | Example 3 | A pet is enjoying a game of fetch with his owner.                     | neutral       |

Table 2: Three examples of sentence pairs with different relationship labels.

with the *NULL* token we inserted. For example, the word “is” in the hypothesis in Example 1 does not correspond to any word in the premise and is therefore aligned with *NULL*. The words “face” and “whiskers” in Example 2 and “owner” in Example 3 are also aligned with *NULL*. Intuitively, if some important content words in the hypothesis cannot find a match in the premise and are therefore aligned with *NULL*, then we should have a higher chance to believe that the relationship label is either contradiction or neutral.

### Values of Gate Vectors

Next, let us look at the values of the learned gate vectors of our *mLSTM* for the three examples. We show these values under the setting where  $d$  is set to 150. Each row of these plots corresponds to one of the 150 dimensions. Again, a darker color indicates a higher value.

An input gate controls whether the input at the current position should be used in deriving the final hidden state of the current position. From the three plots of the input gates of the three examples, we can observe that generally for stop words such as prepositions and articles the input gates have lower values, suggesting that the matches of these words in the hypothesis with the premise are less important for deriving the hidden state of the current position, and hence less important for predicting the final relationship label. On the other hand, content words such as nouns and verbs tend to have higher values of the input gates, which also makes sense because these words are generally more important for determining the relationship label between the two sentences. Overall, the observation with the input gates verifies our assumption that the *mLSTM* helps differentiate the more important word-level matching results from the less important ones.

Next, let us look at the forget gates. Recall that a forget gate controls the importance of the *previous* cell state in deriving the final hidden state of the current position. Higher values of a forget gate indicate that we need to remember the previous cell state and pass it on whereas lower values

indicate that we should probably forget the previous cell. From the three plots of the forget gates of the three examples, we can see that overall the colors are the lightest for Example 1, which is an entailment. Light colors correspond to low values, and in this case they suggest that when the hypothesis is an entailment of the premise, the *mLSTM* tends to forget the previous matching results. On the other hand, for Example 2 and Example 3, which are contradiction and neutral, we see generally darker colors, which correspond to higher values of the forget gates. In particular, in Example 2, we can see that the colors are consistently dark starting from the word “his” in the hypothesis until the end. We believe the explanation is that after the *mLSTM* processes the first three words of the hypothesis, “A cat washes,” it sees that the matchings between “cat” and “dog” and between “washes” and “jumping” are both strong indications of a contradiction, and therefore these matching results need to be remembered until the end of the *mLSTM* for the final softmax-based classifier to make a prediction.

We have also checked the plots for the forget gates of some other sentence pairs, and we observe that generally for entailment, the forget gates have low values, while for contradiction and neutral, the forget gates start to have high values from certain position of the hypothesis. We therefore hypothesize that the way the *mLSTM* works is as follows. It remembers important mismatches, which are useful for predicting the contradiction or the neutral relationship, and forgets good matches. At the end of the *mLSTM*, if no important mismatch is remembered, the final classifier then will likely predict entailment by default. Otherwise, depending on the kind of mismatch remembered, the classifier will predict either contradiction or neutral.

It is also interesting to point out that the values of the input gates seem to have a negative correlation with the values of the forget gates. In other words, at each position of the hypothesis, if the input gate has high values, then the forget gate tends to have low values, and vice versa.

For the output gates, we are not able to draw

any important conclusion except that the output gates seem to be positively correlated with the input gates but they tend to be darker than the input gates.

## 4 Related Work

There has been much work on natural language inference. Shallow methods rely mostly on lexical similarities but are shown to be robust. For example, Bowman et al. (2015) experimented with a lexicalized classifier-based method, which only uses lexical information to extract features used by a classifier, and the method achieves an accuracy of 78.2% on the SNLI corpus. More advanced methods use syntactic structures of the sentences to help matching them. For example, Mehdad et al. (2009) applied syntactic-semantic tree kernels for recognizing textual entailment. Because inference is essentially a logic problem, methods based on formal logic (Clark and Harrison, 2009) or natural logic (MacCartney, 2009) have also been proposed. A comprehensive review on existing work on natural language inference can be found in (Sammons et al., 2011).

The work most relevant to ours is the recently proposed neural attention model-based method by Rocktäschel et al. (2015), which we have detailed in previous sections. Neural attention models have recently been applied to some natural language processing tasks including machine translation (Bahdanau et al., 2014), abstractive summarization (Rush et al., 2015) and question answering (Hermann et al., 2015). Rocktäschel et al. (2015) showed that the neural attention model could help derive a better representation of the premise to be used to match the hypothesis, whereas in our work we also use it to derive representations of the premise that are used to sequentially match the words in the hypothesis.

The Stanford Natural Language Inference corpus is new and so far it has only been used in a few studies. Besides the work by Bowman et al. (2015) themselves and by Rocktäschel et al. (2015), there is another study by Vendrov et al. (2015) in which a *Skip-Thought* model by Kiros et al. (2015) was applied to the NLI task. The authors reported an accuracy of 81.5% on the data set. Because this accuracy is lower than the best performance reported by Rocktäschel et al. (2015) and because our main focus was to examine the effectiveness of our match-LSTM compared with the model by

Rocktäschel et al. (2015), we did not include their study for comparison in our experiments.

## 5 Conclusions

In this paper, we proposed a special LSTM architecture for the task of natural language inference. Based on a recent work by Rocktäschel et al. (2015), we first used neural attention models to derive attention-weighted vector representations of the premise. We then designed a match-LSTM that processes the hypothesis word by word while trying to match the hypothesis with the premise. Specifically the *mLSTM* takes in as input both the current hidden state of the hypothesis and an attention-weighted representation of the premise, and generates an output that represents the matching between the premise and the hypothesis up to the current position. As a result, the last hidden state of the *mLSTM* can be used for predicting the relationship between the premise and the hypothesis.

Experiments on the SNLI corpus showed that the *mLSTM* model outperformed the state-of-the-art performance reported so far on this data set. Moreover, closer analyses on the gate vectors revealed that our *mLSTM* indeed remembers and passes on important matching results, which are typically mismatches that indicate a contradiction or a neutral relationship between the premise and the hypothesis.

## References

- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.
- Dzmitry Bahdanau, HyungHyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *Proceedings of the International Conference on Learning Representations*.
- Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642.
- Peter Clark and Phil Harrison. 2009. An inference-based approach to recognizing entailment. In *Proceedings of the Text Analysis Conference*.



Ido Dagan, Oren Glickman, and Bernardo Magnini. 2005. The PASCAL Recognising Textual Entailment Challenge. In *Proceedings of the PASCAL Challenges Workshop on Recognizing Textual Entailment*.

Oren Glickman, Ido Dagan, and Moshe Koppel. 2005. Web based probabilistic textual entailment. In *Proceedings of the PASCAL Challenges Workshop on Recognizing Textual Entailment*.

Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Su-leyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. In *Advances in Neural Information Processing Systems*, pages 1684–1692.

Diederik Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Ryan Kiros, Yukun Zhu, Ruslan R Salakhutdinov, Richard Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Skip-thought vectors. In *Advances in Neural Information Processing Systems*, pages 3276–3284.

Bill MacCartney, Michel Galley, and Christopher D Manning. 2008. A phrase-based alignment model for natural language inference. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.

Bill MacCartney. 2009. *Natural Language Inference*. Ph.D. thesis, Stanford University.

Yashar Mehdad, Alessandro Moschitti1, and Fabio Massimo Zanzotto. 2009. SemKer: Syntactic/semantic kernels for recognizing textual entailment. In *Proceedings of the Text Analysis Conference*.

Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.

Tim Rocktäschel, Edward Grefenstette, Karl Moritz Hermann, Tomáš Kočiský, and Phil Blunsom. 2015. Reasoning about entailment with neural attention. *arXiv preprint arXiv:1509.06664*.

Alexander M Rush, Sumit Chopra, and Jason Weston. 2015. A neural attention model for abstractive sentence summarization. *arXiv preprint arXiv:1509.00685*.

Mark Sammons, VG Vinod Vydiswaran, and Dan Roth. 2011. Recognizing textual entailment. *Multilingual Natural Language Applications: From Theory to Practice*. Prentice Hall, Jun.

Ivan Vendrov, Ryan Kiros, Sanja Fidler, and Raquel Urtasun. 2015. Order-embeddings of images and language. *arXiv preprint arXiv:1511.06361*.

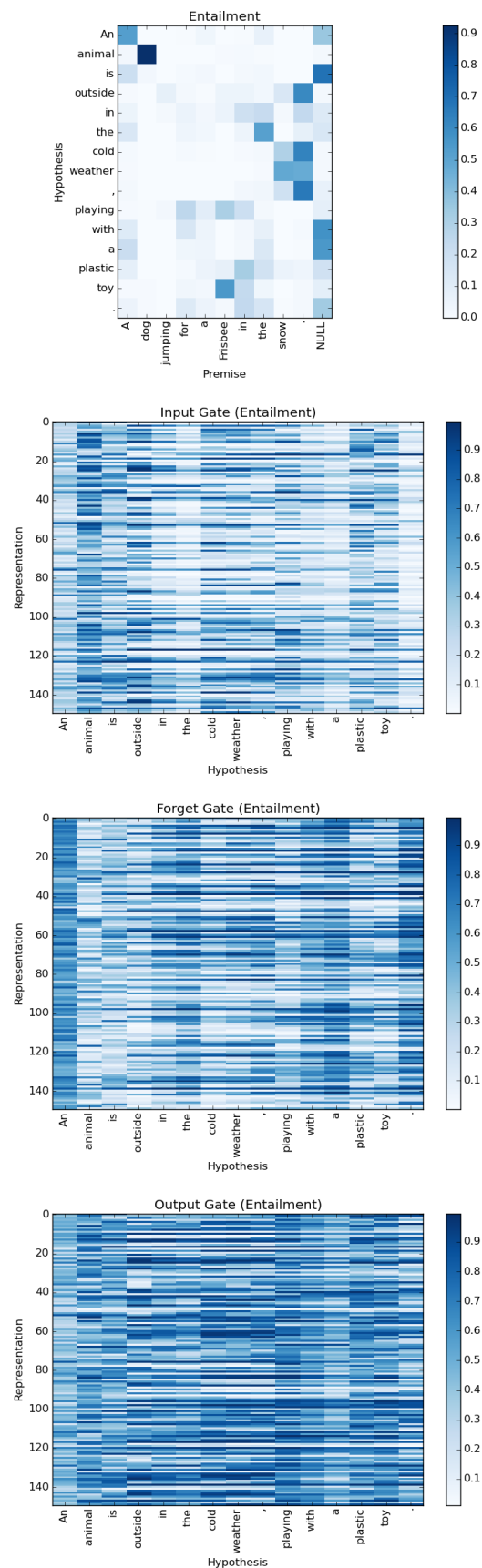


Figure 2: The alignment weights and gate vectors for Example 1 (entailment).

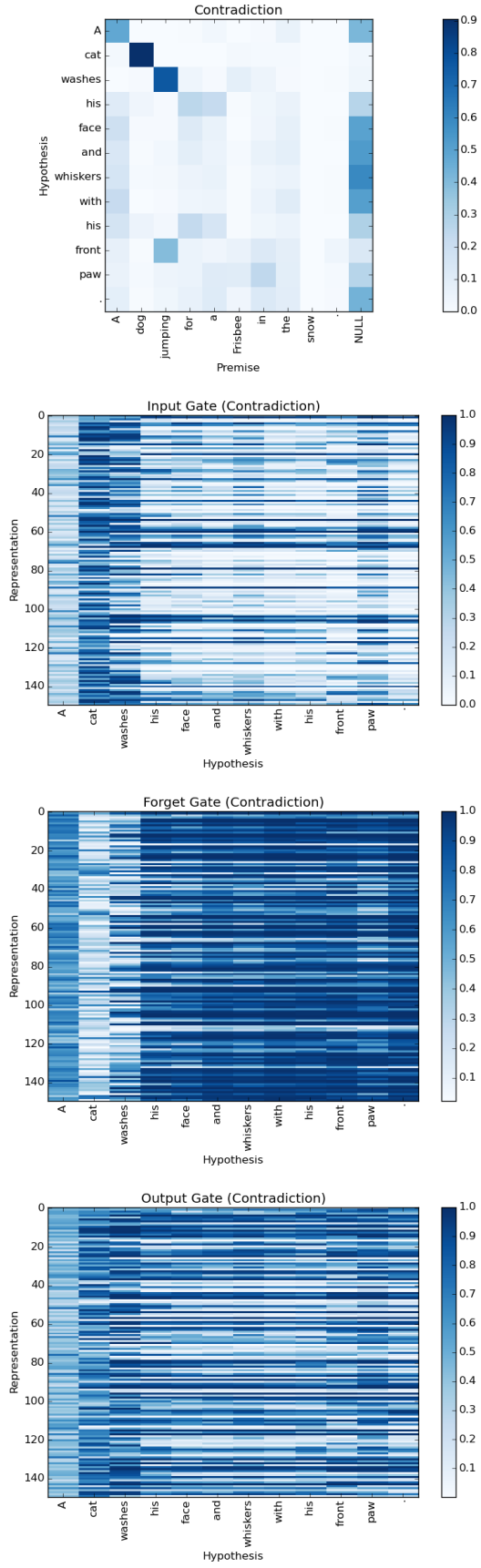


Figure 3: The alignment weights and gate vectors for Example 2 (contradiction).

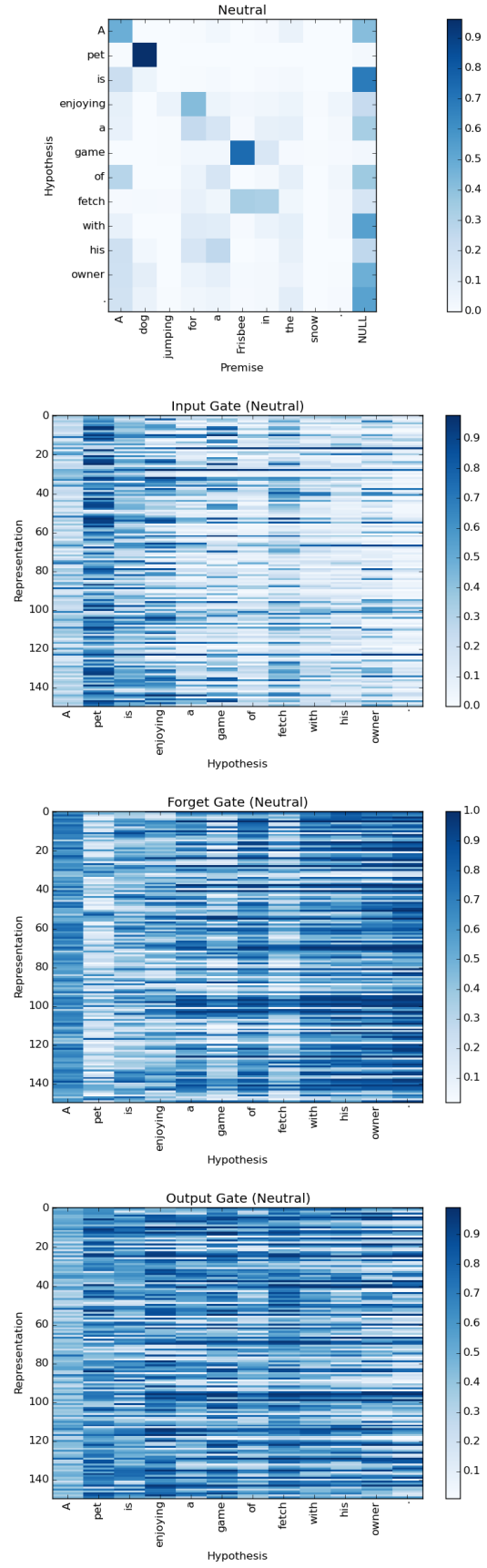


Figure 4: The alignment weights and gate vectors for Example 3 (neutral).