

Attention Based Joint Model with Negative Sampling for New Slot Values Recognition

By: Mulan Hou

houmulan@bupt.edu.cn



CONTENTS

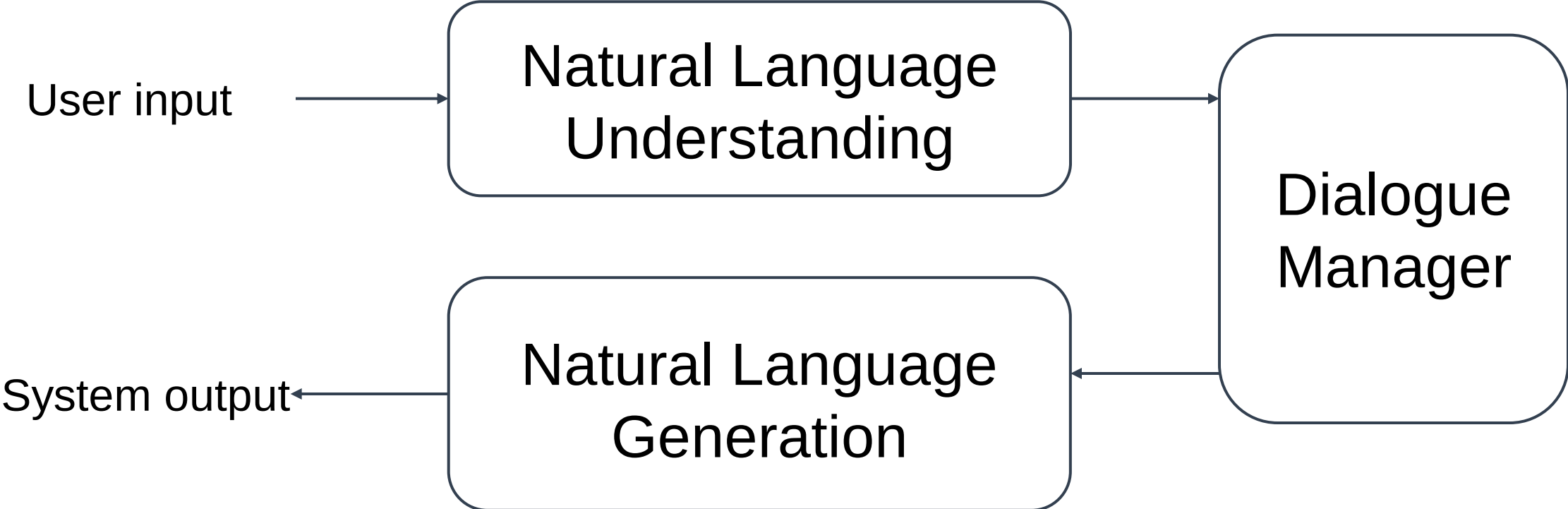
- 1 Introduction
- 2 Related work
- 3 Motivation
- 4 Proposed model
- 5 Experiments
- 6 Conclusion

CHAPTE
R

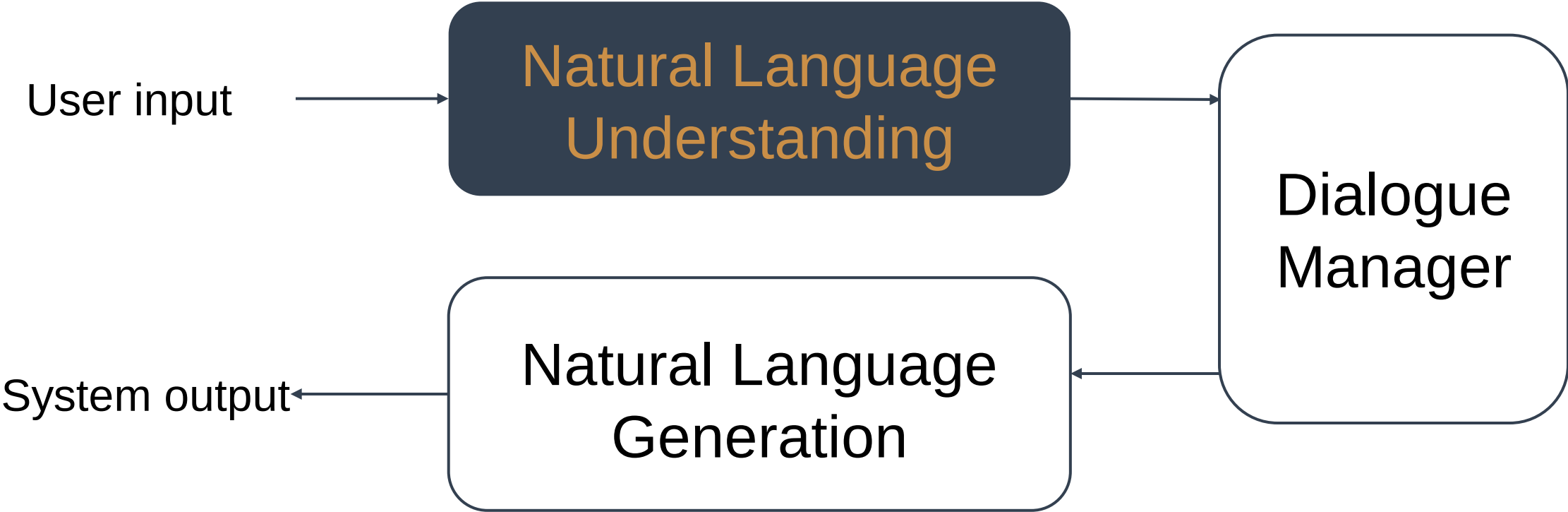
1

Introduction

Introduction



Introduction



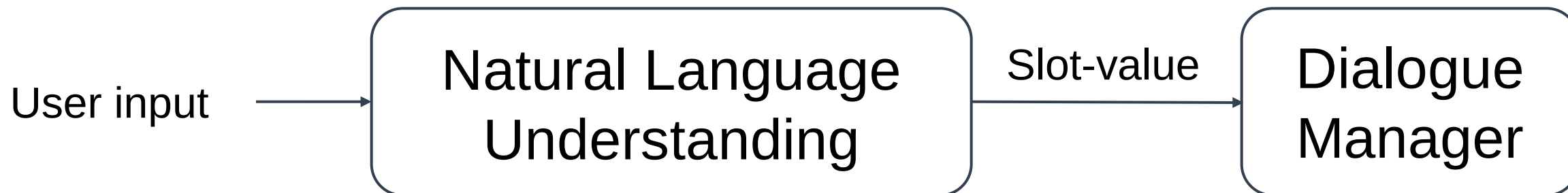
Introduction



Introduction

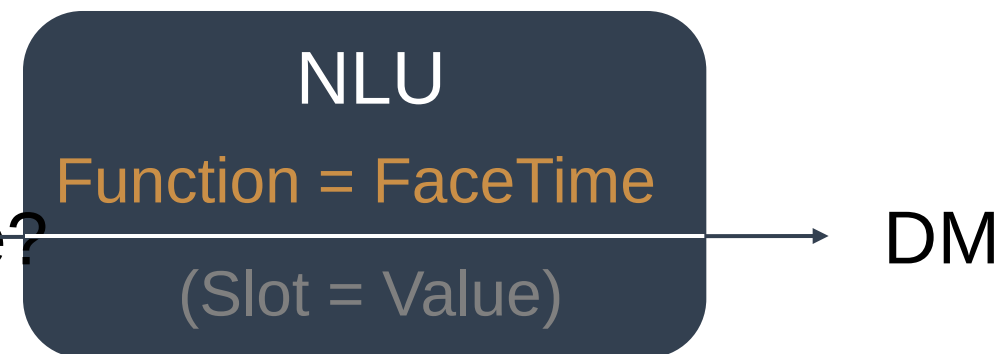


Introduction

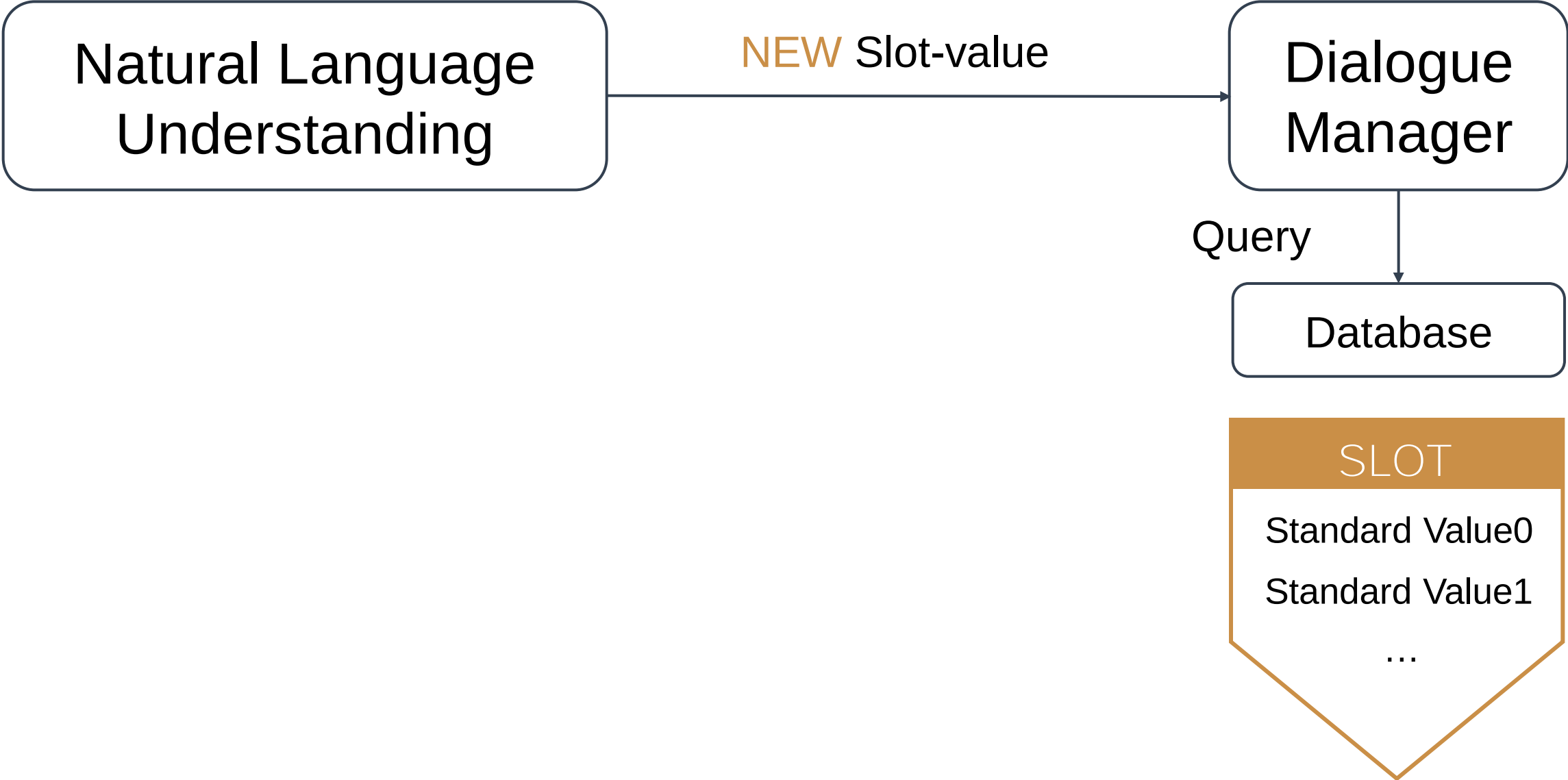


E.g

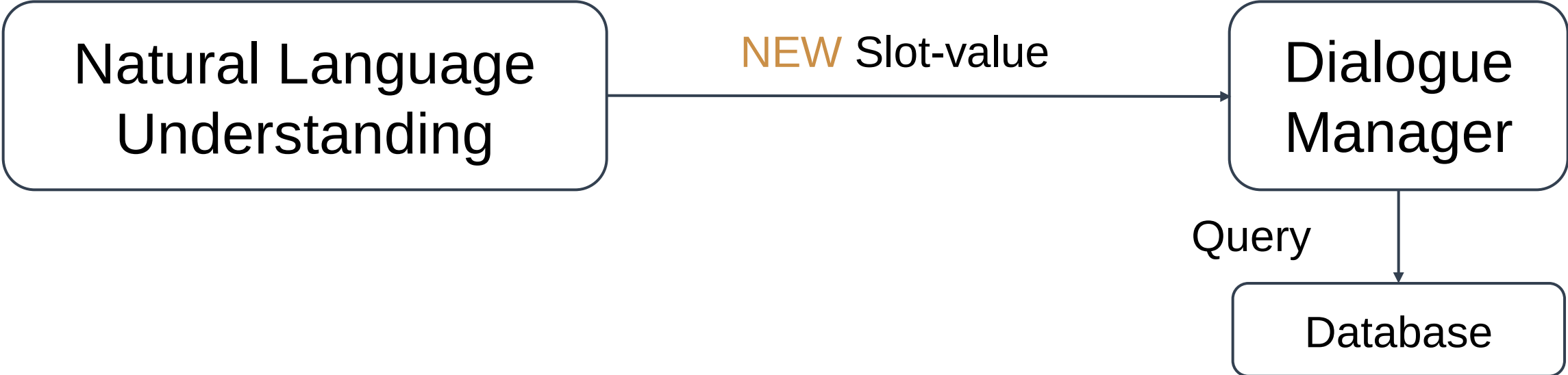
Can I **have a video chat** on my phone?



Introduction

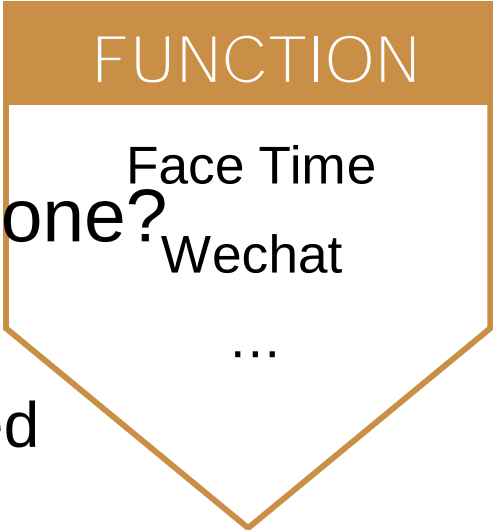
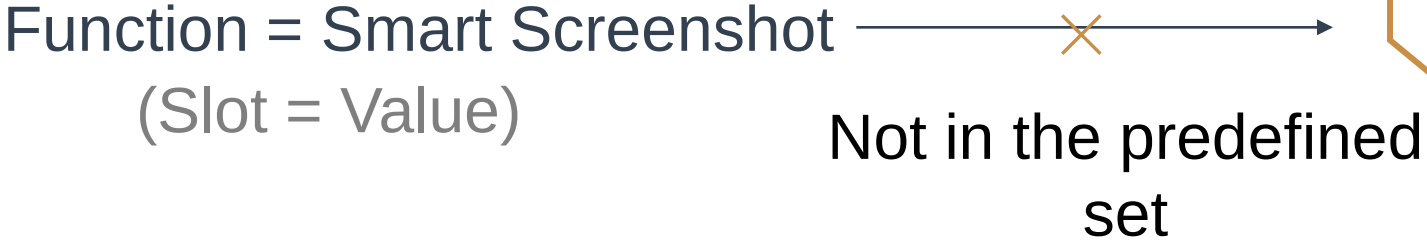


Introduction

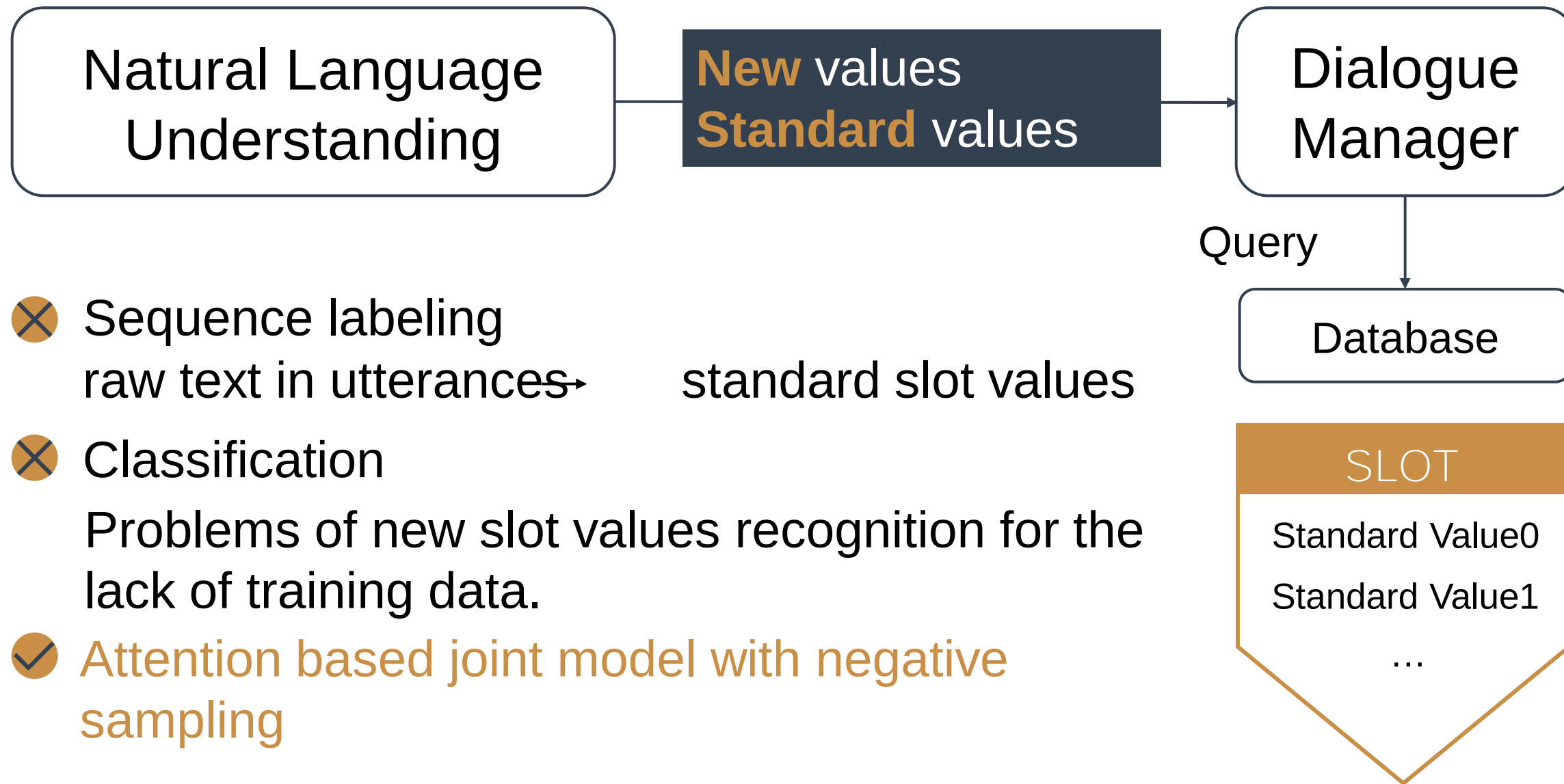


E.g

Can I scroll the screen to have a screenshot on my phone?



Introduction



- ✘ Sequence labeling
raw text in utterances → standard slot values
- ✘ Classification
Problems of new slot values recognition for the lack of training data.
- ✔ Attention based joint model with negative sampling

CHAPTE R

2

Related work

Sequence labeling

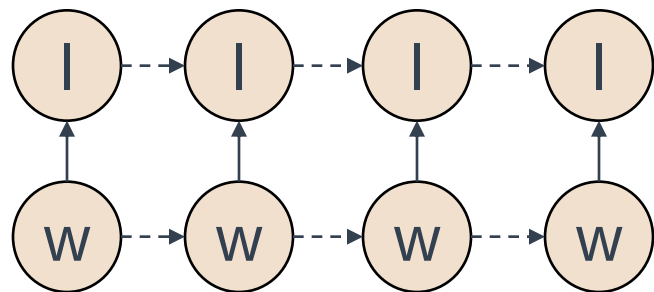
Pipeline based

Classification based

Related work

Sequence labeling

Sequence labels
words



O O B I I I O O O
⊕ ↑ ↑ ↑ ↑ ↑ ↑ ↑ ↑
Can I **have a video chat** on my phone?

Pipeline based methods

Can I **have a video chat** on my phone? 1 extract the raw texts from utterance

have a video chat → FaceTim e 2 normalized the texts into standard slot values

Classification based methods

Can I **have a video chat** on my phone? → FaceTim e

Related work

Sequence labeling

- ⊗ Extra normalization operations

Pipeline based methods

- ⊗ Prone to accumulating errors

Classification based methods

- ⊗
 - Disability to deal with new slot values
 - Losing local information

Related work

Sequence labeling

Xuezhe Ma and Eduard Hovy. *End-to-end sequence labeling via bi-directional lstm-cnns-erf*. 2016

Gokhan Tur, Dilek Hakkani-Tur et al. *Sentence simplification for spoken language understanding*. 2011

Kaisheng Yao, Baolin Peng et al. *Recurrent neural networks for language understanding*. 2013

Kaisheng Yao, Baolin Peng et al. *Spoken language understanding using long short-term memory neural networks*. 2014

Pipeline based methods

F Lefevre. *Dynamic bayesian networks and discriminative classifiers for multi-stage semantic interpretation*. 2007

Peter Zhen, Benjamin Douglas et al. *A speech-driven second screen application for tv program discovery*. 2014

Classification based methods

Rahul Bhagat, Anton Leuski, and Eduard Hovy. *Statistical shallow semantic parsing despite little training data*. 2005

François Mairesse, Milica Gasic et al. *Spoken language understanding from unaligned data using discriminative classification models*. 2009

Ana Mendes Pedro Mota, Luísa Coheur. *Natural language understanding as a classification process: report of initial experiments and results*. 2012

CHAPTE
R

3

Motivation

Motivation

Sequence

classification based methods

- Take local information into consideration
- Obtains normalized slot values directly



Attention based joint model with negative sampling

- Able to deal with new slot values

CHAPTE R

4

Proposed model

Attention based joint model
Negative sampling

Proposed

Attention-based joint
model

model



Model



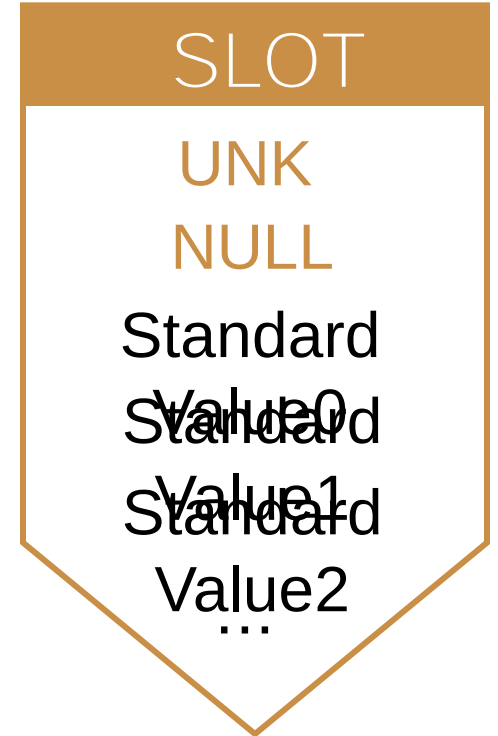
Proposed

Attention-based joint
model

model

Model

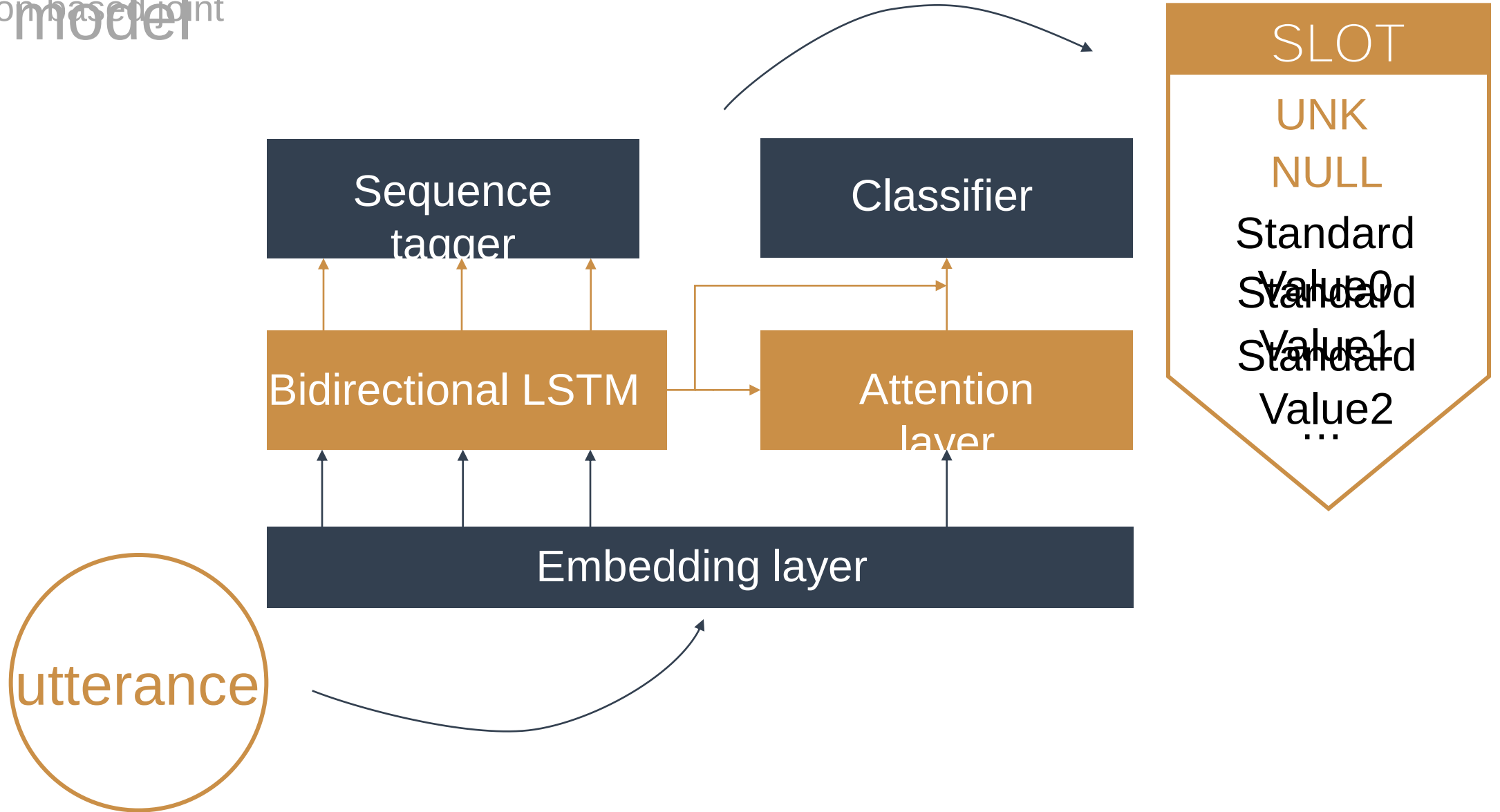
utterance



Proposed

Attention-based joint
model

model



Proposed

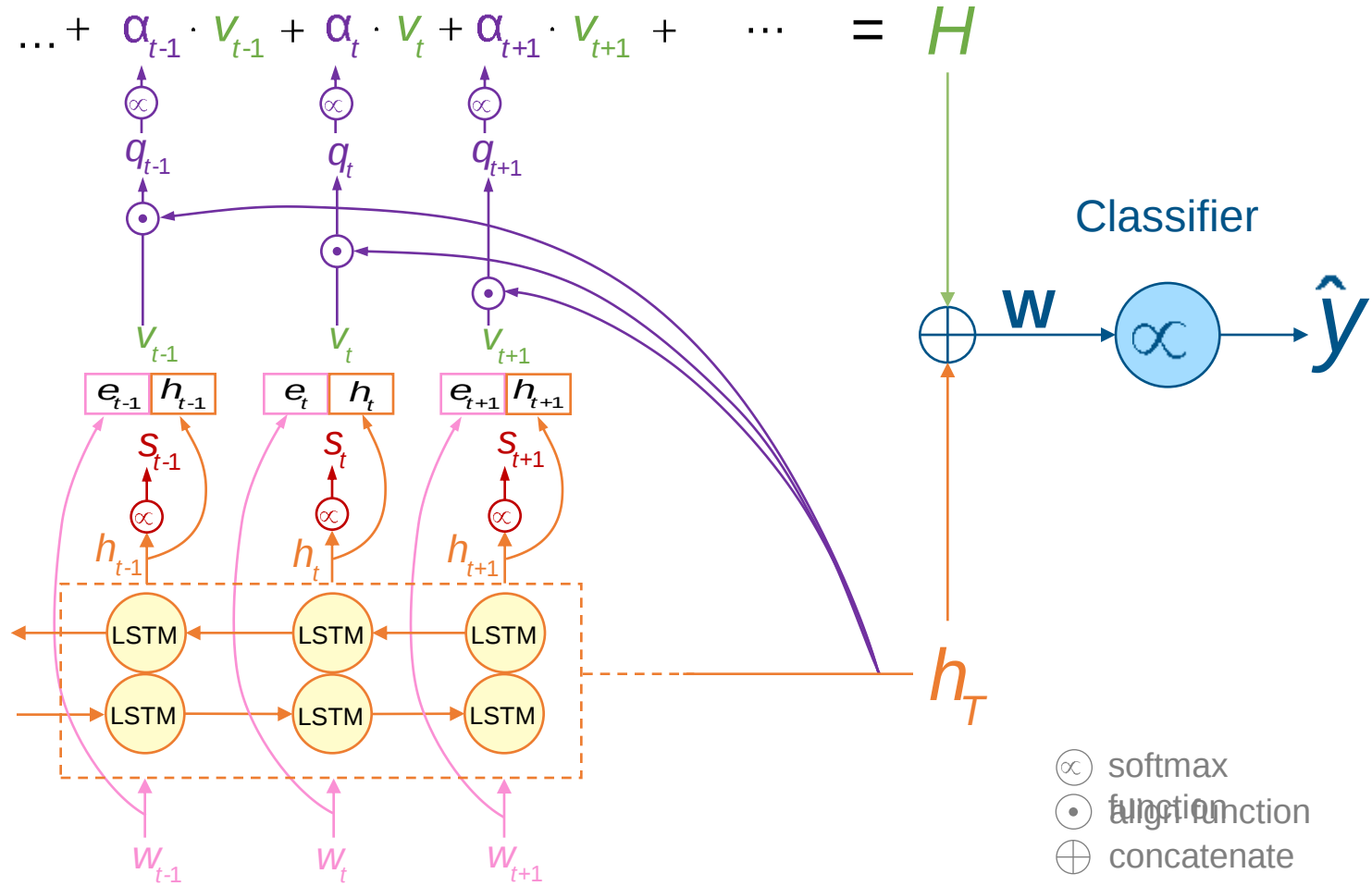
Attention-based joint
model

Attention layer

Sequence tagger

Bidirectional
LSTM

Embedding layer



Proposed

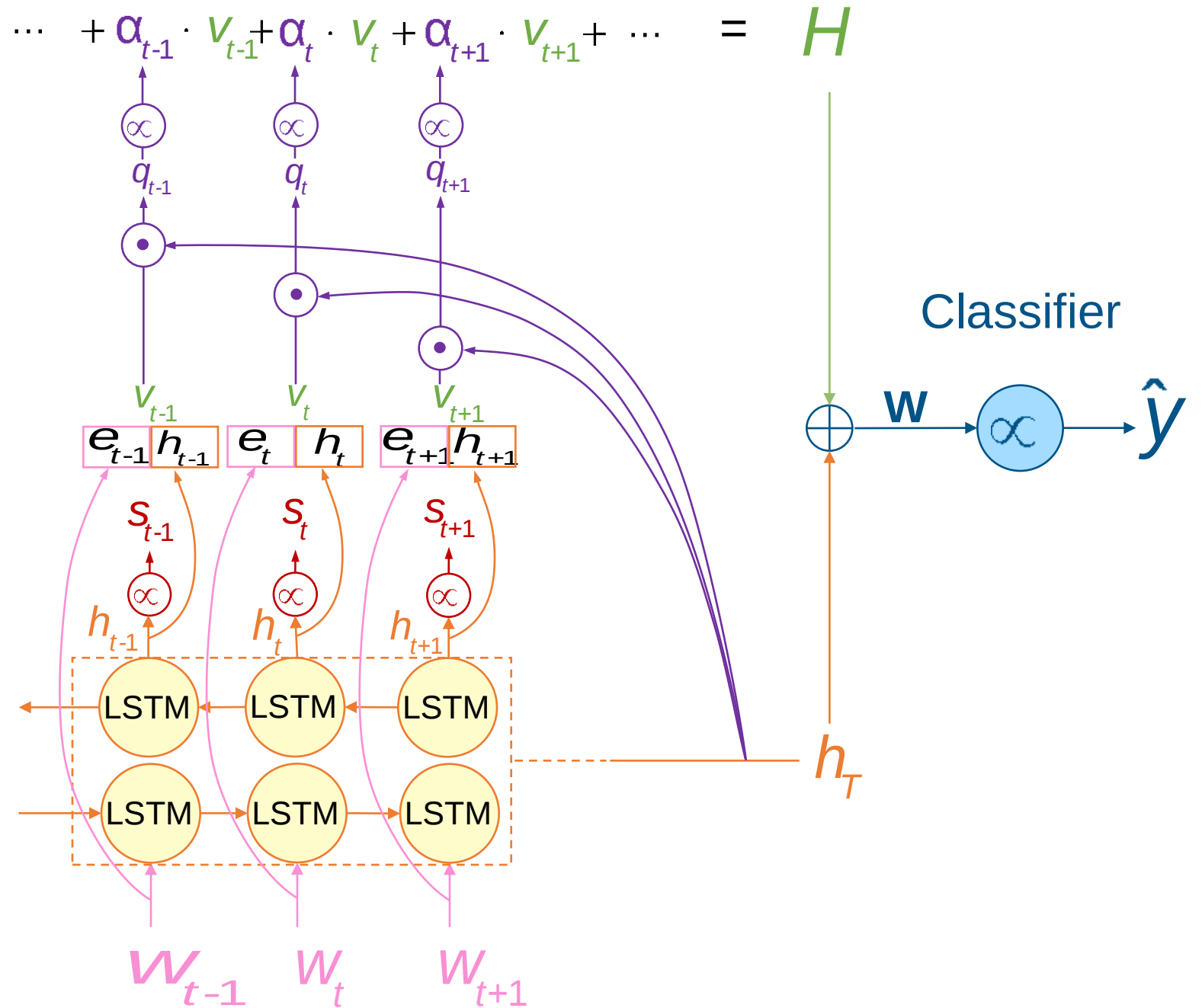
Attention-based joint
model

Attention
layer

Sequence tagger

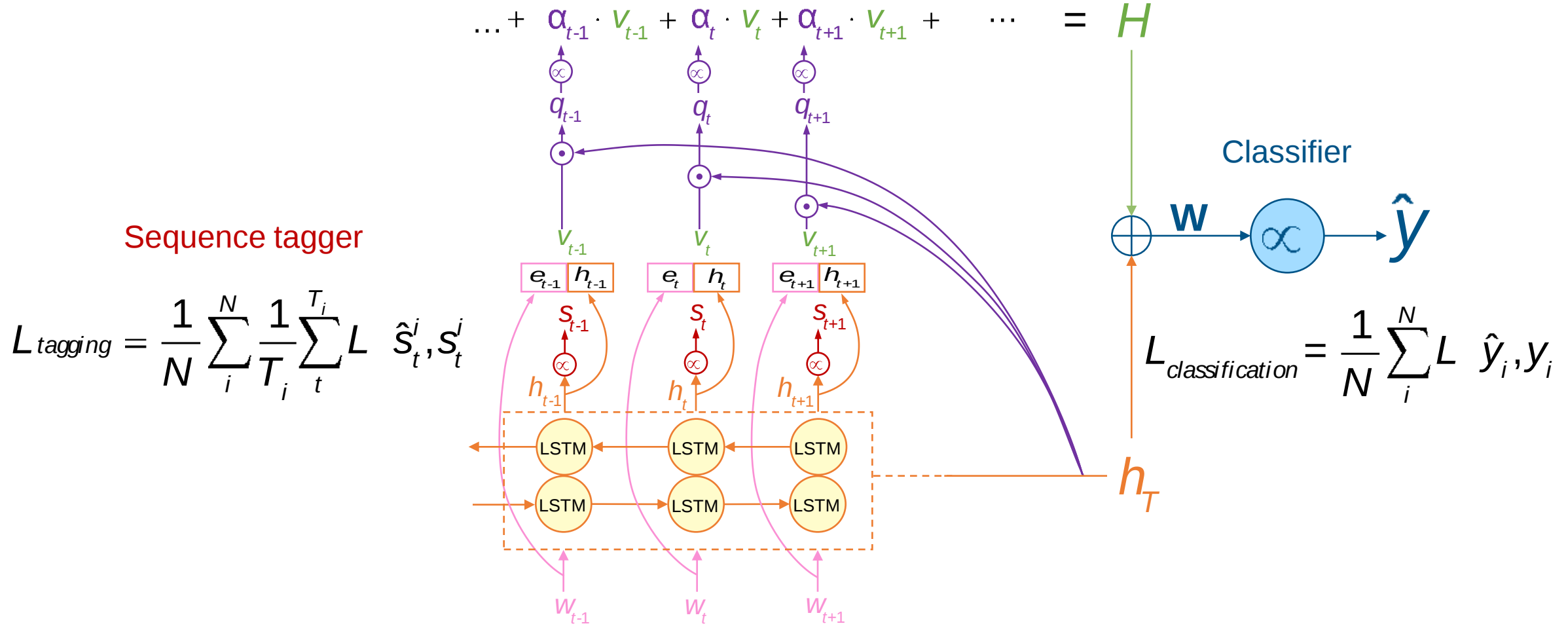
Bidirectional LSTM

Embedding layer



Proposed

Attention-based joint model



$$L = gL_{tagging} + (1 - g)L_{classification}$$

Proposed

Negative sampling
model

Models will fail in recognizing new slot values without corresponding training data

Construct negative samples of the existing slot values to simulate new ones

Existing slot value: Can I have a video chat on my phone?
Slot = FaceTime

New slot value: Can I scroll the screen to have a screenshot on my phone?
Slot = Smart Screenshot

Shared context:

Can I ... on my phone?
Can I Random words on my phone?
Non-value

Proposed

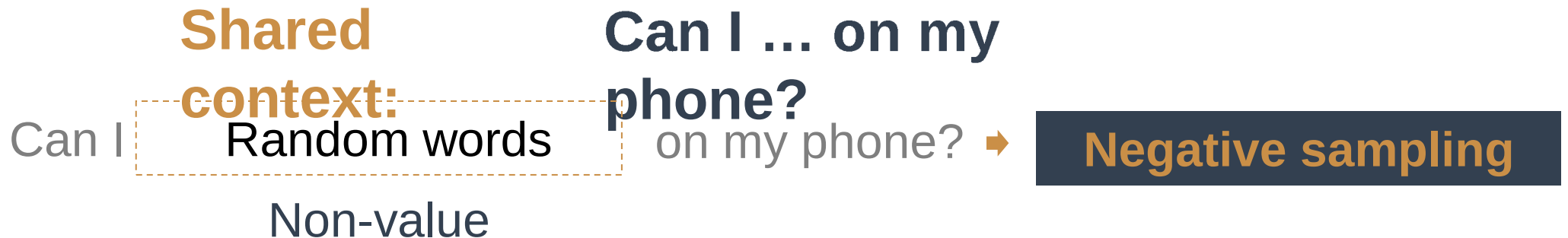
Negative sampling
model

Models will fail in recognizing new slot values without corresponding training data

Construct negative samples of the existing slot values to simulate new ones

Existing slot value: Can I have a video chat on my phone?
Slot = FaceTime

New slot value: Can I scroll the screen to have a screenshot on my phone?
Slot = Smart Screenshot



CHAPTE
R

5

Experiments

Results
Analyses

Experiments

Results

Dataset: DSTC(English) |

Service(Chinese)

DSTC --- an English dataset from a public contest and we use DSTC2 and DSTC3 together. It collects 5510 dialogues about hotels and restaurants booking. Only the slot

Service --- a Chinese dialogue dataset which is mainly about consultation for cell phones and contains a single slot named *'function'*.

Corpus	DSTC			Service			
	train	dev	test	train	dev	test	
Original data	old	2805	937	917	3682	514	1063
	new	0	113	275	0	15	64
	null	2244	840	953	427	64	109
negative samples	561	0	0	736	0	0	
overall size	5610	189	214	4815	593	1236	

Statistics of two datasets

Corpus	DSTC			Service		
	train	dev	test	train	dev	test
old	66	64	65	80	55	67
new	0	21	21	0	13	44

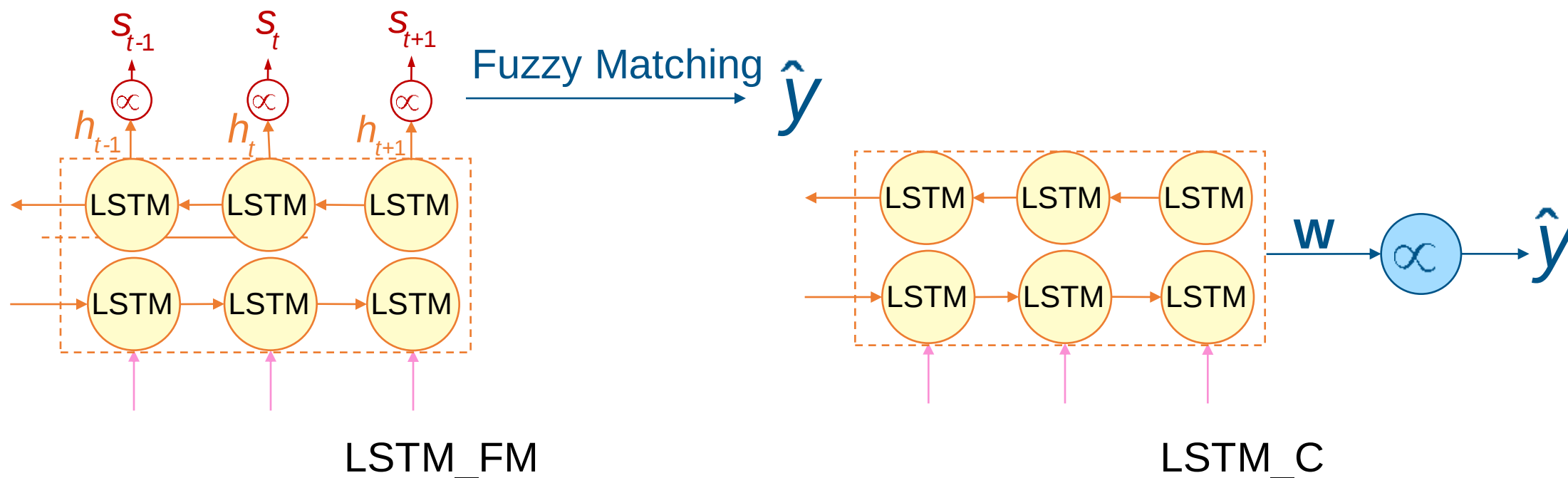
Statistics of value types

Experiments

Results

Baselines: LSTM_FM | LSTM_C, without negative samples

- LSTM_FM: pipeline based method, labeling the words with slot values tags by LSTM model and then normalized them into standard values by Fuzzy Matching.
- LSTM_C: classification based method, encoding the utterance by LSTM model and then use a full-connected layer as a Classifier.



Experiments

Results

Results:

	DSTC				Service			
	all	NEW	OLD	NULL	all	NEW	OLD	NULL
LSTM_FM	0.849 1	0.3063	0.9670	0.8923	0.8981	0.5693	0.9320	0.569 3
LSTM_C	0.829 0	0.0000	0.9249	0.9761	0.9210	0.0000	0.9720	0.964 3
(a) F1-scores of classification for different models	0.969							
AJM_NS(ours)	0.969 2	0.8621	0.9738	0.9822	0.9749	0.7759	0.9881	0.963 3

Experiments

Results

Results:

	DSTC				Service			
	all	NEW	OLD	NULL	all	NEW	OLD	NULL
LSTM_FM	0.849 1	0.3063	0.9670	0.8923	0.8981	0.5693	0.9320	0.569 3
LSTM_C	0.829 0	0.0000	0.9249	0.9761	0.9210	0.0000	0.9720	0.964 3
(a) F1 scores of classification for different models	0.969 2	0.8621	0.9738	0.9822	0.9749	0.7759	0.9881	0.963 3

Experiments

Results

Results:

	DSTC				Service			
	all	NEW	OLD	NULL	all	NEW	OLD	NULL
LSTM_FM	0.849 1	0.3063	0.9670	0.8923	0.8981	0.5693	0.9320	0.569 3
LSTM_C	0.829 0	0.0000	0.9249	0.9761	0.9210	0.0000	0.9720	0.964 3
AJM_NS(ours)	0.966 2	0.8621	0.9738	0.9822	0.9749	0.7759	0.9881	0.963 3

(a) F1 scores of classification for different models

	DSTC			Service		
	all	NEW	OLD	all	NEW	OLD
LSTM_FM	0.854 6	0.2363	0.983 7	0.8850	0.2615	0.926 9
AJM_NS(ours)	0.902 4	0.5684	0.994 6	0.912	0.1219	0.967 3

(b) F1 scores on sequence labeling

Experiments

Results

Results:

	DSTC				Service			
	all	NEW	OLD	NULL	all	NEW	OLD	NULL
LSTM_FM	0.849 1	0.3063	0.9670	0.8923	0.8981	0.5693	0.9320	0.569 3
LSTM_C	0.829 0	0.0000	0.9249	0.9761	0.9210	0.0000	0.9720	0.964 3
AJM_NS(ours)	0.969 2	0.8621	0.9738	0.9822	0.9749	0.7759	0.9881	0.963 3

(a) F1 scores of classification for different models

	DSTC			Service		
	all	NEW	OLD	all	NEW	OLD
LSTM_FM	0.854 6	0.2363	0.983 7	0.8850	0.2615	0.926 9
	0.902 4	0.5684	0.994 6	0.9122 6	0.1219 9	0.967 3

(b) F1 scores on sequence labeling

Experiments

Results

Comparison inside model: Attention mechanism(AJM) & Negative sampling(JM_NS)

	DSTC				Service			
	all	NEW	OLD	NULL	All	NEW	OLD	NULL
Full(AJM_NS)	0.9632	0.8621	0.9738	0.9822	0.9749	0.7759	0.9881	0.9633
-Attention only(JM_NS)	0.9515	0.8129	0.9739	0.9699	0.9700	0.7207	0.9862	0.9585
-NS only(AJM)	0.8247	0.0000	0.9426	0.9492	0.9234	0.0000	0.9761	0.9511

Experiments

Results

Comparison inside model: Attention mechanism(AJM) & Negative sampling(JM_NS)

	DSTC				Service			
	all	NEW	OLD	NULL	All	NEW	OLD	NULL
Full(AJM_NS)	0.9632	0.8621	0.9738	0.9822	0.9749	0.7759	0.9881	0.9633
-NS only(AJM)	0.8247	0.0000	0.9426	0.9492	0.9234	0.0000	0.9761	0.9511

- Confusion metrics

	DSTC									
	NEW OLD NULL				NEW OLD NULL					
AJM	NEW	0	184	91	➔	NEW	225	33	17	AJM_NS
	OLD	0	916	1		OLD	9	908	0	
	NULL	0	9	944		NULL	13	4	936	
	Service									
	NEW OLD NULL				NEW OLD NULL					
AJM	NEW	0	55	9	➔	NEW	45	17	2	AJM_NS
	OLD	0	106	0		OLD	4	105	2	
	NULL	0	3	107		NULL	3	7	105	

Experiments

Results

Comparison inside model: Attention mechanism(AJM) & Negative sampling(JM_NS)

	DSTC				Service			
	all	NEW	OLD	NULL	All	NEW	OLD	NULL
Full(AJM_NS)	0.9632	0.8621	0.9738	0.9822	0.9749	0.7759	0.9881	0.9633
-NS only(AJM)	0.8247	0.0000	0.9426	0.9492	0.9234	0.0000	0.9761	0.9511

- **Classification results based on negative samples**

	DSTC				Service			
	all	NEW	OLD	NULL	all	NEW	OLD	NULL
LSTM_FM_NS	0.857	0.353	0.928	0.924	0.964	0.620	0.900	0.648
	2	6	6	1	2	3	9	8
LSTM_C_NS	0.954	0.826	0.963	0.982	0.968	0.710	0.982	0.981
	3	1	7	2	4	3	5	5
→ JM_NS	0.951	0.812	0.973	0.969	0.970	0.720	0.986	0.958
	5	9	9	9	0	7	2	5

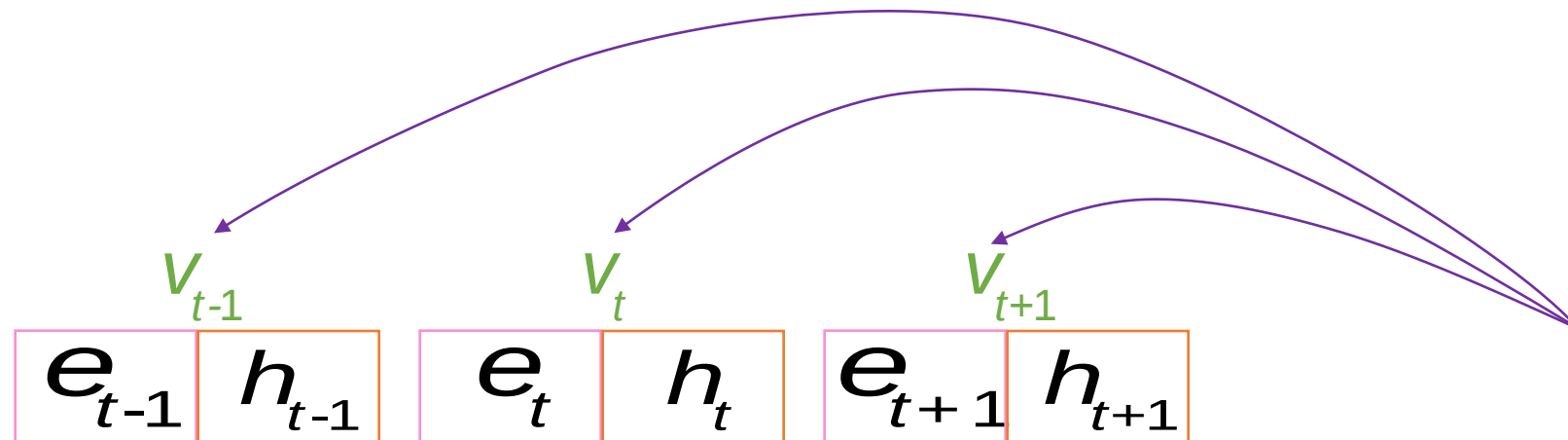
Experiments

Results

Comparison inside model: Attention mechanism(AJM) & Negative sampling(JM_NS)

	DSTC				Service			
	all	NEW	OLD	NULL	All	NEW	OLD	NULL
Full(AJM_NS)	0.9632	0.8621	0.9738	0.9822	0.9749	0.7759	0.9881	0.9633
-Attention only(JM_NS)	0.9515	0.8129	0.9739	0.9699	0.9700	0.7207	0.9862	0.9585

- Attention mechanism



Experiments

Results

Comparison inside model: Attention mechanism(AJM) & Negative sampling(JM_NS)

	DSTC				Service			
	all	NEW	OLD	NULL	All	NEW	OLD	NULL
Full(AJM_NS)	0.9632	0.8621	0.9738	0.9822	0.9749	0.7759	0.9881	0.9633
-Attention only(JM_NS)	0.9515	0.8129	0.9739	0.9699	0.9700	0.7207	0.9862	0.9585

- Attention mechanism

		DSTC										Service											
		i	want	an	indonesia	restauran	t	in	the	north	part	of	town	A7	的	儿	童	模	式	怎么	解		
Full (AJM_NS)	heatmap																						
	True	O	O	O	B-food	O	O	O	O	O	O	O	indonesia	O	O	B-func	I-func	I-func	I-func	O	O	O	儿童模式
	Pred	O	O	O	B-food	O	O	O	O	O	O	O	indonesia	O	O	B-func	I-func	I-func	I-func	O	O	O	儿童模式
-Attention (JM_NS)	True	O	O	O	B-food	O	O	O	O	O	O	indonesia	O	O	B-func	I-func	I-func	I-func	O	O	O	儿童模式	
	Pred	O	O	O	B-food	O	O	O	O	O	O	african	O	O	B-func	I-func	I-func	I-func	O	O	O	指纹解锁	

CHAPTE
R

6

Conclusion

Conclusion

We propose an attention based joint model with negative sampling.

- Maps the utterance into standard slot values directly without extra normalization operations
- Negative sampling for existing values for a certain slot S enables our model to effectively recognize new slot values
- Joint model collaborated by attention mechanism promotes the performance
- Experimental results demonstrate that our model achieves impressive improvements on new slot values with less damage on other sub-datasets

THANK YOU