# End-to-end Gated Self-attentive Memory Network for Dialog Response Selection

Shuo Sun<sup>1\*</sup>, Yik-Cheung Tam<sup>2\*</sup>, Jie Cao<sup>3\*</sup>, Canxiang Yan<sup>2</sup>, Zuohui Fu<sup>4</sup>, Cheng Niu<sup>2</sup>, Jie Zhou<sup>2</sup>

<sup>1</sup> Center for Language and Speech Processing, Johns Hopkins University
<sup>2</sup> WeChat AI - Pattern Recognition Center, Tencent Inc.
<sup>3</sup> School of Computing, University of Utah
<sup>4</sup> Department of Computer Science, Rutgers University



\* Equal contribution. Work done when Shuo Sun, Jie Cao and Zuohui Fu were interns at WeChat AI - Pattern Recognition Center, Tencent inc.

## INTRODUCTION

- Gated Self-attentive Memory Network encodes dialog history and external domain knowledge in an end-to-end trainable manner.
- Novelty is that each utterance in the memory is enhanced with self-attention building the connection between dialog history and external domain.
- Our approach ranks at the second place for both student advising and Ubuntu subtasks integrated with external domain knowledge.



- XGBoost is employed to train a second stage meta classifier.
- Scores from many GSMN models with various hyperparameters are used as features.
- Training instances with mean prediction scores outside the range [0.001,0.95] are filtered out.



### **EXPERIMENTAL RESULTS**



#### **Memory Attention**

- Attention mechanism is used to "retrieve" relevant memory vectors.
- Output is a weighted sum of the memory vectors.

#### **Gated Memory Attention**

- Input and the retrieved memories are combined via a gating mechanism.
- Regulates the degree of enhancement.
- Prevents information overload.



### GATED SELF-ATTENTIVE MEMORY NETWORK



Model	<b>R@1</b>	<b>R@10</b>	<b>R@50</b>	MRF
Dual-LSTM baseline	0.062	0.296	0.728	N/A
HGRU baseline	0.164	0.632	0.922	0.299
SMN w/ 1 hop	0.218	0.642	0.956	0.337
2 hops	0.198	0.620	0.938	0.320
3 hops	0.206	0.648	0.942	0.333
GSMN w/ 1 hop	0.220	0.632	0.954	0.343
2 hops	0.214	0.644	0.960	0.338
3 hops	0.214	0.628	0.956	0.335
1 hop + EK	0.220	0.644	0.956	0.343
2  hops + EK	0.224	0.654	0.944	0.354

Baseline and GSMN results on the Flex advising

dev set. EK denotes external knowledge.

Model	<b>R@1</b>	R@10	R@50	MRR
Dual-LSTM baseline	0.083	0.360	0.804	N/A
SMN w/ 1 hop	0.326	0.671	0.952	0.445
SMN w/ 2 hop	0.337	0.686	0.956	0.455
GSMN w/ 1 hop	0.379	0.733	0.973	0.497
2 hops	0.389	0.755	0.972	0.508
3 hops	0.398	0.761	0.976	0.515

Baseline and GSMN results on the Ubuntu dev set. (Single model)

# **EFFECTIVENESS OF SELF-ATTENTION AT UTTERANCE LEVEL**

Model		<b>R@1</b>		R@5		R@10			MRR			
	LU	AU	Gain	LU	AU	Gain	LU	AU	Gain	LU	AU	Gain
1 hop	0.194	0.220	▲0.026	0.592	0.643	▲0.051	0.910	0.954	▲0.044	0.317	0.343	▲0.026
2 hops	0.208	0.214	▲0.006	0.596	0.644	▲0.048	0.920	0.960	▲0.040	0.326	0.338	▲0.012
3 hops	0.208	0.214	▲0.006	0.616	0.628	▲ $0.012$	0.928	0.956	▲0.028	0.329	0.335	▲0.006
1 hop + EK	0.172	0.220	▲0.048	0.578	0.644	▲0.066	0.910	0.956	▲0.046	0.294	0.343	▲0.049
2  hops + EK	0.186	0.224	▲0.038	0.602	0.654	▲0.052	0.932	0.944	▲0.012	0.319	0.354	

LU: restricted the inputs to only the last utterance of a dialog history. AU: use almost all utterances of a dialogue history (last 6 utterances)

- Significant performance gains across all evaluation metrics and model settings were observed.
- This confirms that modeling co-references among dialog utterances via self-attention was effective.

### **OFFICIAL EVALUATION RESULTS**

Measure	Ubuntu	Advising	Advising
		- Case 1	- Case 2
Recall@1	0.475	0.494	0.18
Recall@10	0.814	0.85	0.562
Recall@50	0.978	0.98	0.94
MRR	0.595	0.6078	0.3069

Measure	Ubuntu	Advising	Advising
		- Case 1	- Case 2
Recall@1	0.504	0.538	0.178
Recall@10	0.827	0.864*	0.608
Recall@50	0.98	0.986	0.944
MRR	0.6172	0.6455*	0 3149

- Short-term memories and long-term memories are stacked sequentially.
- Process is repeated in a multi-hop fashion.
- Weights are not shared across hops.
- Memory enhanced dialog utterances are summed into a single vector.



•



- Gated Self-attentive Memory Network effectively integrates external knowledge and dialog history in an end-to-end fashion.
- For future work, we believe that improving the encoding power of dialog history and external domain knowledge as well as their interaction will be crucial for further performance improvement.