

Incorporating Loose-Structured Knowledge into LSTM with Recall Gate for Conversation Modeling

Zhen Xu¹, Bingquan Liu¹, Baoxun Wang², Chengjie Sun¹, Xiaolong Wang¹

¹School of Computer Science and Technology, Harbin Institute of Technology, Harbin, China

²Application and Service Group, Microsoft, Beijing, China

¹{z xu, bqliu, cjsun, wangxl}@insun.hit.edu.cn

²baoxwang@microsoft.com

Abstract

Modeling human conversations is the essence for building satisfying chat-bots with multi-turn dialog ability. Conversation modeling will notably benefit from domain knowledge since the relationships between sentences can be clarified due to semantic hints introduced by knowledge. In this paper, a deep neural network is proposed to incorporate background knowledge for conversation modeling. Through a specially designed Recall gate, domain knowledge can be transformed into the extra global memory of Long Short-Term Memory (LSTM), so as to enhance LSTM by cooperating with its local memory to capture the implicit semantic relevance between sentences within conversations. In addition, this paper introduces the loose structured domain knowledge base, which can be built with slight amount of manual work and easily adopted by the Recall gate. Our model is evaluated on the context-oriented response selecting task, and experimental results on both two datasets have shown that our approach is promising for modeling human conversations and building key components of automatic chatting systems.

reply almost all kinds of queries instead of completing tasks given by users as the normal dialog system (Mu and Yin, 2010) does. The ability of communicating like a real human is critical for keeping users' activity, thus various additional functions are possible to be introduced during the dialog and even the commercial services can be adopted (See Duer²).

Apparently, the conversation between humans takes *contextual relevance* as the essence with no doubt (Ribeiro et al., 2015), that is, the response should be semantic relevant with both the "direct" question and the history contents. Generally, it is a challenging task for Chat-bots to detect semantic clues of the conversation and provide the context-aware replies. For this purpose, the conversation should be well modeled, so that the semantic continuation and switching of the context can be sensed and the appropriate candidate replies is possible to be further selected.

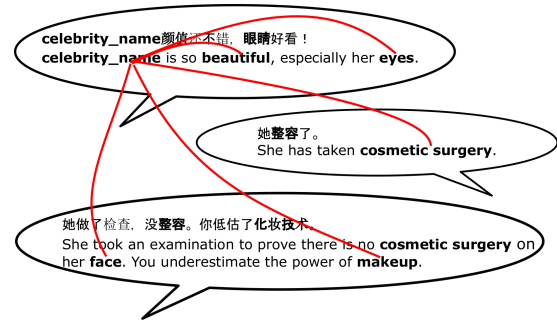


Figure 1: Conversation Example from SNS

1 Introduction

In recent years, the demand on Chat-bots has changed from answering simple questions as a toy to performing smooth open-domain conversations like real humans. Some good explorations have been conducted by the open-domain chat-bots like Clever-bot¹, etc., and their final goal is to

Intuitively, the sequence modeling approaches have great potential to capture the context semantic clues. Indeed, the Recurrent Neural Network (RNN) and Long Short-term Memory (LSTM) (Hochreiter and Schmidhuber, 1997) based methods have already utilized for this

¹<http://www.cleverbot.com/>

²<http://duer.baidu.com/>

task (Sutskever et al., 2014; Cho et al., 2014). The end-to-end learning ability of such models has already brought promising results to some non-trivial problems like sequence-to-sequence mapping (Shang et al., 2015).

Obviously, conversation modeling can benefit from human knowledge, only if the proper strategy is taken to adopt knowledge into machine learning models reasonably. As shown by the example in Figure 1, the semantic clue can be easily captured if the prior knowledge about *celebrity_name* is involved. In human conversations, such background knowledge generally takes the role of *global memory*, which can be naturally recalled by humans at the right moment. Heuristically, it is of great value to simulate the knowledge-supported conversation behavior of human beings, by capturing the acting mechanisms of human knowledge in conversation flows, so as to recognize semantic relevances in conversations more precisely and further find better responses for creating human-like chat-agents.

This paper proposes a deep neural network to address conversation modeling in the given domain. By introducing a new trainable gate to recall the global domain memory, our deep learning model incorporates background knowledge to enhance the sequence semantic modeling ability of LSTM. Methodologically, this paper concentrates on learning the implicit working mechanism of global memory in capturing semantic clues in conversations, including the functions of recalling and incorporating global memory with short-term memory. In addition, our model obtains global memory from the loose structured knowledge base. The knowledge base is built by organizing *entity-attribute* pairs as items which can be absorbed into our model by Recall gate directly. The preparation of such knowledge base requires little amount of manual work, by contrast, building complex-structured knowledge (e.g., graph) is indeed not a trivial task.

2 Related Work

Previous studies on conversation modeling are mainly task-completion oriented, such as act classification, state tracking in spoken field (Jurcicek et al., 2011; Williams et al., 2013; Henderson et al., 2013). With the rapid accumulation of available conversation data from SNS

(e.g., Twitter³, Weibo⁴), data-driven tasks like response ranking or generation have attracted more attention in this field (Serban et al., 2015b).

Early studies on conversation modeling focus on *one turn conversation* including one query and a response, thus heuristically the Information Retrieval (IR) or Statistical Machine Translation (SMT) based methods are taken to get responses. IR systems are built on movie scripts (Banchs and Li, 2012) or subtitles (Ameixa et al., 2014) to select responses related to the given questions. Ritter et al. (2011) introduce SMT to generate responses by taking query-response pairs as parallel corpus and get better performance than the IR approaches. To improve the readability of responses generated by SMT, Ji et al. (2014) propose an IR framework by introducing learning to rank strategy with the outputs of SMT as matching features.

The multi-turn conversation modeling task becomes even more challenging after taking history contents into consideration, thus complex approaches with better modeling ability are actually needed. Since it is natural to consider the conversation as a sequence of short texts, some studies introduce Neural Network based Sequence-to-Sequence (S2S) framework to address the conversation modeling issue. Vinyals and Le (2015) predict the next sentence based on the given one or several previous sentences with S2S directly. Serban et al. (2015a) and Shang et al. (2015) take response generation as the decoding process with the given distributed representation generated from previous sentences by the encoding process. The difference between (Serban et al., 2015a) and (Shang et al., 2015) lies in that Shang et al. (2015) add a feedback attention mechanism in the encoder-decoder framework.

Comparing with the generating strategies, response selecting approach is an option that avoids low readability naturally. Sordoni et al. (2015) propose an RNN model to evaluate the relevance between contexts, questions and their candidate response, and take the relevance as the feature to generate response. Lowe et al. (2015) and Kadlec et al. (2015) also present a ranking strategy on the Ubuntu Dialogue Corpus. In their work, an affinity model is used to measure the relevance between the context and a candidate reply,

³<https://twitter.com/>

⁴<http://weibo.com/>

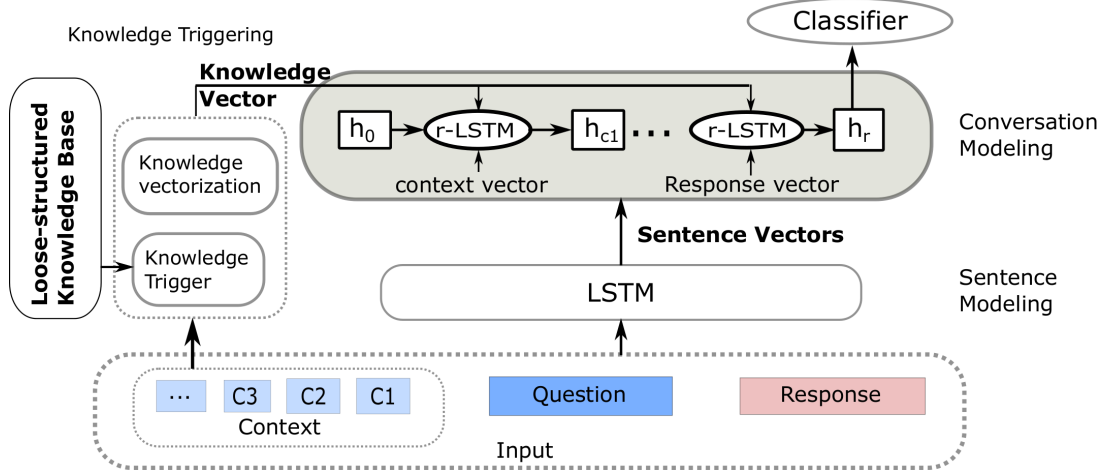


Figure 2: The Framework of Conversation Modeling by LSTM with Recall Gate(r-LSTM).

and the relevance score is take for response selection.

As pointed by Vinyals and Le (2015), the lack of general background knowledge is an obvious limitation of the current conversation modeling approaches. Noticing that previous studies rarely take background knowledge smoothly in the conversation modeling architectures, this paper tries to explore the effect of the background knowledge to the Neural Network based models.

3 LSTM with Recall Gate for Conversation Modeling

As mentioned in Section 1, our motivation of conversation modeling is to provide reliable evidence to detect users' context-aware queries and further select best answers based on the conversation history, for building automatic dialog systems. Apparently, the utterance⁵ in a conversation is very likely to be semantically relevant with both the history **context** and the **background knowledge**. The **context** is made up of the sequence of utterances appearing in conversation prior to the query and response. The semantic relevance between utterances implicitly performs as the clue for conversation modeling, meanwhile, background knowledge offers essential hints to enhance the effect of semantic clues. In this section, we propose a deep learning framework taking the specially designed Recall gate to smoothly integrate knowledge hints into LSTM for capturing such semantic clues.

⁵This paper takes utterance as the alias of the sentence in conversations. The usage of this word follows (Hurford et al., 2007)

3.1 Architecture

Figure 2 illustrates our conversation modeling framework. The proposed architecture is composed of sentence modeling, knowledge triggering, and conversation modeling components, the details of which will be given as follows:

Conversation Modeling: Since a conversation can be considered as a sequence of utterances that could be mapped into dense semantic vectors by the sentence modeling component, basically we take LSTM based methodology to address the conversation modeling task. As Figure 2 shows, there are two kinds of inputs: background knowledge and utterance vectors (including context, query and candidate response). Obviously, for the traditional LSTM, it is difficult to absorb the knowledge vectors properly.

In order to make background knowledge well involved into the deep neural network, this paper proposes a special Recall gate to enhance LSTM for conversation modeling. Inspired by our observation of human memory, the Recall-Gate is designed to convert the loose-structured domain knowledge to the global memory, which cooperates with the local memory in the cell of LSTM to provide evidence to judge whether an utterance is related to the dialog history or not. The details of LSTM with Recall gate(r-LSTM) will be given in Subsection 3.2.

In our work, the target of conversation modeling is to select better responses oriented to the given conversations (composed of the context and a query), so as to promote user's satisfaction in human-machine conversation. For this goal, a **binary classifier** is set to get the classifying confi-

dence and decide whether a candidate response is related to the current conversation or not.

Sentence modeling: As the basis of the entire task, sentence modeling aims to provide reasonable semantic representations as inputs of the upper layers. The RNN based methods, especially LSTM, have got good performance on sentence modeling recently, thus our framework adopts LSTM as the component mapping utterances into the real-valued semantic space (Sutskever et al., 2014; Cho et al., 2014). The process of sentence modeling introduced in this paper is shown in Figure 3.

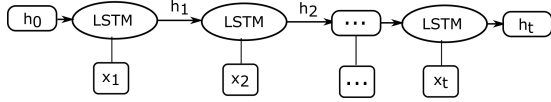


Figure 3: Sentence Modeling Based on LSTM.

The last hidden state h_t is always taken as a summary of the sentence, thus we adopt h_t as the representation of each utterance in conversation. In our implementation, the LSTM is first pre-trained as a language model for initialization, after that, it is tuned according to the end-to-end training of the network in Figure 2

Knowledge Triggering: As shown in the Figure 2, contexts will trigger the related background knowledge from the loose-structured knowledge base (KB), and the related knowledge will be taken as an input signal of our model and transformed into global memory. We take a simple knowledge trigger method, assuming that each entity can be correctly located in the KB.

For the given context, the process of generating related knowledge vector is as follows: First, the context is mapped to a bag of entities by matching words to a pre-defined vocabulary described in section 3.3. After that, attributes in pairs got from the previous step are ranked by their frequencies. We select top $N(= 10, 20, \dots)$ attributes as background knowledge of the given context. Finally the pairs are mapped to dense vectors (Sukhbaatar et al., 2015) by Equation 1:

$$kb = \sum_{i=1}^{i=N} \text{Embedding}(\text{attribute}_i) \quad (1)$$

where kb stands for the knowledge vector to be absorbed by our model. The embeddings of attributes are pre-trained on a open-domain corpus.

3.2 LSTM Cell with Recall Gate

Apparently, prior knowledge is of great value in the process of sequence understanding (Ovchinnikova, 2012). Nevertheless, integrating knowledge into deep learning architectures is still a challenging work. Since knowledge is indeed a global signal, and it is unwise to simply take knowledge as an additional utterance in conversation understanding, which is supported by the experimental results in Section 4. To make domain knowledge perform as a global memory for greater effect, we design a special model component for LSTM.

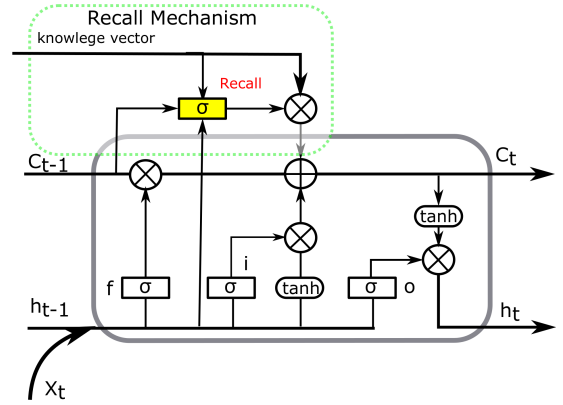


Figure 4: Illustration of the LSTM Cell with Recall Gate.

Figure 4 shows the detail of the LSTM cell with Recall gate (r-LSTM cell). Basically, our model absorbs the knowledge item embedding converted from the triggering results (as described by Subsection 3.1), and drives it cooperating with the previous sequence state and the present network input to summarize the current status for the next time-step, by adding a new recalling mechanism. Given previous hidden state h_{t-1} , previous memory c_{t-1} , current input x_t and background knowledge kb , we define the output of Recall gate as follows:

$$r_t = \sigma(W_{ri}[h_{t-1}, x_t] + W_{rc}c_{t-1} + W_{rk}kb + b_r) \quad (2)$$

where W_{ri} , W_{rc} and W_{rk} indicate the connection weights from current input(including previous state and x_t), memory cell and background knowledge respectively, and b_r is bias of the recall gate.

Generated by Recall gate, r_t indicates the proportion of global knowledge kb taking effect in the determination of current memory cell c_t . This pro-

cedure is described by Equation 3:

$$c_t = f_t * c_{t-1} + i_t * c_{input} + r_t * kb \quad (3)$$

The current hidden state h_t can be obtained by:

$$h_t = o_t * \tanh(c_t) \quad (4)$$

where o_t indicates the output gate of LSTM. The computations of forget gate, input gate, output gate and input cell are the same as normal LSTM. From the equations above it can be seen that in our methodology, external domain knowledge plays a significant role in generating the memory cell of LSTM and influencing the final hidden state consequently.

In the conversation model illustrated by Figure 2, the r-LSTM reads one utterance per time-step and generate a hidden state to represent the current conversation under the condition of global knowledge. With the multiplication to r_t shown in Equation 3, global knowledge kb including semantic clues takes effect on memory cell. In this way, the semantic relevance between utterances in conversation process could be sensed well by r-LSTM. As the higher-level global memory rather than the shallow input, external knowledge will be utilized more effectively in r-LSTM based conversation model.

Comparing with the normal LSTM, the r-LSTM could be more powerful in capturing semantic clues with better effective support of external knowledge because of the recall mechanism of Recall gate. In other words, the Recall-Gate makes r-LSTM to perform better on incorporating external knowledge. In contrast to r-LSTM, it is difficult for the general LSTM to involve external knowledge effectively.

3.3 Loose-Structured Knowledge Base

It is known that building a complex-structured knowledge (e.g., WordNet⁶, Yago⁷) requires large amount of human work. Obviously, an easily-built knowledge base is quite valuable for applications. This paper introduces loose-structured knowledge base composed of items with a flexible format “entity-attribute”, in which the *attributes* can be either other entities or the related keywords.

The extraction process of entity-attribute pairs is described as following: (1) For a given domain,

we extract entity or attributes from the domain-specific corpus by using statistic metrics such as tf-idf, entropy etc. After that, KL-divergence between domain-specific and general corpora to filter plain words and get entities and attributes for special domain; (2) According to the vocabulary composed of entities and attributes, counting the frequency of “entity-attribute” pair with a slide window. (3) The final knowledge base is obtained according to the frequency oriented statistics.

4 Experiment

In this part, our conversation modeling approach will be evaluated on *context-oriented best response selection* task, which is considered as a binary classification problem as illustrated by Figure 2, that is, the goal of our model is to give higher confidence to context related responses.

4.1 Experimental Setup

DataSet: We execute experiments on two datasets in special domains: Baidu TieBa Corpus and Ubuntu Corpus⁸. The Tieba dataset is composed of multi-turn conversations in the *celebrity* domain, extracted from crawled Tieba threads after several matching and filtering work. The Ubuntu corpus includes conversations collected from Ubuntu chat rooms focusing on technical supports.

For training the classifier, we adapt sampling method and dataset constructing strategy in Lowe et al. (2015) to generate training and testing set. For each positive sample in the training set, one negative sample is prepared. By contrast, for every positive sample in testing and validation set, there are 9 negative ones generated correspondingly. Finally, we construct a training set including 1 million conversations for both Tie and Ubuntu, and sample 50,000 conversations for validation and testing respectively.

Figure 5 shows the distributions of turn numbers of conversations in Tieba and Ubuntu respectively. Our model needs contexts to trigger background knowledge, so we select conversations with turn-number ranging from 3 to 7 for experiments.

Evaluation Metric: Since our model is applied in a classification problem, and the output confidence scores can be also used for ranking the can-

⁶<http://wordnet.princeton.edu/>

⁷<http://www.mpi-inf.mpg.de/departments/databases-and-information-systems/research/yago-naga/yago/>

⁸We have uploaded the Ubuntu dataset to https://www.dropbox.com/s/2fdn26rj6h9bpvl/ubuntu_data.zip?dl=0 and the Tieba corpus is confidential in the authors’ working organization.

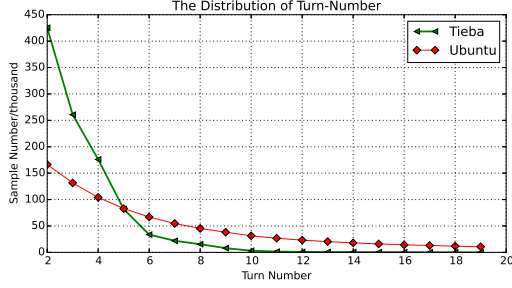


Figure 5: The Distribution of Turn Number.

didates, we take **Accuracy** to measure the performance of classifiers and introduce **Recall@k** to evaluate the ranking ability of the approaches. The metric Recall@k is applied in the work of Lowe et al. (2015) for the response selection task.

Baseline: To illustrate the performance of our model, we introduce five methods as the baselines, The descriptions of which are given as follows:

MLP and *LSTM*: Multi-Layer Perceptron(*MLP*) is the first baseline in our experiment. For our task, this model is designed to take the concatenated sentence vectors of both contexts and candidate responses as input, and outputs the binary classification results. *LSTM* is an intuitive method for our task, because this model can be directly used to model conversations by taking each utterance embedding as the input of each time-step.

It should be noted that neither of the baselines above considers external knowledge for prediction, and the architecture of *LSTM* is similar with Figure 2 but r-LSTM is not involved.

Affinity Model: As assumed by Lowe et al. (2015), context vectors and response vectors generate by RNN or LSTM can be aligned according to a relevance matrix, which can be formulated as:

$$p(r|c) = \sigma(c^T M r + b)$$

where c and r are the embedding results of context and response given by RNN or LSTM, and M indicates the relevance matrix. This methodology achieves the state-of-art results on the Ubuntu corpus without domain knowledge, thus the reason for introducing this baseline is to show the impact of reasonable knowledge integration.

MLP+kb and *LSTM+kb*: To examine the effectiveness of domain knowledge in conversation modeling, we incorporate knowledge into the baselines *MLP* and *LSTM* above in a trivial way, named as *MLP+kb* and *LSTM+kb*. *MLP+kb*

adopts knowledge vector by padding it to conversation vector(concatenated vectors of utterances), and *LSTM+kb* takes knowledge embeddings as the input of the first time-step.

Model Arch and Training Details:

Pre-training of word embeddings: For Ubuntu corpus, we take word vectors from (Lowe et al., 2015) as the initialization of embeddings. For Tieba corpus, the pre-trained character vectors on threads are utilized. All of the approaches (including baselines) in our experiments update word embeddings during the training procedure. The detailed dimension settings are given in Table 1.

dimension of	Ubuntu	Tieba
word embedding	300	100
sentence vector	200	100
knowledge vector	200	100
conversation vector	200	100

Table 1: Dimension settings of each part in the models.

For all the models, the training process runs until the loss of validation set in current iteration is larger than the previous one. All the models are trained on a K40m GPU with 12G memory.

4.2 Results and Analysis

Table 2 and Table 3 list the experimental results of all the approaches for the context-oriented response selection task, on Chinese Tieba and English Ubuntu corpus respectively.

From Table 2 and Table 3, it can be observed that our model (r-LSTM) notably outperforms the baseline methods on both datasets by all the metrics as expected. Basically, we ascribe the promotions to the fact that our model effectively captures semantic clues with the sequence modeling architecture of r-LSTM, meanwhile, benefited from smoothly incorporating the loose structured knowledge of the Recall Gate (see Figure 4), r-LSTM is able to detect implicit semantic hints in both conversation histories and candidate responses. In addition, the results also shows that vectors generated by LSTM based sentence modeling can represent sentences well as described in subsection 3.1.

The results of Affinity Model shown in Table 3 is lower than those reported in Lowe et al. (2015). There are to possible explanations for this observation: First, our experimental dataset extracted

Method	Acc	1 in 2 R@1	1 in 10 R@1	2 in 10 R@2	3 in 10 R@3	5 in 10 R@5
MLP	0.6834	0.7215	0.3118	0.4845	0.6080	0.7848
MLP+kb	0.6978	0.7660	0.3572	0.5440	0.6670	0.8240
LSTM	0.6890	0.7713	0.3696	0.5600	0.6810	0.8316
LSTM+kb	0.7140	0.7993	0.4121	0.6041	0.7172	0.8565
Affinity Model	RNN	0.6826	0.7286	0.3210	0.4896	0.6109
	LSTM	0.6925	0.7825	0.3814	0.5685	0.6804
r-LSTM	0.7260	0.8172	0.4348	0.6361	0.7527	0.8815

Table 2: Experiment Results on Tieba Corpus.

Method	Acc	1 in 2 R@1	1 in 10 R@1	2 in 10 R@2	3 in 10 R@3	5 in 10 R@5
MLP	0.6837	0.7548	0.3858	0.5343	0.6563	0.8060
MLP+kb	0.7123	0.8126	0.4840	0.6248	0.7361	0.8645
LSTM	0.7162	0.8332	0.5418	0.6957	0.7437	0.8863
LSTM+kb	0.7322	0.8679	0.6043	0.7535	0.8087	0.9119
Affinity Model	RNN	0.6850	0.7618	0.3963	0.5430	0.6732
	LSTM	0.7238	0.8542	0.5846	0.7214	0.7864
r-LSTM	0.7518	0.8892	0.6493	0.7853	0.8578	0.9324

Table 3: Experiment Results on Ubuntu Corpus.

from the raw corpus are different from that used in (Lowe et al., 2015), even though the extraction method is the same as Lowe et al. (2015). Thus the change of data distribution may influence the performance of this model; Second, the average length of contexts (by word) in the dataset used in this paper is 57.81, by contrast, in (Lowe et al., 2015) the average context length is 108.93. As mentioned in section 1, semantic clues in context are essential for response selection, so more context could offer more semantic information. We will open the dataset to other researchers for further comparison.

Ignoring RNN Affinity model, the baseline methods shown in the tables can be classified into two categories: MLP based methods and LSTM based ones. It’s clear that LSTM based methods perform better than the MLP based approaches, which indicates the ability of capturing semantic clues plays significant role in our task, and obviously sequence models have congenital advantage. Heuristically, explicit or implicit semantic relationships between utterances exist in human-to-human conversations to keep the continuity of conversations. Due to the “memory cell”, LSTM can model the long-dependent semantic relationships effectively, so approaches based on LSTM make better understanding of conversations.

As shown in Table 2 and 3, methods taking account of background knowledge get 3%-7% improvements comparing with the ones without any external knowledge on the best response selection task. This phenomenon indicates that knowledge is one of the primary factors for conversation un-

derstanding. In detail, there are two aspects influencing the improvement of approaches incorporating domain knowledge: (1) The quality and quantity of knowledge. Benefited from the characteristics of loose structured knowledge, we can build a knowledge base with large amount of domain specified information and update it with little manual work. Even the quality of such loose structured knowledge is lower than the knowledge base built with much human efforts (like Knowledge Graph), our model can still effectively utilize it. (2) The fusion strategy between knowledge and utterances in conversation modeling. Our model r-LSTM obtains up to 2-4% improvement by all the metrics than “LSTM+kb” on both datasets. This demonstrates the importance of strategies for incorporating external knowledge in conversation modeling. Moreover, there are semantic gaps between loose structured knowledge and sentences, it’s troublesome for methods like “LSTM+kb”, to overcome these gaps for methods taking knowledge as an additional utterance directly.

The *Recall@3* of our model on Tieba and Ubuntu dataset is 75.27% and 85.78% respectively, that is, most best responses can be recalled in the top three candidate responses (ranked by confidence). This can be attributed to the recall mechanism of r-LSTM involving background knowledge as global memory. As the global input, external domain knowledge can be recalled by r-LSTM and influence the local memory at right moments in the whole conversation process. The introduction of background knowledge provide essential semantic hints thus enhances the ability to

context	I like celebrity-name madly! I like him as well! His performance on Spring Festival Gala is very real! His model on Spring Festival is very comical. But I still like him.		
LSTM+kb	r-LSTM	Label	Response
0.7959	0.8676	1	Me too. This uncle is really handsome !
0.8655	0.7555	0	I always feel that he is old-fashioned.
0.3531	0.1172	0	So embarrassed tonight.
context	Is celebrity-name still a super-star? Her acting is poor, isn't it? celebrity-name remains at TV acting level, her cinematic feeling is so weak.		
0.6772	0.9078	1	Comparing with actresses in the same period, her acting and appearance are both good.
0.9278	0.8217	0	I like celebrity-name-1 , the real actress .
0.5598	0.0419	0	I do not think so, because they just took it out for a walk.

Table 4: Samples of best response selection. Entities are anonymous during translation.

detect the semantic relevances between sentences.

Due to the limitation of quality and average length of contexts in Tieba, the models' overall performances are a little lower than those on the Ubuntu, as shown in Table 2 and 3. According to the distribution of turn numbers shown in Figure 5, there are few conversations with more than 5 turns in the Tieba corpus. Consequently, there are less history contents involving semantic clues for selecting response as mentioned in section 1.

4.3 Case Study

In order to further show the working mechanism in our model intuitively, we give two cases comparing our r-LSTM and LSTM+kb in Table 4. For better understanding, we translate both the contexts and the candidate responses into English from Chinese. As shown by Table 4, each case contains the context composed of three utterances and three candidate responses with labels and predicted scores. The label 1 indicates the true response to the given context. Scores represent the confidence of responses as the best one.

From the table it can be seen that r-LSTM gives highest score to best response, while "LSTM+kb" offers unsatisfied results. Both methods shown in 4 involved background knowledge, so this result is caused by the effectiveness of utilizing knowledge. In our r-LSTM architecture, background knowledge is considered as global memory and can be recalled in the conversation modeling process, and our Recall Gate could also build semantic relationships of utterances validly. Taking knowledge as input directly, LSTM+kb makes less use of knowledge in the modeling process because of the semantic gap between sentences and knowledge items. Furthermore, based on the ob-

servation of human conversations, it is reasonable to consider background knowledge as global signal of the neural network, instead of an additional input as LSTM+kb does. The large range of confidence scores given by r-LSTM also shows that our model does well on both selecting best responses and recognizing inappropriate responses.

In both the two samples, the second response also gets high confidence score, and LSTM+kb even takes it as the best one. Actually, such responses are general responses without any key points. Even though they can be taken to reply the queries, such general answers are not promising because less dialog turns are expected after them. Thus distinguishing general and best responses is of great value for chat-agents, and our model have potential in this scenario.

5 Conclusions and Future Work

In this paper, we proposed a deep learning architecture to incorporate loose structured knowledge for end-to-end conversation modeling. Our approach shows good potential on the context-oriented response selecting task.

The contributions of this paper can be summarized as follows: (1) By investigating the influence of domain knowledge on conversation understanding, we present a LSTM framework with a designed Recall gate to utilize knowledge smoothly and effectively. Transforming knowledge into global memory, the Recall gate enables LSTM to integrate global memory into sequential local memory to conduct enhanced conversation modeling. (2) To guarantee the flexibility of domain knowledge base for practical usage, this paper introduces the loose-structured knowledge base organized as "entity-attribute" pairs, which can be

extracted easily from raw corpora. The knowledge can be directly absorbed by the Recall gate after being triggered according to dialog contexts.

Our future study will be conducted along the following directions: first, we will continue exploring the working mechanism of global memory to refine our Recall gate; Second, the utilization of our model for open-domain conversation modeling will be investigated.

References

- [Ameixa et al.2014] David Ameixa, Luisa Coheur, Pedro Fialho, and Paulo Quaresma. 2014. Luke, i am your father: dealing with out-of-domain requests by using movies subtitles. In *Intelligent Virtual Agents*, pages 13–21. Springer.
- [Banchs and Li2012] Rafael E Banchs and Haizhou Li. 2012. Iris: a chat-oriented dialogue system based on the vector space model. In *Proceedings of the ACL 2012 System Demonstrations*, pages 37–42. Association for Computational Linguistics.
- [Cho et al.2014] Kyunghyun Cho, Bart van Merriënboer, Çaglar Gülçehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, Doha, Qatar; A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1724–1734.
- [Henderson et al.2013] Matthew Henderson, Blaise Thomson, and Steve Young. 2013. Deep neural network approach for the dialog state tracking challenge. In *Proceedings of the SIGDIAL 2013 Conference*, pages 467–471.
- [Hochreiter and Schmidhuber1997] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.
- [Hurford et al.2007] James R Hurford, Brendan Heasley, and Michael B Smith. 2007. *Semantics: a coursebook*. Cambridge University Press.
- [Ji et al.2014] Zongcheng Ji, Zhengdong Lu, and Hang Li. 2014. An information retrieval approach to short text conversation. *CoRR*, abs/1408.6988.
- [Jurcicek et al.2011] Filip Jurcicek, Simon Keizer, Milica Gašić, Francois Mairesse, Blaise Thomson, Kai Yu, and Steve Young. 2011. Real user evaluation of spoken dialogue systems using amazon mechanical turk. In *Proceedings of INTERSPEECH*, volume 11.
- [Kadlec et al.2015] Rudolf Kadlec, Martin Schmid, and Jan Kleindienst. 2015. Improved deep learning baselines for ubuntu corpus dialogs. *CoRR*, abs/1510.03753.
- [Lowe et al.2015] Ryan Lowe, Nissan Pow, Iulian Serban, and Joelle Pineau. 2015. The ubuntu dialogue corpus: A large dataset for research in unstructured multi-turn dialogue systems. *CoRR*, abs/1506.08909.
- [Mu and Yin2010] Yanhua Mu and YiXin Yin. 2010. Task-oriented spoken dialogue system for humanoid robot. In *Multimedia Technology (ICMT), 2010 International Conference on*, pages 1–4. IEEE.
- [Ovchinnikova2012] Ekaterina Ovchinnikova. 2012. Natural language understanding and world knowledge. In *Integration of World Knowledge for Natural Language Understanding*, pages 15–37. Springer.
- [Ribeiro et al.2015] Eugénio Ribeiro, Ricardo Ribeiro, and David Martins de Matos. 2015. The influence of context on dialogue act recognition. *arXiv preprint arXiv:1506.00839*.
- [Ritter et al.2011] Alan Ritter, Colin Cherry, and William B Dolan. 2011. Data-driven response generation in social media. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 583–593. Association for Computational Linguistics.
- [Serban et al.2015a] Iulian V Serban, Alessandro Sordoni, Yoshua Bengio, Aaron Courville, and Joelle Pineau. 2015a. Hierarchical neural network generative models for movie dialogues. *arXiv preprint arXiv:1507.04808*.
- [Serban et al.2015b] Iulian Vlad Serban, Ryan Lowe, Laurent Charlin, and Joelle Pineau. 2015b. A survey of available corpora for building data-driven dialogue systems. *CoRR*, abs/1512.05742.
- [Shang et al.2015] Lifeng Shang, Zhengdong Lu, and Hang Li. 2015. Neural responding machine for short-text conversation. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing, ACL 2015, July 26-31, Beijing, China, Volume 1: Long Papers*, pages 1577–1586.
- [Sordoni et al.2015] Alessandro Sordoni, Michel Galley, Michael Auli, Chris Brockett, Yangfeng Ji, Margaret Mitchell, Jian-Yun Nie, Jianfeng Gao, and Bill Dolan. 2015. A neural network approach to context-sensitive generation of conversational responses. *arXiv preprint arXiv:1506.06714*.
- [Sukhbaatar et al.2015] Sainbayar Sukhbaatar, Jason Weston, Rob Fergus, et al. 2015. End-to-end memory networks. In *Advances in Neural Information Processing Systems*, pages 2431–2439.
- [Sutskever et al.2014] Ilya Sutskever, Oriol Vinyals, and Quoc VV Le. 2014. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112.

[Vinyals and Le2015] Oriol Vinyals and Quoc V. Le. 2015. A neural conversational model. *CoRR*, abs/1506.05869.

[Williams et al.2013] Jason Williams, Antoine Raux, Deepak Ramachandran, and Alan Black. 2013. The dialog state tracking challenge. In *Proceedings of the SIGDIAL 2013 Conference*, pages 404–413.