

# Recurrent Neural Networks with External Memory for Language Understanding

Baolin Peng<sup>1</sup>, Kaisheng Yao<sup>2</sup>

<sup>1</sup> The Chinese University of Hong Kong

<sup>2</sup> Microsoft Research

blpeng@se.cuhk.edu.hk, kaisheny@microsoft.com

## Abstract

Recurrent Neural Networks (RNNs) have become increasingly popular for the task of language understanding. In this task, a semantic tagger is deployed to associate a semantic label to each word in an input sequence. The success of RNN may be attributed to its ability to memorize long-term dependence that relates the current-time semantic label prediction to the observations many time instances away. However, the memory capacity of simple RNNs is limited because of the gradient vanishing and exploding problem. We propose to use an external memory to improve memorization capability of RNNs. We conducted experiments on the ATIS dataset, and observed that the proposed model was able to achieve the state-of-the-art results. We compare our proposed model with alternative models and report analysis results that may provide insights for future research.

**Index Terms:** Recurrent Neural Network, Language Understanding, Long Short-Term Memory, Neural Turing Machine

## 1. Introduction

Neural network based methods have recently demonstrated promising results on many natural language processing tasks [1, 2]. Specifically, recurrent neural networks (RNNs) based methods have shown strong performances, for example, in language modeling [3], language understanding [4], and machine translation [5, 6] tasks.

The main task of a language understanding (LU) system is to associate words with semantic meanings [7–9]. For example, in the sentence "Please book me a ticket from Hong Kong to Seattle", a LU system should tag "Hong Kong" as the departure-city of a trip and "Seattle" as its arrival city. The widely used approaches include conditional random fields (CRFs) [8, 10], support vector machine [11], and, more recently, RNNs [4, 12].

A RNN consists of an input, a recurrent hidden layer, and an output layer. The input layer reads each word and the output layer produces probabilities of semantic labels. The success of RNNs can be attributed to the fact that RNNs, if successfully trained, can relate the current prediction with input words that are several time steps away. However, RNNs are difficult to train, because of the gradient vanishing and exploding problem [13]. The problem also limits RNNs' memory capacity because error signals may not be able to back-propagated far enough.

There have been two lines of researches to address this problem. One is to design learning algorithms that can avoid gradient exploding, e.g., using gradient clipping [14], and/or gradient vanishing, e.g., using second-order optimization methods [15]. Alternatively, researchers have proposed more advanced model architectures, in contrast to the simple RNN that

uses, e.g., Elman architecture [16]. Specifically, the long short-term memory (LSTM) [17, 18] neural networks have three gates that control flows of error signals. The recently proposed gated recurrent neural networks (GRNN) [6] may be considered as a simplified LSTM with fewer gates.

Along this line of research on developing more advanced architectures, this paper focuses on a novel neural network architecture. Inspired by the recent work in [19], we extend the simple RNN with Elman architecture to using an external memory. The external memory stores the past hidden layer activities, not only from the current sentence but also from past sentences. To predict outputs, the model uses input observation together with a content retrieved from the external memory. The proposed model performs strongly on a common language understanding dataset and achieves new state-of-the-art results.

This paper is organized as follows. We briefly describe background of this research in Sec. 2. Section 3 presents details of the proposed model. Experiments are in section 4. We relate our research with other works in Sec. 5. Finally, we have conclusions and discussions in Sec. 6.

## 2. Background

### 2.1. Language understanding

A language understanding system predicts an output sequence with tags such as named-entity given an input sequence words. Often, the output and input sequences have been aligned. In these alignments, an input may correspond to a null tag or a single tag. An example is given in Table 1.

book	a	flight	from	Hong Kong	to	Seattle
-	-	-	-	Dpt-city	-	Arv-city

Table 1: An example of language understanding. Label names have been shortened to fit. Many words are labeled null or '-'.

Given a  $T$ -length input word sequence  $x_1^T$ , a corresponding output tag sequence  $y_1^T$ , and an alignment  $A$ , the posterior probability  $p(y_1^T | A, x_1^T)$  is approximated by

$$p(y_1^T | x_1^T) \approx \prod_{t=1}^T p(y_t | x_{t-k}^{t+k}), \quad (1)$$

where  $k$  is the size of a context window and  $t$  indexes the positions in the alignment.

## 2.2. Simple recurrent neural networks

The above posterior probability can be computed using a RNN. A RNN consists of an input layer  $x_t$ , a hidden layer  $h_t$ , and an output layer  $y_t$ . In Elman architecture [16], hidden layer activity  $h_t$  is dependent on both the input  $x_t$  and also recurrently on the past hidden layer activity  $h_{t-1}$ .

Because of the recurrence, the hidden layer activity  $h_t$  is dependent on the observation sequence from its beginning. The posterior probability is therefore computed as follows

$$\begin{aligned} p(y_1^T | x_1^T) &\approx \prod_{t=1}^T p(y_t | x_1^t) \\ &= \prod_{t=1}^T p(y_t | h_t, x_t) \end{aligned} \quad (2)$$

where the output  $y_t$  and hidden layer activity  $h_t$  are computed as

$$\begin{aligned} y_t &= g(h_t), \\ h_t &= \sigma(x_t, h_{t-1}). \end{aligned} \quad (3)$$

In the above equation,  $g(\cdot)$  is softmax function and  $\sigma(\cdot)$  is sigmoid or tanh function. The above model is denoted as simple RNN, to contrast it with more advanced recurrent neural networks described below.

## 2.3. Recurrent neural networks using gating functions

The current hidden layer activity  $h_t$  of a simple RNN is related to its past hidden layer activity  $h_{t-1}$  via the nonlinear function in Eq. (4). The non-linearity can cause errors back-propagated from  $h_t$  to explode or to vanish. This phenomenon prevents simple RNN from learning patterns that are spanned with long time dependence [14].

To tackle this problem, long short-term memory (LSTM) neural network was proposed in [17] with an introduction of memory cells, linearly dependent on their past values. LSTM also introduces three gating functions, namely input gate, forget gate and output gate. We follow a variant of LSTM in [18].

More recently, a gated recurrent neural network (GRNN) [6] was proposed. Instead of the three gating functions in LSTM, it uses two gates.

One is a reset gate  $r_t$  that relates a candidate activation with the past hidden layer activity  $h_{t-1}$ ; i.e.,

$$\hat{h}_t = \tanh(W_{xh}x_t + W_{hh}(r_t \odot h_{t-1})) \quad (5)$$

where  $\hat{h}_t$  is the candidate activation.  $W_{xh}$  and  $W_{hh}$  are the matrices relate the current observation  $x_t$  and the past hidden layer activity.  $\odot$  is element-wise product.

The second gate is an update gate  $z_t$  that interpolates the candidate activation and the past hidden layer activity to update the current hidden layer activity; i.e.,

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \hat{h}_t. \quad (6)$$

These gates are usually computed as functions of the current observation  $x_t$  and the past hidden layer activity; i.e.,

$$r_t = \sigma(W_{xr}x_t + W_{hr}h_{t-1}) \quad (7)$$

$$z_t = \sigma(W_{xz}x_t + W_{hz}h_{t-1}) \quad (8)$$

where  $W_{xr}$  and  $W_{hr}$  are the weights to observation and to the past hidden layer activity for the reset gate.  $W_{xz}$  and  $W_{hz}$  are similarly defined for the update gate.

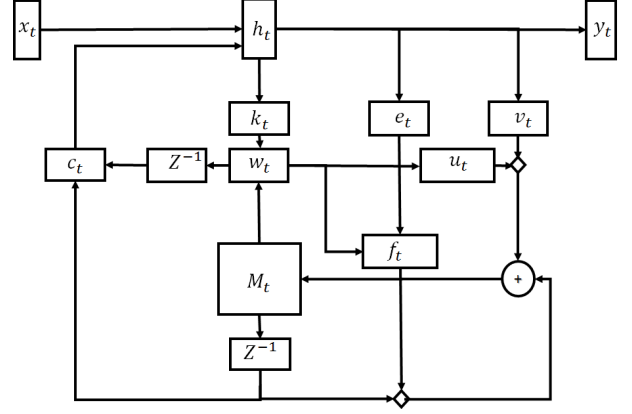


Figure 1: The RNN-EM model. The model reads input  $x_t$  and outputs  $y_t$ . Its hidden layer activity  $h_t$  depends on the input and the model's memory content retrieved in  $c_t$ .  $f_t$  and  $u_t$  are the forget and update gates.  $k_t$ ,  $e_t$  and  $v_t$  each denote key, erase and new content vector.  $M_t$  is the external memory.  $w_t$  is the weight and it is a function of  $k_t$  and  $M_t$ .  $Z^{-1}$  denotes a time-delay operator. The diamond symbol  $\diamond$  denotes diagonal matrix multiplication.

## 3. The RNN-EM architecture

We extend simple RNN in this section to using external memory. Figure 1 illustrates the proposed model, which we denote it as RNN-EM. Same as with the simple RNN, it consists of an input layer, a hidden layer and an output layer. However, instead of feeding the past hidden layer activity directly to the hidden layer as with the simple RNN, one input to the hidden layer is from a content of an external memory. RNN-EM uses a weight vector to retrieve the content from the external memory to use in the next time instance. The element in the weight vector is proportional to the similarity of the current hidden layer activity with the content in the external memory. Therefore, content that is irrelevant to the current hidden layer activity has small weights. We describe RNN-EM in details in the following sections. All of the equations to be described are with their bias terms, which we omit for simplicity of descriptions. We implemented RNN-EM using Theano [20, 21].

### 3.1. Model input and output

The input to the model is a dense vector  $x_t \in R^{d \times 1}$ . In the context of language understanding,  $x_t$  is a projection of input words, also known as word embedding.

The hidden layer reads both the input  $x_t$  and a content  $c_t$  vector from the memory. The hidden layer activity is computed as follows

$$h_t = \sigma(W_{ih}x_t + W_c c_t) \quad (9)$$

where  $\sigma(\cdot)$  is tanh function.  $W_{ih} \in R^{p \times d}$  is the weight to the input vector.  $c_t \in R^{m \times 1}$  is the content from a read operation to be described in Eq. (15).  $W_c \in R^{p \times m}$  is the weight to the content vector.

The output from this model is fed into the output layer as follows

$$y_t = g(W_{ho}h_t) \quad (10)$$

where  $W_{ho}$  is the weight to the hidden layer activity and  $g(\cdot)$  is softmax function.

Notice that in case of  $c_t = h_{t-1}$ , the above model is simple RNN.

### 3.2. External memory read

RNN-EM has an external memory  $M_t \in R^{m \times n}$ . It can be considered as a memory with  $n$  slots and each slot is a vector with  $m$  elements. Similar to the external memory in computers, the memory capacity of RNN-EM may be increased if using a large  $n$ .

The model generates a key vector  $k_t$  to search for content in the external memory. Though there are many possible ways to generate the key vector, we choose a simple linear function that relates hidden layer activity  $h_t$  as follows

$$k_t = W_k h_t \quad (11)$$

where  $W_k \in R^{m \times p}$  is a linear transformation matrix. Our intuition is that the memory should be in the same space of or affine to the hidden layer activity.

We use cosine distance  $K(u, v) = \frac{u \cdot v}{\|u\| \|v\|}$  to compare this key vector with contents in the external memory. The weight for the  $c$ -th slot  $M_t(:, c)$  in memory  $M_t$  is computed as follows

$$\hat{w}_t(c) = \frac{\exp \beta_t K(k_t, M_t(:, c))}{\sum_q \exp \beta_t K(k_t, M_t(:, q))} \quad (12)$$

where the above weight is normalized and sums to 1.0.  $\beta_t$  is a scalar larger than 0.0. It sharpens the weight vector when  $\beta_t$  is larger than 1.0. Conversely, it smooths or dampens the weight vector when  $\beta_t$  is between 0.0 and 1.0. We use the following function to obtain  $\beta_t$ ; i.e.,

$$\beta_t = \log(1 + \exp(W_\beta h_t)) \quad (13)$$

where  $W_\beta \in R^{1 \times p}$  maps the hidden layer activity  $h_t$  to a scalar.

Importantly, we also use a scalar coefficient  $g_t$  to interpolate the above weight estimate with the past weight as follows:

$$w_t = (1 - g_t)w_{t-1} + g_t \hat{w}_t \quad (14)$$

This function is similar to Eq. (6) in the gated RNN, except that we use a scalar  $g_t$  to interpolate the weight updates and the gated RNN uses a vector to update its hidden layer activity.

The memory content is retrieved from the external memory at time  $t - 1$  using

$$c_t = M_{t-1} w_{t-1}. \quad (15)$$

### 3.3. External memory update

RNN-EM generates a new content vector  $v_t$  to be added to its memory; i.e.,

$$v_t = W_v h_t \quad (16)$$

where  $W_v \in R^{m \times p}$ . We use the above linear function based on the same intuition in Sec. 3.2 that the new content and the hidden layer activity are in the same space of or affine to each other.

RNN-EM has a forget gate as follows

$$f_t = 1 - w_t \odot e_t \quad (17)$$

where  $e_t \in R^{n \times 1}$  is an erase vector, generated as  $e_t = \sigma(W_{he} h_t)$ . Notice that the  $c$ -th element in the forget gate is zero only if both read weight  $w_t$  and erase vector  $e_t$  have their  $c$ -th element set to one. Therefore, memory cannot be forgotten if it is not to be read.

Method	F1 score
CRF [26]	92.94
simple RNN [4]	94.11
CNN [27]	94.35
LSTM [28]	94.85
GRNN	94.82
RNN-EM	95.25

Table 2: F1 scores (in %) on ATIS.

RNN-EM has an update gate  $u_t$ . It simply uses the weight  $w_t$  as follows

$$u_t = w_t. \quad (18)$$

Therefore, memory is only updated if it is to be read.

With the above described two gates, the memory is updated as follows

$$M_t = \text{diag}(f_t)M_{t-1} + \text{diag}(u_t)v_t \quad (19)$$

where  $\text{diag}(\cdot)$  transforms a vector to a diagonal matrix with diagonal elements from the vector.

Notice that when the number of memory slots is small, it may have similar performances as a gated RNN. Specifically, when  $n = 1$ , Eqs. (19) and (6) are qualitatively similar.

## 4. Experiments

### 4.1. Dataset

In order to compare the proposed model with alternative modeling techniques, we conducted experiments on a well studied language understanding dataset, Air Travel Information System (ATIS) [22–24]. The training part of this dataset consists of 4978 sentences and 56590 words. There are 893 sentences and 9198 words for test. The number of semantic label is 127, including the common null label. We use lexicon-only features in experiments.

### 4.2. Comparison with the past results

The input  $x_t$  in RNN-EM has a window size of 3, consisting of the current input word and its neighboring two words. We use the AdaDelta method to update gradients [25]. The maximum number of training iterations was 50. Hyper parameters for tuning included the hidden layer size  $p$ , the number of memory slots  $n$ , and the dimension for each memory slot  $m$ . The best performing RNN-EM had 100 dimensional hidden layer and 8 memory slots with 40 dimensional memory slot.

Table 2 lists performance in F1 score of RNN-EM, together with the previous best results of alternative models in the literature. Since there are no previous results from GRNN, we use our own implementation of it for this study. These results are optimal in their respective systems. The previous best result was achieved using LSTM. A change of 0.38% of F1 score from LSTM result is significant at the 90% confidence level. Results in Table 2 show that RNN-EM is significantly better than the previous best result using LSTM.

### 4.3. Analysis on convergence and averaged performances

Results in the previous sections were obtained with models using different sizes. This section further compares neural network models given that they have approximately the same number of parameters, listed in Table 3. We use AdaDelta [25] gradient update method for all these models. Figure 2 plots their

Model	hidden layer dimension	# of Parameters
simple RNN	115	$\approx 7.4 * 10^3$
LSTM	50	$\approx 7.5 * 10^3$
GRNN	60	$\approx 7.4 * 10^3$
RNN-EM <sup>†</sup>	100,40 × 8	$\approx 7.3 * 10^3$

<sup>†</sup> 100 dimensional hidden layer, 40 dimensional slot with 8 slots.

Table 3: The size of each neural network models.

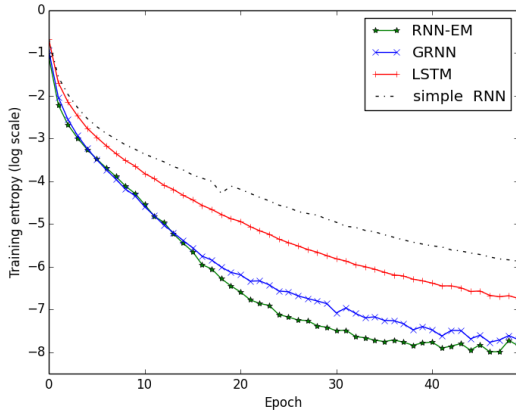


Figure 2: Convergence of training entropy. The entropy value has been converted to its logarithm.

training set entropy with respect to iteration numbers. To better illustrate their convergences, we have converted entropy values to their logarithms. The results show that RNN-EM converges to lower training entropy than other models. RNN-EM also converges faster than the simple RNN and LSTM.

We further repeated ATIS experiments for 10 times with different random seeds for these neural network models. We evaluated their performances after their convergences. Table 4 lists their averaged F1 scores, together with their maximum and minimum F1 scores. A change of 0.12% is significant at the 90% confidence level, when comparing against LSTM result. Results in Table 4 show that RNN-EM, on average, significantly outperforms LSTM. The best performance by RNN-EM is also significantly better than the best performing LSTM.

#### 4.4. Analysis on memory size

The size of the external memory  $M_t$  is proportional to the number of memory slots  $n$ . We fixed the dimension of memory slots to 40 and varied the number of slots. Table 5 lists their test set F1 scores. The best performing RNN-EM was with  $n = 8$ . Notice that RNN-EM with  $n = 1$  performed better than the simple RNN with 94.09% F1 score in Table 4. This can be explained as using gate functions in Eqs. (17) and (18) in RNN-EM, which are absent in simple RNNs. RNN-EM with  $n = 1$  also performed similarly as the gated RNN with 94.70% F1 score in Table 4, partly because of these gate functions.

Memory capacity may be measured using training set entropy. Table 5 shows that training set entropy is decreased initially with  $n$  increased from 1 to 8, showing that the memory capacity of the RNN-EM is improved. However, the entropy is increased with  $n$ s further increased. This suggests that memory

Method	Max	Min	Averaged
simple RNN	94.09	93.64	93.8
LSTM	94.81	94.62	94.73
GRNN	94.70	94.32	94.61
RNN-EM	95.22	94.71	94.96

Table 4: The maximum, minimum and averaged F1 scores (in %) by neural network models.

slot number $n$	1	2	4	8	16
F1 score	94.67	94.87	94.91	95.22	94.75
entropy $\times 10^3$	2.23	1.96	1.91	1.90	2.05
slot number $n$	32	64	128	256	512
F1 score	94.87	94.77	94.57	94.84	94.53
entropy $\times 10^3$	2.16	2.30	2.36	3.43	6.10

Table 5: Test set F1 scores (in %) and training set entropy by RNN-EM with different slot numbers.

capacity of RNN-EM cannot be increased simply by increasing the number of slots.

## 5. Related works

The RNN-EM is along the same line of research in [19, 29] that uses external memory to improve memory capacity of neural networks. Perhaps the closest work is the Neural Turing Machine (NTM) work in [19], which focuses on those tasks that require simple inference and has proved its effectiveness in copy, repeat and sorting tasks. NTM requires complex models because of these tasks. The proposed model is considerably simpler than NTM and can be considered as an extension of simple RNN. Importantly, we have shown through experiments on a common language understanding dataset the promising results from using the external memory architecture.

## 6. Conclusions and discussions

In this paper, we have proposed a novel neural network architecture, RNN-EM, that uses external memory to improve memory capacity of simple recurrent neural networks. On a common language understanding task, RNN-EM achieves new state-of-the-art results and performs significantly better than the previous best result using long short-term memory neural networks. We have conducted experiments to analyze its convergence and memory capacity. These experiments provide insights for future research directions such as mechanisms of accessing memory contents and methods to increase memory capacity.

## 7. Acknowledgement

The authors would like to thank Shawn Tan and Kai Sheng Tai for useful discussions on NTM structure and implementation.

## 8. References

- [1] Y. Bengio, R. Ducharme, P. Vincent, and C. Janvin, "A neural probabilistic language model," *Journal of Machine Learning Research*, vol. 3, pp. 1137–1155, 2003.
- [2] R. Collobert and J. Weston, "A unified architecture for natural language processing: deep neural networks with multitask learning," in *ICML*, 2008, pp. 160–167.

- [3] T. Mikolov, M. Karafiát, L. Burget, J. Cernocký, and S. Khudanpur, "Recurrent neural network based language model," in *INTERSPEECH*, 2010, pp. 1045–1048.
- [4] K. Yao, G. Zweig, M. Hwang, Y. Shi, and D. Yu, "Recurrent neural networks for language understanding," in *INTERSPEECH*, 2013, pp. 2524–2528.
- [5] J. Devlin, R. Zbib, Z. Huang, T. Lamar, R. M. Schwartz, and J. Makhoul, "Fast and robust neural network joint models for statistical machine translation," in *ACL*, 2014, pp. 1370–1380.
- [6] K. Cho, B. van Merriënboer, Ç. Gülçehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using RNN encoder-decoder for statistical machine translation," in *EMNLP*, 2014, pp. 1724–1734.
- [7] W. Ward *et al.*, "The cmu air travel information service: Understanding spontaneous speech," in *Proceedings of the DARPA Speech and Natural Language Workshop*, 1990, pp. 127–129.
- [8] C. Raymond and G. Riccardi, "Generative and discriminative algorithms for spoken language understanding," in *INTERSPEECH*, 2007, pp. 1605–1608.
- [9] R. de Mori, "Spoken language understanding: a survey," in *ASRU*, 2007, pp. 365–376.
- [10] J. D. Lafferty, A. McCallum, and F. C. N. Pereira, "Conditional random fields: Probabilistic models for segmenting and labeling sequence data," in *ICML*, 2001, pp. 282–289.
- [11] T. Kudo and Y. Matsumoto, "Chunking with support vector machines," in *NAACL*, 2001.
- [12] G. Mesnil, Y. Dauphin, K. Yao, Y. Bengio, L. Deng, D. Hakkani-Tur, X. He, L. Heck, G. Tur, D. Yu, and G. Zweig, "Using recurrent neural networks for slot filling in spoken language understanding," *IEEE/ACM Trans. on Audio, Speech, and Language Processing*, vol. 23, no. 3, pp. 530–539, 2015.
- [13] Y. Bengio, P. Y. Simard, and P. Frasconi, "Learning long-term dependencies with gradient descent is difficult," *IEEE Transactions on Neural Networks*, vol. 5, no. 2, pp. 157–166, 1994.
- [14] R. Pascanu, T. Mikolov, and Y. Bengio, "On the difficulty of training recurrent neural networks," in *ICML*, 2013, pp. 1310–1318.
- [15] J. Martens and I. Sutskever, "Training deep and recurrent networks with hessian-free optimization," in *Neural Networks: Tricks of the Trade - Second Edition*, 2012, pp. 479–535.
- [16] J. Elman, "Finding structure in time," *Cognitive science*, vol. 14, no. 2, pp. 179–211, 1990.
- [17] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [18] A. Graves, A. Mohamed, and G. E. Hinton, "Speech recognition with deep recurrent neural networks," in *ICASSP*, 2013, pp. 6645–6649.
- [19] A. Graves, G. Wayne, and I. Danihelka, "Neural Turing machines," *CoRR*, vol. abs/1410.5401, 2014.
- [20] F. Bastien, P. Lamblin, R. Pascanu, J. Bergstra, I. J. Goodfellow, A. Bergeron, N. Bouchard, and Y. Bengio, "Theano: new features and speed improvements," Deep Learning and Unsupervised Feature Learning NIPS 2012 Workshop, 2012.
- [21] J. Bergstra, O. Breuleux, F. Bastien, P. Lamblin, R. Pascanu, G. Desjardins, J. Turian, D. Warde-Farley, and Y. Bengio, "Theano: a CPU and GPU math expression compiler," in *Proceedings of the Python for Scientific Computing Conference (SciPy)*, Jun. 2010, oral Presentation.
- [22] D. Dahl, M. Bates, M. Brown, W. Fisher, K. Hunnicke-Smith, D. Pallett, C. Pao, A. Rudnicky, and E. Shriberg, "Expanding the scope of the ATIS task: The ATIS-3 corpus," in *Proceedings of the workshop on Human Language Technology*. Association for Computational Linguistics, 1994, pp. 43–48.
- [23] Y.-Y. Wang, A. Acero, M. Mahajan, and J. Lee, "Combining statistical and knowledge-based spoken language understanding in conditional models," in *COLING/ACL*, 2006, pp. 882–889.
- [24] G. Tur, D. Hakkani-Tr, and L. Heck, "What's left to be understood in ATIS?" in *IEEE Workshop on Spoken Language Technologies*, 2010.
- [25] M. D. Zeiler, "ADADELTA: An adaptive learning rate method," *arXiv:1212.5701*, 2012.
- [26] G. Mesnil, X. He, L. Deng, and Y. Bengio, "Investigation of recurrent-neural-network architectures and learning methods for language understanding," in *INTERSPEECH*, 2013.
- [27] P. Xu and R. Sarikaya, "Convolutional neural network based triangular CRF for joint intent detection and slot filling," in *ASRU*, 2013, pp. 78–83.
- [28] K. Yao, B. Peng, Y. Zhang, D. Yu, G. Zweig, and Y. Shi, "Spoken language understanding using long short-term memory neural networks," in *IEEE SLT*, 2014.
- [29] J. Weston, S. Chopra, and A. Bordes, "Memory networks," *submitted to ICLR*, vol. abs/1410.3916, 2015. [Online]. Available: <http://arxiv.org/abs/1410.3916>