

# Top-K Attention Mechanism for Complex Dialogue System

Chang-Uk Shin<sup>a</sup>, Jeong-Won Cha<sup>b\*</sup>

<sup>a</sup> Department of Eco-friendly Offshore plant FEED Engineering, Changwon National University, Republic of Korea

<sup>b</sup> Department of Computer Engineering, Changwon National University, Republic of Korea

{papower1, jcha}@changwon.ac.kr

## Abstract

Nowadays, natural language processing tasks such as dialogue modeling, question answering, sentence classification are usually attempted with a sequence or combination of Recurrent Neural Network(RNN), Convolutional Neural Network(CNN), and attention mechanism(Attention). Dialogue modeling using RNN encodes given history of the dialogue with RNN. Then the representation of the history is passed into the generator or classifier module to create or classify the system utterance. Most of the operation in RNN cannot be parallelize because the operation needs the operation result of previous timestep. Thus, the total inference time is relatively longer than CNN and Attention-based models. In this paper, we use CNN, Attention, and Pointer Network to model the 'Next Utterance Classification' task. We use Recall @ K, which measures the correct answer in the top K among the 100 candidates listed for the performance measure. The proposed system achieves R @ 1 20.92%, which exceeds the baseline performance. In this paper, we show that it is possible to extract useful information from long utterance history without RNN and to solve the next classification problem based on the information.

## Introduction

The dialogue system should be able to understand the user's utterance using the information of the given utterance history and the user's profile, and consequently to be able to solve the user's needs. The main purpose of the dialogue system research is to extract useful features from previous utterances and write the following utterances based on them.

As the natural language processing researches using the artificial neural network progresses actively, the artificial neural network model is also applied to the dialogue modeling task. At first, a chitchat model which receives one utterance by using RNN encoder-decoder structure and

responds appropriately was studied. After that, research continues how to encode the conversation and how to write the corresponding response.

The difficulty of dialog system modeling varies according to the given environment such as domains and rules. The more domains you are trying to learn, the larger the size of the dataset you need to learn the conversation model. Also, learning an artificial neural network model from a large amount of data sets is sometimes very tricky. For example, processing for a relatively long sequence with a small frequency, processing of words occurring at a low frequency, and time required for actual learning is increased in proportion to the size of the data set. Also, in the case of RNN, it needs the hidden state of the previous timestep when performing every timestep operation. Therefore, many operations cannot be performed in parallel, which degrades performance.

In this paper, we point out the computational efficiency of RNN and design and exploit the next utterance classification model using only CNN, top-k attention, and pointer network. We note the given conversation history by word, characterize it by convolution of the result and candidate, and finally, enter into the pointer network to select the final candidate.

All operations employ only CNN, attention, and pointer network so that there is no difference in reasoning time according to input length. Experiments on the ubuntu dialogue dataset of the DSTC7 track1 sentence selection with the proposed structure show that R @ 1 achieves 20.92%, which exceeds the performance of the dual encoder based model provided by the baseline, which is 8.32%.

Nonetheless, it took only 14% of the time to deduce the Dual Encoder. We demonstrate performance and reasoning time through experiments and verify the effectiveness of the proposed top-K attention.

## Related Research

[Lowe 2016a, Lowe 2016b] pointed out the inadequacy of the evaluation of the generation-based dialogue model and

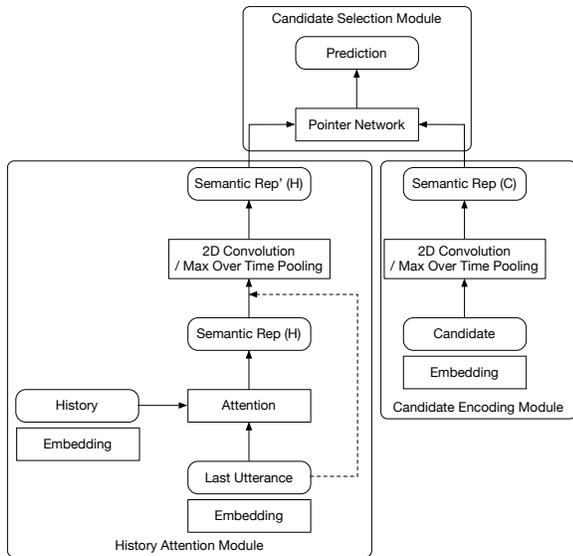
proposed the classification-based task NUC. The performance of the generation-based dialog system was measured based on the similarity of the generated system utterance to the pre-written utterance. However, this has been found to have a very weak correlation with actual human evaluation.

NUC is a task that analyzes the given context and selects the utterances that follow from the given candidates. They asserted the following points as advantages of this evaluation: the ability to control the degree of difficulty by adjusting the number of candidates, the intuitiveness of the performance measure, the ability to use the actual dialog system as an extension, can be controlled to provide strong constraints on the output, and that the response can be guaranteed to be fluent.

[Lowe 2016a] also pointed out that a large amount of corpus is essential for successfully modeling conversation tasks and has released a large-scale conversational data set collected in ubuntu IRC logs. The dataset is a multi-turn conversation of 1 million and is modeled as NUC. Since several speakers can speak at the same time in one chat room, they extracted only two-person conversations over three turns.

## Dataset

The dataset used for the experiment is the ubuntu dialogue dataset provided in the sentence selection track of DSTC7. It was quite preprocessed at the time of distribution, but I once again performed lemmatize using Stanford tagger [Toutanova and Manning 2003]. We collect package lists registered in the apt repository, tokenize package names, and delexicalized the URLs and file paths to URL / PATH



using regular expressions. Nevertheless, the size of the

Figure 1 : The proposed model architecture

dictionary is very large, and there are many restrictions on learning. Therefore, we changed all to the unknown, leaving only the top 10,000 words in frequency. The overall coverage of the top 10,000 words by frequency was 98.10%.

## Proposed Model

The proposed structure consists of three modules. The first is a history attention module that extracts optimal features using the relationship between the utterance history and the final utterance, the second is the candidate encoding module that encodes the given candidates, and the last is the candidate selection module that selects the final answer based on the results of the two modules.

### Top-K History Attention Module

The History attention module uses the previous  $m$  utterances and final utterances to extract the optimal features. We use  $m$  utterances to remove sentence separations and concatenate them word by word. We tested  $m$  by changing it to 1, 3, 5, 7, 10 and got the highest performance at 10. However, when we raised  $m$  to 10 or more, the memory usage was too large to test.

For efficient computation, we adopted a top- $k$  attention mechanism. We also calculate the degree of mutual concentration using one element and one sequence, as in the existing attention mechanism, in top- $k$  attention. The difference is that we do not perform weighted sum by multiplying all weights by all elements. Instead, we perform a weighted sum by extracting top  $k$  elements from the sequence that we believe are most closely related to the current element.

$$f_{ij} = W^{(2)}(\tanh(W^{(1)}(l_i; h_j))) \quad (1)$$

$$w_{ij} = \exp(f_{ij}) / (\sum_{k=0}^m \{\exp(f_{ik})\}) \quad (2)$$

$$t(x, X, k) = \begin{cases} x, & \text{if } x \geq \max(X, k) \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

$$a_j = \sum_{i=0}^n \{t(w_{ij}, w_{i,k}) h_i\} \quad (4)$$

where,  $l_i$  is the  $i$ -th word in the last utterance,  $h_j$  is the  $j$ -th word in the history utterance. In Eq. (1), each word in the history is connected to each word in the last utterance and input to the MLP that estimates the attention weight. The result,  $f$ , is taken by softmax and then multiplied by all words in the history to create an attended representation  $h$  [Equation 2, 4]. In this case, we do not use all of the results of softmax, but we set it so that only the top  $k$  values can be left through [Equation 3]. We call this operation

top-k attention. In this paper,  $k$  is set to 10 for all experiments.

We use the top-k attention mechanism to reduce the number of computations considerably and to characterize only the words most closely related to each word in the final utterance. In this way, we can cope with noise more robustly than existing attention mechanisms that receive attention from all words, so it is useful to extract only the information that our model needs.

In Experiment 5 to be described later, the performance is again improved by using the semantic representation (H), which is the result of the history attention module, and the word embedding of the last utterance, again (shown by a dotted line in Fig. (Huang, Liu, and Maaten, 2018), (Srivastava, Greff and Schmidhuber, 2015), (He et al., 2015)). In this paper, it can be understood that the embedding of the original word and the information extracted from the history are used together as the variance representation of the final utterance.

We obtain a single utterance variance representation by inputting the final attention and history attention results to the 2D CNN and the max-over-time-pooling layer. In all experiments, we set filter sizes to 2, 3 and 4, and the number of filters per each filter size to 32. Therefore, the size of the variance expression after performing max-over-time pooling of the equation is 96.

#### Candidate Encoding Module

In the experiment, we must choose one among the 100 candidate utterances. We enter each candidate into the 2D CNN structure described above to acquire a distributed representation. We use the same structure of CNN, but we do not share the weight.

#### Candidate Selection Module

We encode utterance distributed representations and candidates, respectively, which have been supplemented by utterance history up to the present. Two distributed representations are 96 dimensions. A structure for selecting one of the given candidates using a pointer network has been applied to machine reading comprehension(Wang, 2017). In this paper, we adopt the structure, connect the utterance distributed representation as many as the number of candidates, and select a candidate using the pointer network.

$$f_i = W^{(2)}(\tanh(W^{(1)}(a; c_i))) \quad (5)$$

$$p_i = \frac{\exp(f_i)}{\sum_{\{j=0\}} \exp(f_j)} \quad (6)$$

$$\hat{p} = \operatorname{argmax}(p_0, \dots, p_{100}) \quad (7)$$

where  $a$  denotes utterance distributed representation, and  $c$  is candidate. We perform attention between  $a$  and  $c$  using MLP and submit the candidate with the highest attention value as a prediction.

## Experiments

### Experimental Settings

In the experiment, we used a dataset of DSTC7 tack1. We randomly extracted 8/9 of the distributed train datasets and used them as the training dataset and the remaining 1/9 as the validation dataset. We used the distributed development dataset as a test dataset.

We used the performance of the dual encoder provided by DSTC7 as the baseline. Experiments were performed according to known experimental procedures. Dual Encoder is an artificial neural network structure proposed to perform NUC in (Rowe 2015). The Dual Encoder binds two RNNs with an RNN encoding the context and an RNN encoding the candidate. The distributed representation resulting from the RNN is entered into the regression module to model the degree to which it is matched appropriately.

Then, submit the candidate with the highest number among the given candidates as the final correct answer. We randomly sampled one of the 99 negative samples to prevent bias during learning and used CommonCrawl word embedding as pertained embedding.

We set the hyper-parameters as follows. The number of 2D CNN filters was 32, and the size of the filters was 2, 3, and 4. These are the same for history and candidate operations. We used Adam as an optimizer and set the learning-rate to 0.001 and epsilon to 0.01. All embedding sizes were set to 25. We then loaded the pretrained glove twitter 25d embedding and fine-tuned it.

Table 1 : The experimental results

Experiment	R@1	R@2	R@5	R@10	R@50	Inference time
Baseline	8.32%	13.36%	24.26%	35.98%	80.04%	2898ms
1 : without attention	10.76%	16.70%	27.80%	38.12%	80.10%	66ms
2 : without top-K attention	5.84%	9.90%	17.86%	27.40%	71.62%	239ms
3 : proposed architecture	14.56%	20.82%	31.96%	43.14%	83.38%	263ms
4 : multiple embedding	19.42%	27.92%	40.18%	51.38%	89.24%	374ms
5 : skip-connection	<b>20.92%</b>	<b>30.10%</b>	<b>43.04%</b>	<b>54.76%</b>	<b>91.40%</b>	394ms

Table 2 : The trend of performance by the length of history utterance

M	R@1	R@2	R@5	R@10	R@50	Inference time
3	12.82%	19.00%	29.44%	39.92%	81.44%	125ms
5	14.36%	21.04%	31.98%	42.16%	83.00%	177ms
7	14.30%	<b>20.90%</b>	<b>32.42%</b>	42.62%	<b>83.46%</b>	219ms
10 (Experiment 3)	<b>14.56%</b>	20.82%	31.96%	<b>43.14%</b>	83.38%	263ms
Without top-K attention	5.84%	9.90%	17.86%	27.40%	71.62%	239ms

### Experimental Results & Analysis

Table 1 shows the experimental results of the known baseline, the experimental results of the baseline we have tested and the experimental results of the proposed system. Experiment 1 is the experimental result using the simplest structure that selects the given candidate using only the last input utterance. In other words, we can encode the final utterance with the word 2D CNN and choose one of the candidates by pointer network value. Since we do not use history, we do the smallest amount of computation in the experiments we have presented, and nevertheless, achieve performance above the baseline.

Experiment 2 is an experiment in which history attention is added to Experiment 1. Providing additional information to the model can improve the performance of the model, and dialogue history is good qualities of conversation modeling. Experiment 2 performed all the words and attention of history. However, Experiment 1 consistently outperformed Experiment 2 from R @ 1 to R @ 50. We interpreted this result as meaningless to extract attention to all the words in history and not to extract appropriate information because it considers the attention of unnecessary words. Therefore, we apply the top-K attention to extract only relevant information and experiment 3 is applied.

Experiment 3 is an experiment in which top-k attention is added. We verify the effectiveness of the proposed top-k attention in this experiment. Experimental results showed a significant improvement in performance of R @ 1 by 3.80%p increase. The argmax operation to find top-k has been added, so the prediction time is somewhat longer. Nevertheless, it is still 1/10 of the baseline.

Experiments 4 and 5 are tuning experiments performed to improve the performance of the proposed method further. Experiment 4 is an experiment with improved embedding. We used glove word embedding together with word embedding learning with skip-gram and randomly initialized word embedding. All 3 word embedding sizes were set to 25. We have learned skip-gram word vectors by inputting all given learning corpus. In this way, the glove word vector learned by twitter can represent general knowledge in a vast domain, and the skip-gram word vector can express meaning in a specific domain. We also added 25-dimensional random initialized word embedding to allow the model to capture information that cannot be obtained in

context. And in our experiments, the best performance was achieved when all 3 word-embeddings were set to be fine-tuned, and the difference was about 5% at R @ 1.

Experiment 5 is an experiment in which the structure is tuned once again to the structure improved by Experiment 4. Only the results of attention calculation of the last utterance were input to the next 2D CNN layer. In Experiment 5, we used the word embedding of the final utterance by skip-connection with the final utterance attention result. It's a simple extension, but we've got a 1.50%p improvement in R @ 1 and a 2.18%p improvement in R @ 2. This is similar to the module proposed by Highway Networks (Srivastava, Greff and Schmidhuber, 2015). Semantically, it can be expressed as considering both the word embedding of the last utterance itself and the history attended by the previous history.

The samples in training corpus have average 5.49 history utterances (first quartile 3, median 5, and third quartile 7). The higher performance can be obtained by providing more history as input. However, providing many features will cause overfitting problem so that the performance will also decrease. 'Ubuntu dialogue corpus' has up to 75 histories, so setting the upper limit is a more appropriate modeling method. Therefore, we set the limit of the input history and did experiments. The results are summarized in [Table 2]. From the experimental results, it can be seen that as the number of history utterances increases, the performance improves and the inference time increases. Also, providing the history over the upper limit makes the model overfit, which shows the performance decline.

### Conclusion

Attempts to model human dialogue are still being researched. Also, RNN-based natural language processing research is continuing. Although the RNN based architecture can extract information regardless of the length of the input, parallelization is limited because the operation must be performed after the operation result of the previous timestep is prepared. Therefore, the computation time is proportional to the length of the input.

In this paper, we tried to solve 'next utterance classification' problem with word level 2D CNN, attention, and pointer network. We use attention just for  $k$  words which are executed between history and last utterance and which have the highest attention value. And we limited the num-

ber of input history. Experimental results on the DSTC7 track1 dataset show that the proposed model achieves approximately 10 times faster inferencing time than the dual encoder, which is the base model.

The top-k attention proposed in this paper extracts  $k$  words that are highly related to conventional attention and performs attention. We predicted that this could extract more useful information for reasoning than conventional attention and could achieve better performance.

Using the top-K attention and history length proposed in this paper, we search only the information within the threshold determined by the hyper-parameter. If we can adaptively change it, we can achieve better performance.

## References

- Ryan Lowe, Nissan Pow, Julian V. Serban and Joelle Pineau. 2016. The Ubuntu Dialogue Corpus: A Large Dataset for Research in Unstructured Multi-Turn Dialogue Systems. arXiv:1506.08909v3.
- Ryan Lowe, Julian V. Serban, Mike Noseworthy, Laurent Charlin and Joelle Pineau. 2016. On the Evaluation of Dialogue Systems with Next Utterance Classification. In Proceedings of the SigDial 2016, 264-269.
- Kristina Toutanova and Christopher D. Manning. 2003. Feature-Rich Part-of-Speech Tagging with a Cyclic Dependency Network. In Proceedings of HLT-NAACL, 252-259.
- Gao Huang, Zhuang Liu and Laurens van der Maaten. 2018. Densely Connected Convolutional Networks. arXiv:1608.06993v5.
- Rupesh Kumar Srivastava, Klaus Greff, Jürgen Schmidhuber. 2015. Highway Networks. arXiv: 1505.00387v2.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun, 2015. Deep Residual Learning for Image Recognition. arXiv:1512.03385v1.
- Wenhui Wang, Nan Yang, Furu Wei, Baobao Chang, Ming Zhou. 2017. Gated Self-Matching Networks for Reading Comprehension and Question Answering. Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, 189-198.