

**YUN-NUNG (VIVIAN) CHEN**

**[HTTP://VIVIANCHEN.IDV.TW](http://vivianchen.idv.tw)**

**HAKKANI-TUR, TUR, GAO, DENG**

Microsoft  
**Research**

**End-to-End Memory Networks with Knowledge Carryover  
for Multi-Turn Spoken Language Understanding**

**INTERSPEECH 2016**



**臺灣大學**

**National Taiwan University**



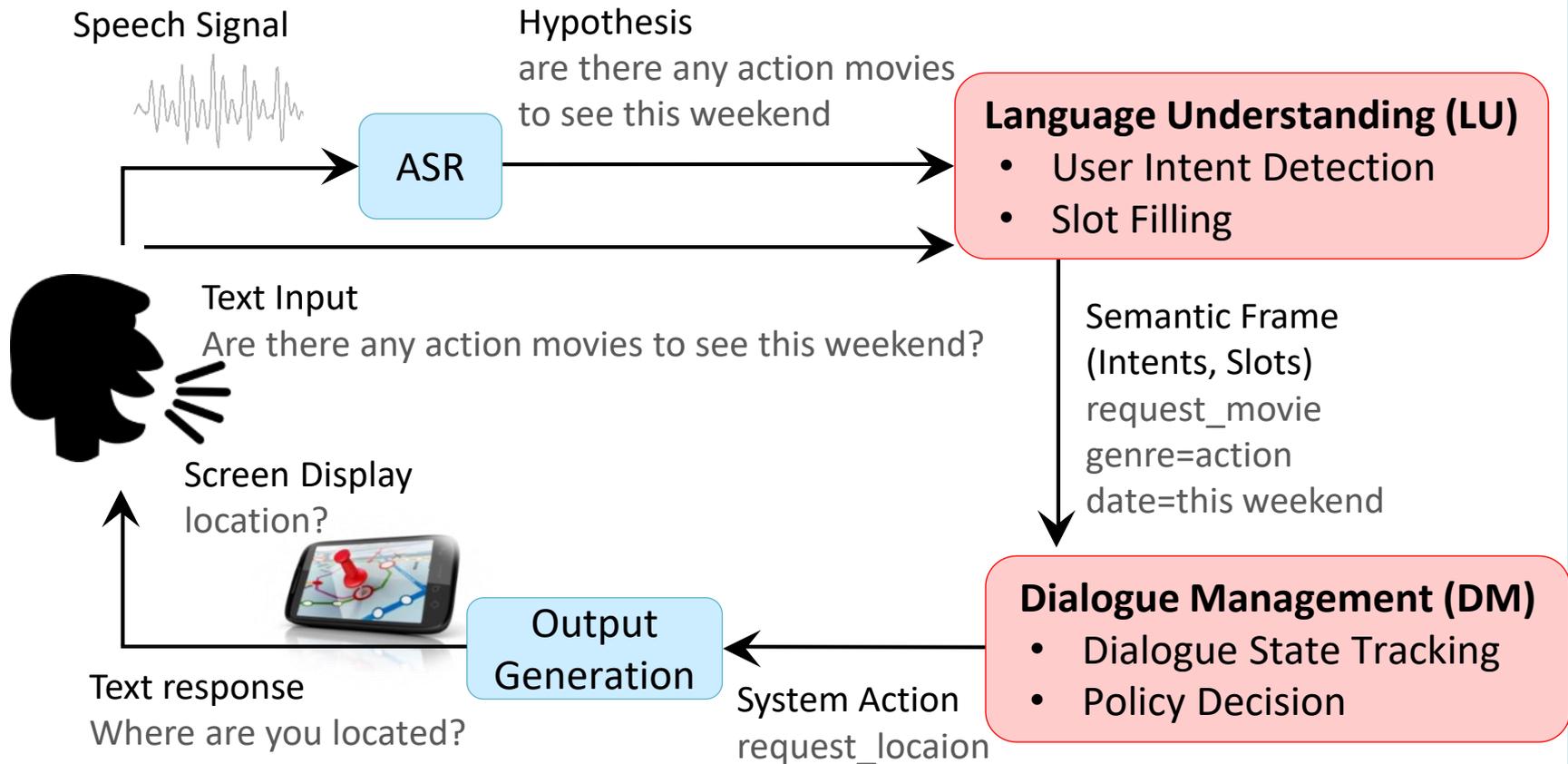
# Outline

- Introduction
    - Spoken Dialogue System
    - Spoken/Natural Language Understanding (SLU/NLU)
  - Contextual Spoken Language Understanding
    - Model Architecture
    - End-to-End Training
  - Experiments
  - Conclusion & Future Work
- 

# Outline

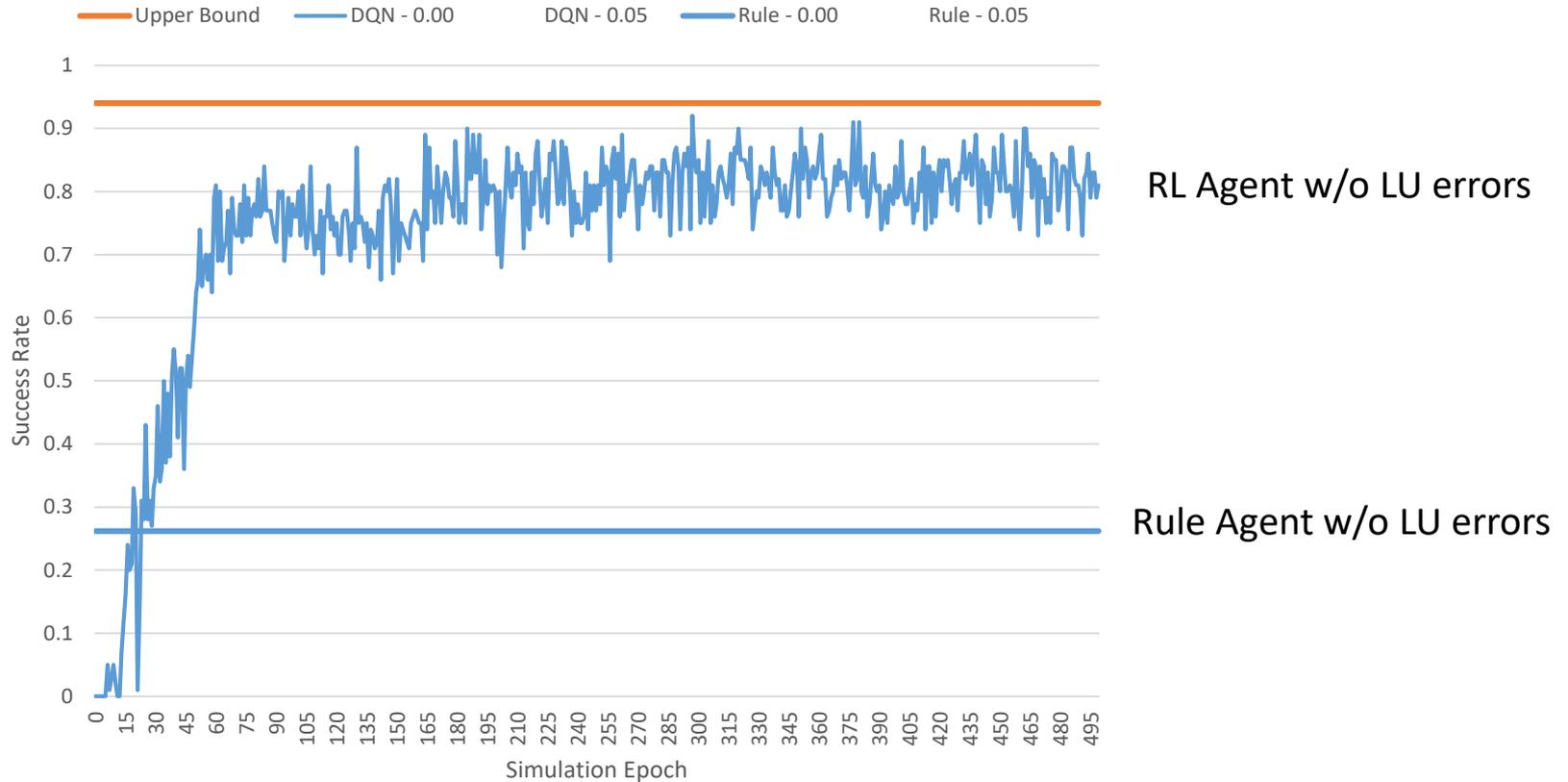
- Introduction
  - Spoken Dialogue System
  - Spoken/Natural Language Understanding (SLU/NLU)
  - Contextual Spoken Language Understanding
  - Model Architecture
  - End-to-End Training
  - Experiments
  - Conclusion & Future Work

# Dialogue System Pipeline



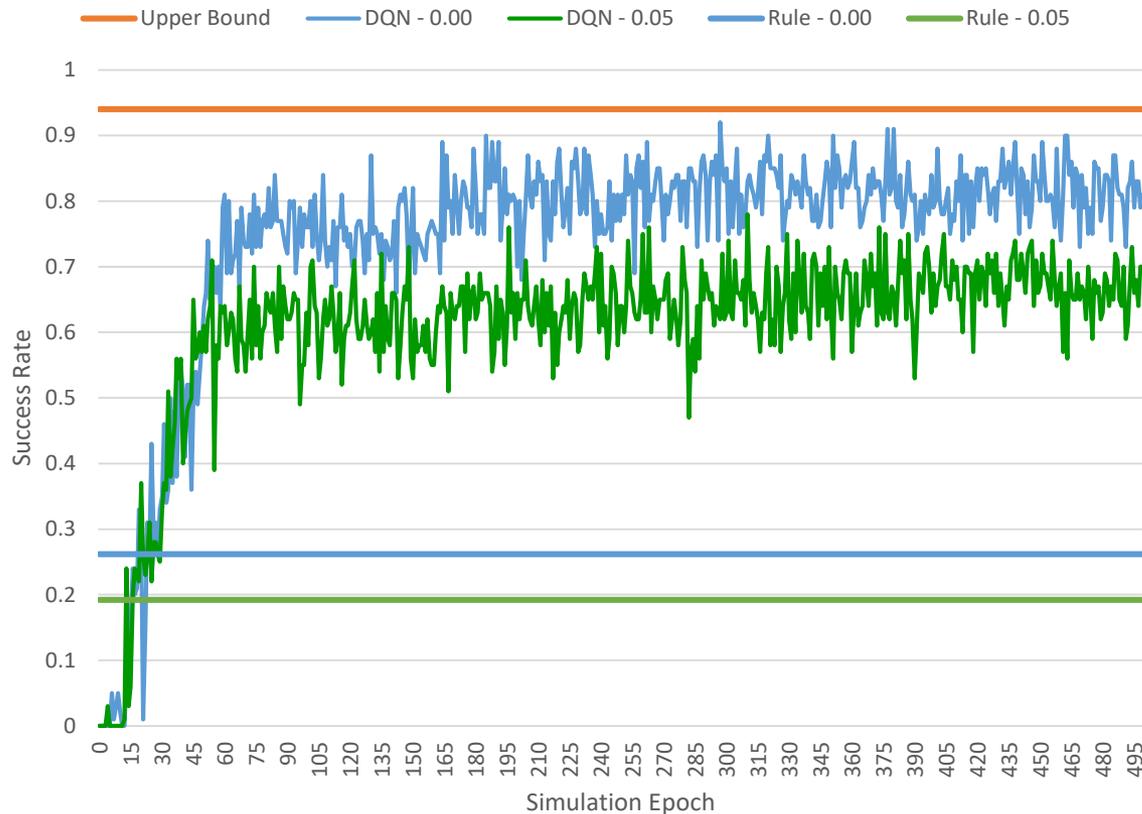
# LU Importance

## Learning Curve of System Performance



# LU Importance

Learning Curve of System Performance

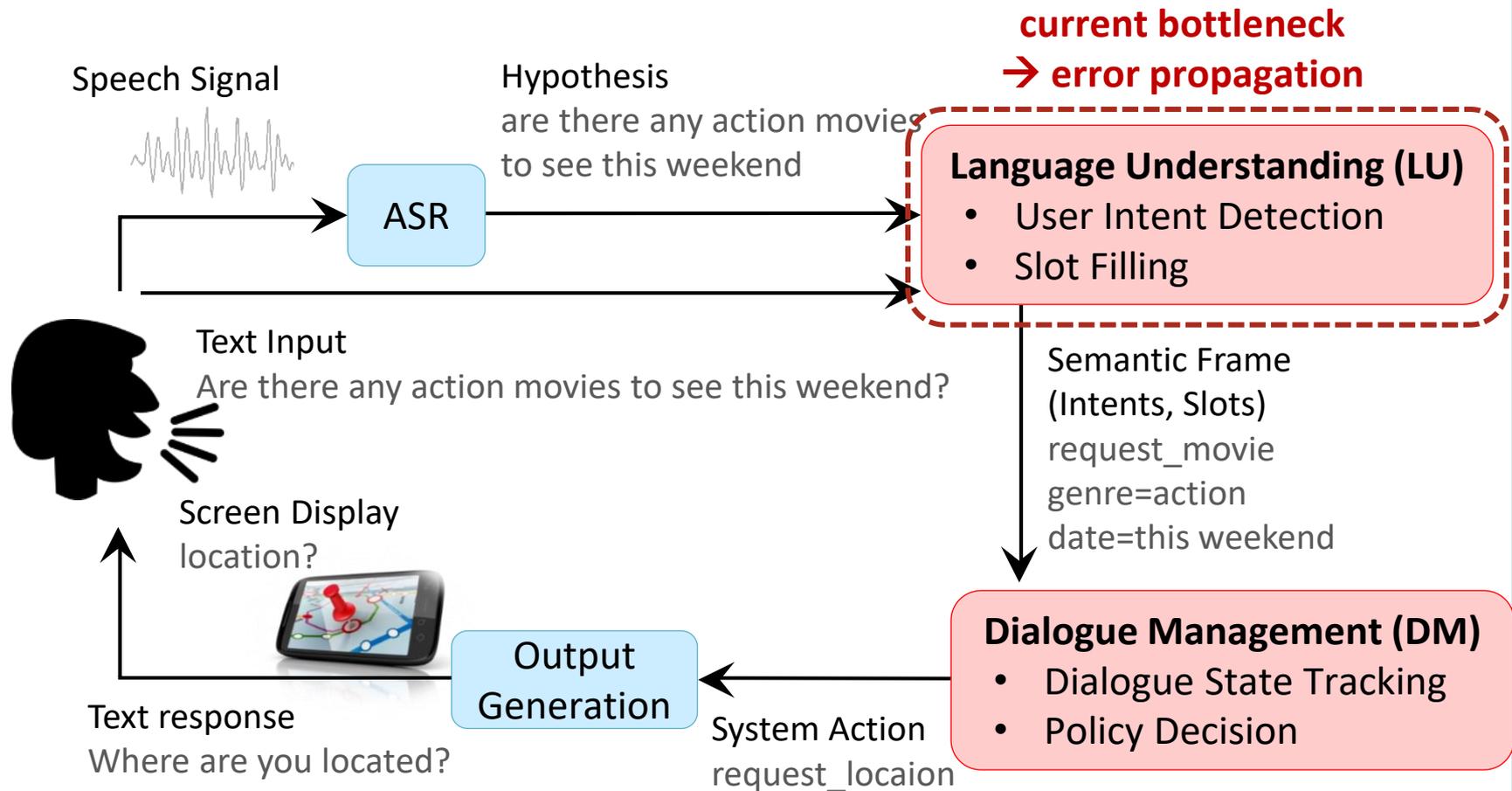


RL Agent w/o LU errors  
 RL Agent w/ 5% LU errors  
 >5% performance drop

Rule Agent w/o LU errors  
 Rule Agent w/ 5% LU errors

The system performance is sensitive to LU errors, for both rule-based and reinforcement learning agents.

# Dialogue System Pipeline



SLU usually focuses on understanding single-turn utterances

The understanding result is usually influenced by  
1) local observations 2) global knowledge.

# Spoken Language Understanding

Domain Identification → Intent Prediction → Slot Filling

$D$  communication

$I$  send\_email



$U$  just sent email to bob about fishing this weekend

$S$  O O O O ↓ O ↓ ↓ ↓

B-contact\_name B-subject I-subject I-subject

→ send\_email(contact\_name="bob", subject="fishing this weekend")

Single Turn

$U_1$  send email to bob

$S_1$  B-contact\_name

→ send\_email(contact\_name="bob")

$U_2$  are we going to fish this weekend

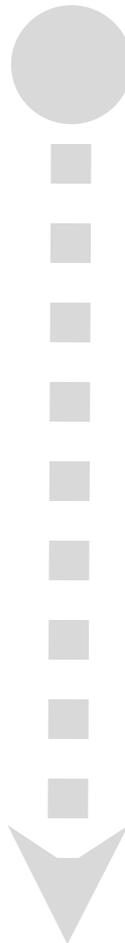
$S_2$  B-message ↓ I-message ↓ I-message ↓ I-message ↓ I-message

I-message I-message I-message

→ send\_email(message="are we going to fish this weekend")

Multi-Turn

# Outline

- Introduction
    - Spoken Dialogue System
    - Spoken/Natural Language Understanding (SLU/NLU)
  - **Contextual Spoken Language Understanding**
    - Model Architecture
    - End-to-End Training
  - Experiments
  - Conclusion & Future Work
- 

# MODEL ARCHITECTURE

## 1. Sentence Encoding

$$m_i = \text{RNN}_{\text{mem}}(x_i)$$

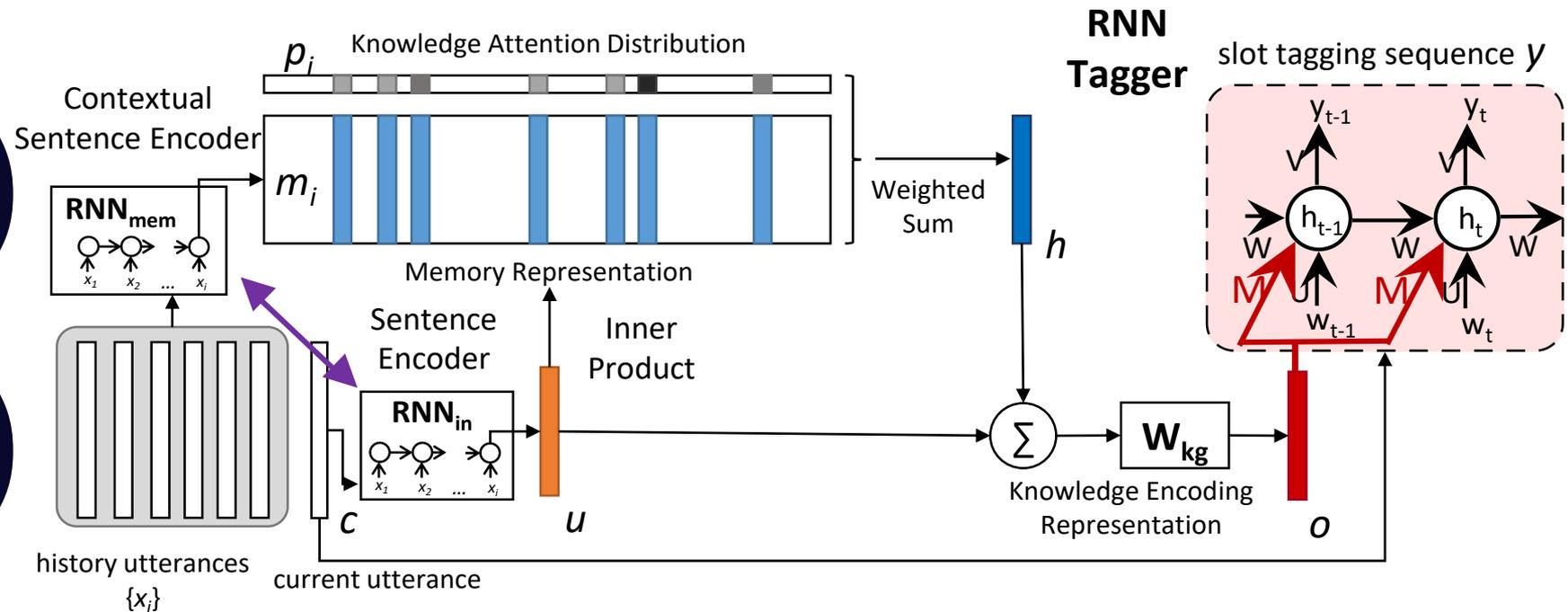
$$u = \text{RNN}_{\text{in}}(c)$$

## 2. Knowledge Attention

$$p_i = \text{softmax}(u^T m_i)$$

## 3. Knowledge Encoding

$$h = \sum_i p_i m_i \quad o = W_{\text{kg}}(h + u)$$



Idea: additionally incorporating contextual knowledge during slot tagging

# MODEL ARCHITECTURE

## 1. Sentence Encoding

$$m_i = \text{RNN}_{\text{mem}}(x_i)$$

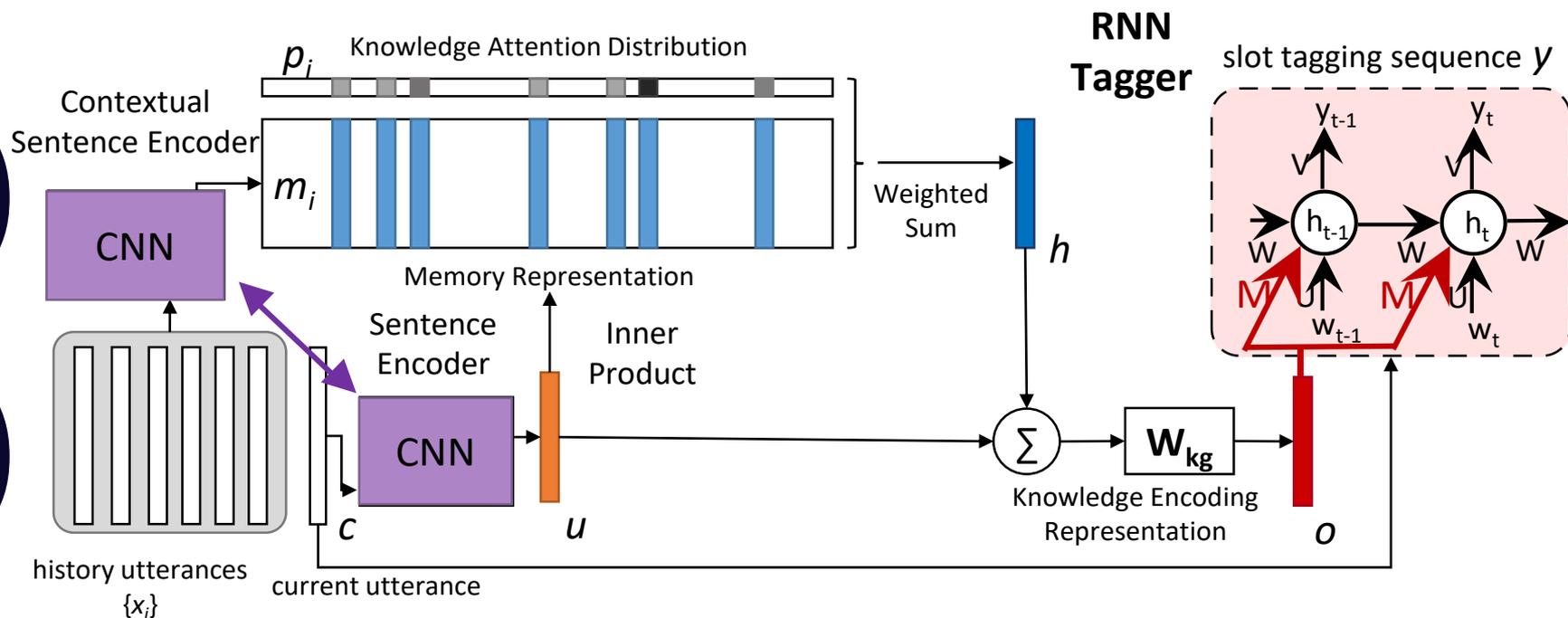
$$u = \text{RNN}_{\text{in}}(c)$$

## 2. Knowledge Attention

$$p_i = \text{softmax}(u^T m_i)$$

## 3. Knowledge Encoding

$$h = \sum_i p_i m_i \quad o = W_{\text{kg}}(h + u)$$



Idea: additionally incorporating contextual knowledge during slot tagging

# END-TO-END TRAINING

- Tagging Objective

$$\mathbf{y} = \text{RNN}(o, \mathbf{c})$$

slot tag sequence

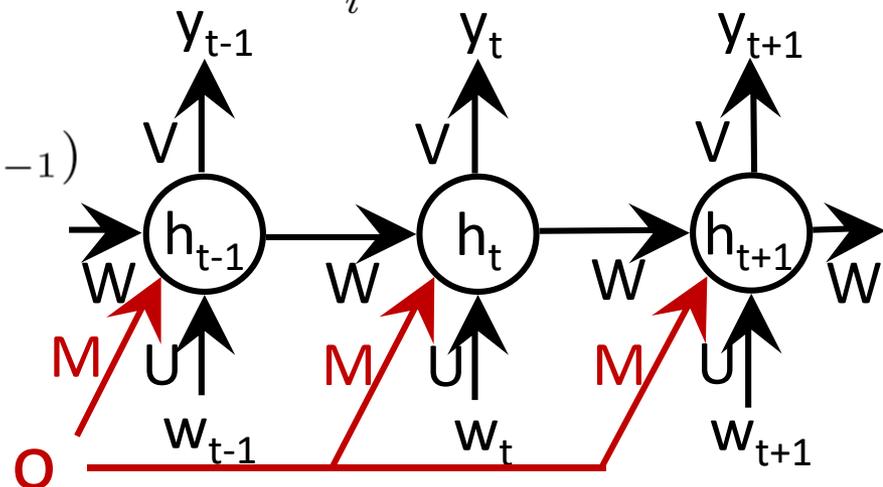
contextual utterances & current utterance

$$p(\mathbf{y} | \mathbf{c}) = p(\mathbf{y} | w_1, \dots, w_T) = \prod_i p(y_i | w_1, \dots, w_i).$$

- RNN Tagger

$$h_t = \phi(Mo + Ww_t + Uh_{t-1})$$

$$\hat{y}_t = \text{softmax}(Vh_t)$$



Automatically figure out the attention distribution without explicit supervision

# Outline

- Introduction
    - Spoken Dialogue System
    - Spoken/Natural Language Understanding (SLU/NLU)
    - Contextual Spoken Language Understanding
      - Model Architecture
      - End-to-End Training
  - Experiments
  - Conclusion & Future Work
- 

# EXPERIMENTS



- Dataset: Cortana communication session data
  - GRU for all RNN
  - adam optimizer
  - embedding dim=150
  - hidden unit=100
  - dropout=0.5

Model	Training Set	Knowledge Encoding	Sentence Encoder	First Turn	Other	Overall
RNN Tagger	single-turn	x	x	<b>60.6</b>	16.2	25.5

The model trained on single-turn data performs worse for non-first turns due to mismatched training data

# EXPERIMENTS



- Dataset: Cortana communication session data
  - GRU for all RNN
  - adam optimizer
  - embedding dim=150
  - hidden unit=100
  - dropout=0.5

Model	Training Set	Knowledge Encoding	Sentence Encoder	First Turn	Other	Overall
RNN Tagger	single-turn	x	x	<b>60.6</b>	16.2	25.5
	multi-turn	x	x	55.9	<b>45.7</b>	<b>47.4</b>

Treating multi-turn data as single-turn for training performs reasonable

# EXPERIMENTS



- Dataset: Cortana communication session data
  - GRU for all RNN
  - adam optimizer
  - embedding dim=150
  - hidden unit=100
  - dropout=0.5

Model	Training Set	Knowledge Encoding	Sentence Encoder	First Turn	Other	Overall
RNN Tagger	single-turn	x	x	60.6	16.2	25.5
	multi-turn	x	x	55.9	45.7	47.4
Encoder-Tagger	multi-turn	current utt (c)	RNN	57.6	56.0	56.3
	multi-turn	history + current (x, c)	RNN	<b>69.9</b>	<b>60.8</b>	<b>62.5</b>

Encoding current and history utterances improves the performance but increases the training time

# EXPERIMENTS

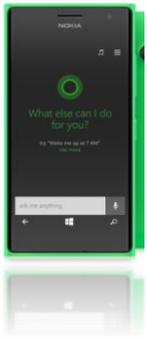


- Dataset: Cortana communication session data
  - GRU for all RNN
  - adam optimizer
  - embedding dim=150
  - hidden unit=100
  - dropout=0.5

Model	Training Set	Knowledge Encoding	Sentence Encoder	First Turn	Other	Overall
RNN Tagger	single-turn	x	x	60.6	16.2	25.5
	multi-turn	x	x	55.9	45.7	47.4
Encoder-Tagger	multi-turn	current utt (c)	RNN	57.6	56.0	56.3
	multi-turn	history + current (x, c)	RNN	69.9	60.8	62.5
Proposed	multi-turn	history + current (x, c)	RNN	<b>73.2</b>	<b>65.7</b>	<b>67.1</b>

Applying memory networks significantly outperforms all approaches with much less training time

# EXPERIMENTS



- Dataset: Cortana communication session data
  - GRU for all RNN
  - adam optimizer
  - embedding dim=150
  - hidden unit=100
  - dropout=0.5

Model	Training Set	Knowledge Encoding	Sentence Encoder	First Turn	Other	Overall
RNN Tagger	single-turn	x	x	60.6	16.2	25.5
	multi-turn	x	x	55.9	45.7	47.4
Encoder-Tagger	multi-turn	current utt (c)	RNN	57.6	56.0	56.3
	multi-turn	history + current (x, c)	RNN	69.9	60.8	62.5
Proposed	multi-turn	history + current (x, c)	RNN	<b>73.2</b>	<b>65.7</b>	<b>67.1</b>
	multi-turn	history + current (x, c)	CNN	<b>73.8</b>	<b>66.5</b>	<b>68.0</b>

CNN produces comparable results for sentence encoding with shorter training time

# Outline

- Introduction
    - Spoken Dialogue System
    - Spoken/Natural Language Understanding (SLU/NLU)
  - Contextual Spoken Language Understanding
    - Model Architecture
    - End-to-End Training
  - Experiments
  - Conclusion & Future Work
- 

# Conclusion

- The proposed *end-to-end* memory networks **store contextual knowledge**, which can be exploited dynamically based on an **attention model** for **manipulating knowledge carryover** for multi-turn understanding
- The end-to-end model performs the **tagging** task instead of **classification**
- The experiments show the *feasibility* and *robustness* of modeling knowledge carryover through memory networks

# Future Work

- Leveraging not only **local observation** but also **global knowledge** for better language understanding
  - Syntax or semantics can serve as global knowledge to guide the understanding model
  - “Knowledge as a Teacher: Knowledge-Guided Structural Attention Networks,” arXiv preprint [arXiv: 1609.03286](https://arxiv.org/abs/1609.03286)

# Q & A

THANKS FOR YOUR ATTENTION!

The code will be available at  
<https://github.com/yvchen/ContextualSLU>

