

Supplemental Material

A Working Memory Model for Task-oriented Dialog Response Generation

Xiuyi Chen^{1,2,3}, Jiaming Xu^{1,2*} and Bo Xu^{1,2,3,4}

¹Institute of Automation, Chinese Academy of Sciences (CASIA). Beijing, China

²Research Center for Brain-inspired Intelligence, CASIA

³University of Chinese Academy of Sciences

⁴Center for Excellence in Brain Science and Intelligence Technology, CAS. China

{chenxiuyi2017, jiaming.xu, xubo}@ia.ac.cn

Task	bAbI Dialog Dataset					DSTC2
	1	2	3	4	5	
Train Dialogs			1000			1618
Val Dialogs			1000			500
Test Dialogs			1000			1117
Vocabulary Size			3747			1229
Avg. User turns	4	6.5	6.4	3.5	12.9	6.8
Avg. Sys. turns	6	9.5	9.9	3.5	18.4	9.2
Avg. KB entries	0	0	24	7	23.7	39.3
Avg. Sys. words	6.3	6.2	7.2	5.7	6.5	10.2
Pointer Ratio	24%	53%	46%	19%	60%	48%
Pointer Hist. Ratio	24%	53%	43%	11%	56%	41%
Pointer KB Ratio	0%	0%	3%	8%	3%	7%

Table 2: Dataset statistics for the bAbI and DSTC2 dialog datasets.

1 Hyper-parameters

The hyper-parameters for all models are optimized on the validation sets; values for the best performing models are given in the Table 1. The learning rate is equal to 0.001 with a decay rate of 0.5 and the random seed is fixed to “666”.

2 Dataset

We conduct experiments on the simulated bAbI Dialogue dataset and the Dialog State Tracking Challenge 2 (DSTC2). And their statistics are reported in Table 2.

The bAbI dialog is synthetically generated by a simulator in the context of restaurant reservation. This dataset consists of five tasks, where tasks 1 to 4 are issuing API calls, updating API calls, displaying options, providing extra information and task 5 combines tasks 1-4 to generate full dialogs. An extra benefit of this dataset is that it provides two test sets for each task: one is generated with the same KB as the train set, and the other is termed as the OOV test set with a different KB. Hence, the OOV test set is much harder than the other set, and needs more reasoning ability.

DSTC2 is a dataset on restaurant reservation, extracted from real human-bot dialogs. In order to

Symbol	Definition
x_i or y_i	a token in the dialog history or system response
$\$$	a special token used as a sentinel (Madotto et al., 2018)
X	$X = \{x_1, \dots, x_n, \$\}$, the dialog history
Y	$Y = \{y_1, \dots, y_m\}$, the expected response
b_i	one KB tuple, actually the corresponding entity
B	$B = \{b_1, \dots, b_l, \$\}$, the KB tuples
$PTRE$	$PTRE = \{ptr_{E,1}, \dots, ptr_{E,m}\}$, dialog pointer index set. supervised information for copying words in dialog history
$PTRS$	$PTRS = \{ptr_{S,1}, \dots, ptr_{S,m}\}$, KB pointer index set. supervised information for copying entities in KB tuples

Table 3: Notation Table.

evaluate end-to-end systems, we here use the refined version from Bordes et al. (2017), which merely uses the raw text of the dialogs without state annotations. Furthermore, this dataset, as derived from a real-word system, presents linguistic diversity and conversational abilities.

3 Data Pre-processing

Given a dialog between a user (u) and a system (s), we represent the dialog utterances as $\{(u_1, s_1), \dots, (u_j, s_j), \dots, (u_t, s_t)\}$ where t denotes the number of turns in the dialog and $j = 1, \dots, t$. Note that there is a background KB information organized in the (subject, relation, object) format. Hence, we define each KB tuple as b_i and the whole KB information containing l tuples as $B = \{b_1, \dots, b_l, \$\}$, where $\$$ is a special token used as a sentinel. Having the dialog history and world knowledge in mind, humans can easily respond with the “correct” utterance, and we define the expected system response $Y = \{y_1, \dots, y_m\}$ as the sequence of words in the utterance s_j . For the current dialog turn j , we define the dialog history as a sequence of tokens $X = \{x_1, \dots, x_n, \$\}$ of $\{u_1, s_1, \dots, s_{j-1}, u_j\}$. Different from most previous works, we do not mix KB tuples with the dialog history. Hence, we can separately store the dialog history and KB information, and the two sentinels are used as a hard gate to control which

* Corresponding Author

Model-Task	T1	T2	T3	T4	T5	DSTC2
WMM2Seq	128-0.1-16	128-0.1-8	128-0.2-64	256-0.3-128	128-0.1-128	256-0.1-8
WMM2Seq+CNN	128-0.1-8	256-0.3-16	128-0.2-32	256-0.3-16	128-0.1-128	256-0.3-8
WMM2Seq+biGRU	128-0.2-32	256-0.3-64	128-0.1-32	256-0.3-64	256-0.2-16	64-0.2-16
WMM2Seq+biGRU (H1)	128-0.1-8	128-0.2-64	256-0.2-32	128-0.2-16	256-0.2-64	256-0.1-64

Table 1: Selected hyper-parameters in each datasets for different models. Each cell in the table is in the (HDD-DR-BSZ) format.

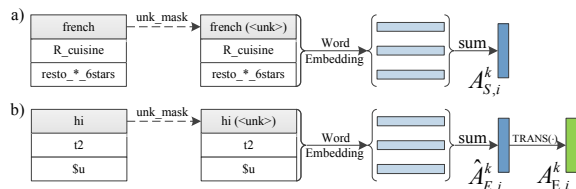


Figure 1: The Contents in the Episodic and Semantic Memories. We store the word representation of dialog history and KB tuples, separately and differently.

distribution to use at each decoding step as the paper illustrated.

To provide the model with an accurate guidance of how to activate the long-term memories, we follow Madotto et al. (2018) and further define two pointer index sets ¹ ($PTR_E = \{ptr_{E,1}, \dots, ptr_{E,m}\}$ and $PTR_S = \{ptr_{S,1}, \dots, ptr_{S,m}\}$) as the supervised information for copying words of dialog history or entities in KB information. We define each pointer index set as follows:

$$ptr_{E,i} = \begin{cases} \max(z) & \text{if } \exists z \text{ s.t. } y_i = u_z \\ n + 1 & \text{otherwise} \end{cases}, \quad (1)$$

$$ptr_{S,i} = \begin{cases} \max(z) & \text{if } \exists z \text{ s.t. } y_i = u_z \\ l + 1 & \text{otherwise} \end{cases}, \quad (2)$$

where $u_z \in X$ is the dialog history in Eq. 1 or the KB tuples in Eq. 2. And $n + 1$ and $l + 1$ are both the sentinel position indexes because the sentinel is at the end of the dialog history (with length n) or KB information (with length l). The idea behind Eq. 1 or Eq. 2 is that we can obtain the positions of where to copy by matching the target text with the dialog history or KB information. The symbols are defined in Table 3.

4 Memory Content

We store word-level content X into the MemNN encoder and E-MemNN. Furthermore, we incorporate additional temporal information and speaker

¹Note, all variables belonging to the episodic memory are with subscript E, and semantic memory are with subscript S.

information into each token of X to capture the sequential dependencies. For example, “hello t1 \$u” means that a user speaks the token “hello” at the turn $i = 1$ in the dialog. As shown in the lower part of Figure 1, we first randomly mask some tokens in X to simulate the OOV situation in dialog history; then we sum the three word embeddings to obtain the token representation; finally, we apply $TRANS(\cdot)$ function to get the context-aware token representation which is stored into the MemNN encoder and E-MemNN. We adopt a (subject, relation, object) representation of KB information and use S-MemNN to memorize each KB tuple as the token in X without context-aware transformation (see the upper part of Figure 1). Please note that we only load the related KB tuples into the semantic memory.

5 Human Evaluation

Systems are evaluated in terms of appropriateness and humanlikeness on a scale from 1 to 5, and the guidance of human evaluation scores are listed in Table 4. Furthermore, scores about the quality of the generated responses are given by two human subjects and the number of each score are shown in Figure 2.

References

- Antoine Bordes, Y-Lan Boureau, and Jason Weston. 2017. Learning end-to-end goal-oriented dialog. In *International Conference on Learning Representations*.
- Andrea Madotto, Chien-Sheng Wu, and Pascale Fung. 2018. Mem2seq: Effectively incorporating knowledge bases into end-to-end task-oriented dialog systems. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1468–1478. Association for Computational Linguistics.

Appropriateness	Quality
1	Wrong grammar, wrong logic, wrong dialogue flow, and wrong entity provided
2	Poor grammar, logic and entity provided
3	Noticeable mistakes about grammar or logic or entity provided but acceptable
4	Correct dialogue flow, logic and grammar but has slight mistakes in entity provided
5	Correct grammar, correct logic, correct dialogue flow, and correct entity provided
Humanlikeness	Quality
1	The utterance is 0% like what a person will say
2	The utterance is 25% like what a person will say
3	The utterance is 50% like what a person will say
4	The utterance is 75% like what a person will say
5	The utterance is 100% like what a person will say

Table 4: The quality of a response according to the appropriateness and humanlikeness.

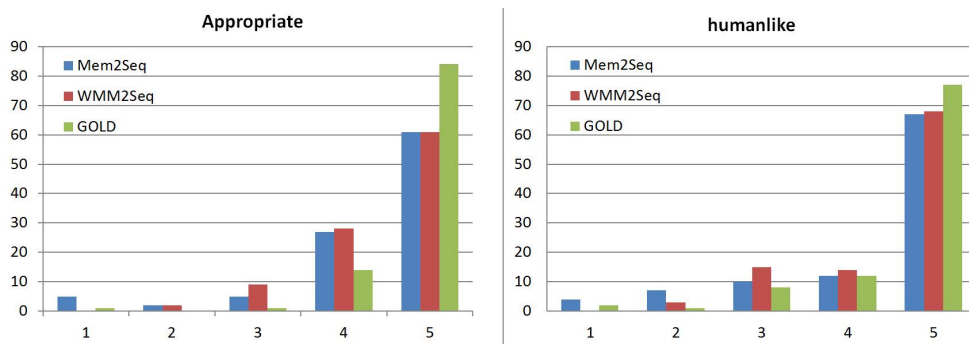


Figure 2: Appropriateness and human-likeness scores according to 100 dialogue samples from DSTC2 test set.