

Improving Pointer-Generator Network with Keywords Information for Chinese Abstractive Summarization

Xiaoping Jiang, Po Hu, Liwei Hou and Xia Wang
School of Computer Science, Central China Normal University,
Wuhan 430079, China
{jiangxp, houliwei }@mails.ccnu.edu.cn
phu@mail.ccnu.edu.cn, wangx_cathy@163.com

Our current topic

- ▶ Some problems with attention-based Seq2Seq model
 - Inaccurate content
 - Insufficient summary
 - Repetitive words
 - Rare or unknown words

Our current topic

▶ Pointer-generator network

- Inaccurate content
- Insufficient summary
- Repetition words
- Rare or unknown words

Pointer-Generator Network

See, A., Liu, P.J., Manning, C.D.: Get to the Point: Summarization with Pointer-Generator Networks. ACL (2017)

Our current topic

▶ Our method

- Inaccurate content
- Insufficient summary
- Repetition words
- Rare or unknown words

Keywords
Attention
Mechanism

Topic coverage

Pointer-generator network

► Attention distribution

$$\begin{aligned} e_i^t &= v^T \tanh(W_h h_t + W_s s_t + W_c c_i^t + b_{attn}) \\ a^t &= \text{softmax}(e^t) \end{aligned}$$

► Generation probability

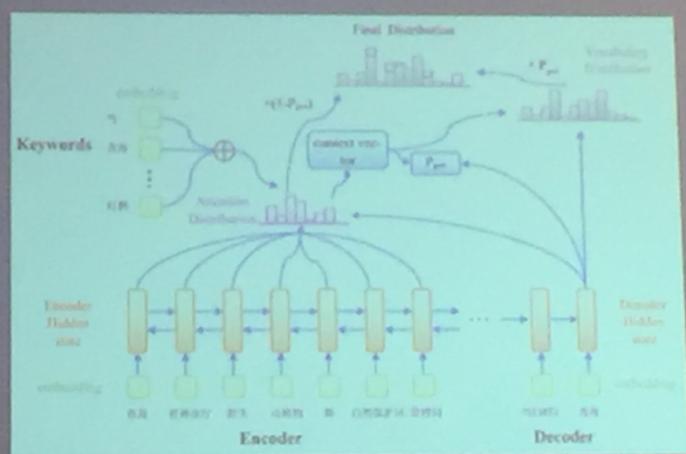
$$\begin{aligned} P_{gen} &= \sigma(w_h^T h_t^* + w_s^T s_t + w_x^T x_t + b_{ptr}) \\ P(w) &= p_{gen} P_{vocab}(w) + (1 - p_{gen}) \sum_{i:w_i} a_i^t \end{aligned}$$

Keywords attention mechanism

► Attention distribution

$$e_i^t = v^T \tanh(W_h h_i + W_s s_t + W_c c_i^t + W_t t + b_{attn})$$

$$\tau = \sum_{i=1}^d k_i$$



Implementation details

- Extract news article's keywords
Method: TextRank word_num: 8
- Save {article, summary, keywords} into a file
- Train and evaluate the model until the loss value is stable and converge
- Run beam search decoding for testing set
beam_size: 6

Reference code: <https://github.com/abisee/pointer-generator>

Experimental Results

Models	R-1	R-2	R-L
TextRank	35.78	23.58	29.60
Pointer-Generator	43.44	29.46	37.66
Pointer-Generator (character)	38.01	24.43	32.32
TextRank + Pointer-Generator	42.79	28.93	37.33
Our Model	44.60	30.46	38.83

Table 2: Comparison results on the NLPCC 2018 validation dataset

Dataset

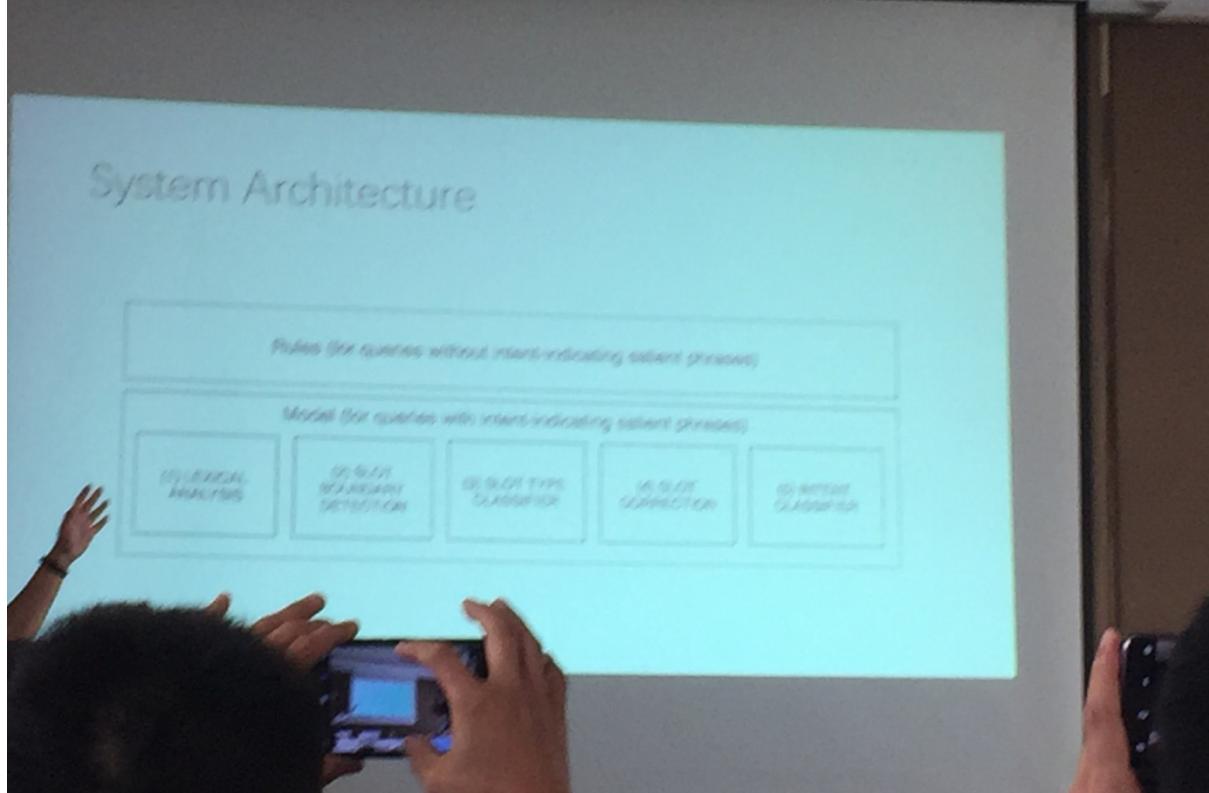
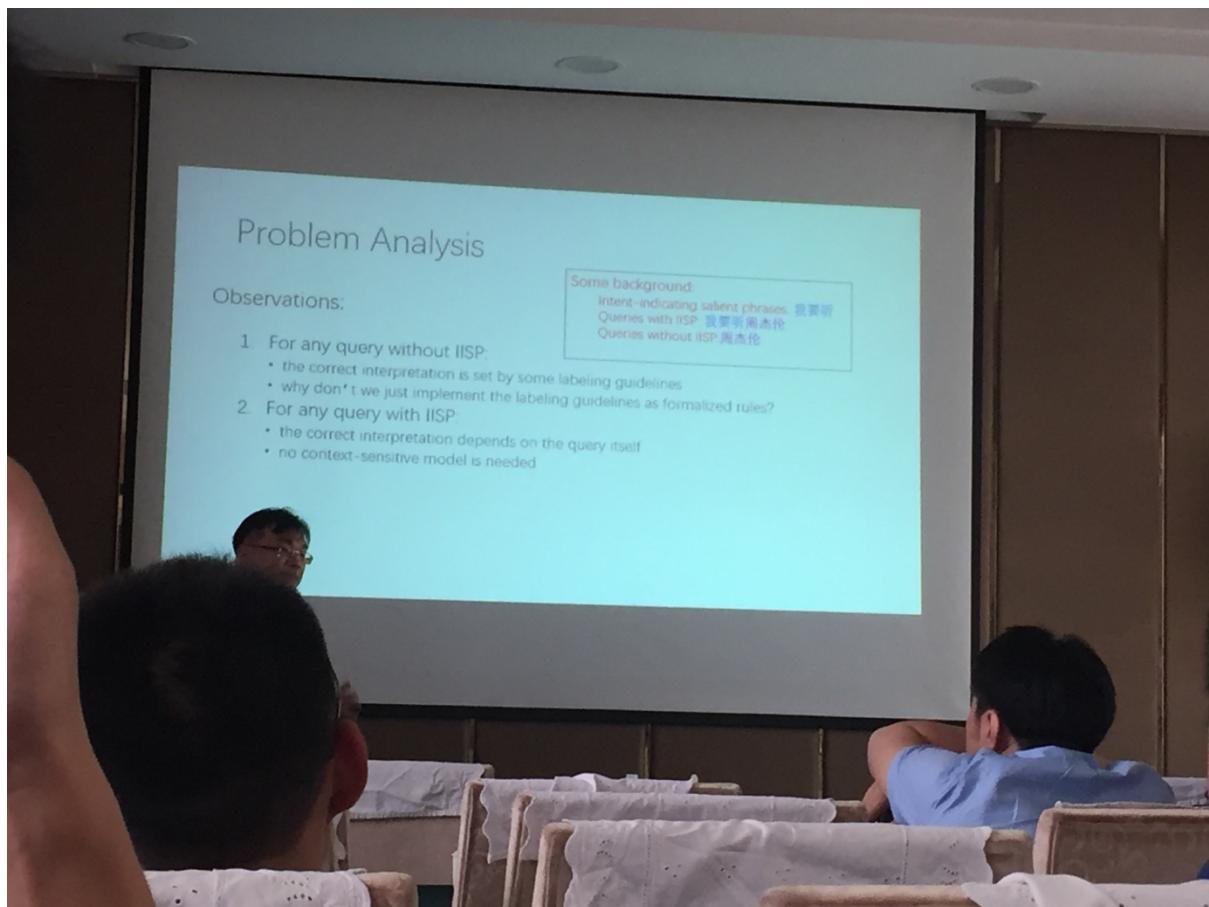
- ▶ NLPCC 2018 shared task3
- ▶ Chinese single document summarization dataset

Length	Training Set (50000)		Validation Set (2000)		Testing Set (2000)	
	Document	Summary	Document	Summary	Document	Summary
Min	32	6	35	8	7	-
Max	13186	85	8871	44	7596	-
Average	579.4	25.9	579.9	25.8	426.2	-

Table 1: The statistics of the datasets

The Sogou System for the NLPCC-18 Shared Task on Spoken Language Understanding

Henry Li 李志瀛
SOGOU Voice Interaction Technology Center
搜狗语音交互技术中心



The SLU Model Pipeline

1. lexical analysis 导航 去 清华 科技 馆
2. slot boundary detection 导航 去 清华 科技 馆
 B I L
3. slot type classifier 导航 去 (清华 科技 馆) type=DESTINATION
4. slot correction 导航 去 (清华 科技 馆) type=DESTINATION
5. intent classifier (导航 去 (DESTINATION 清华 科技 馆)) intent=NAVIGATION

The SLU Model Pipeline

- | | | |
|---------|----------------------------|--|
| | 1. lexical analysis | 导航 去 清华 科技 馆 |
| CRF | 2. slot boundary detection | 导航 去 清华 科技 馆
B I L |
| LR | 3. slot type classifier | 导航 去 (清华 科技 馆) type=DESTINATION |
| | 4. slot correction | 导航 去 (清华 科技 馆) type=DESTINATION |
| XGBoost | 5. intent classifier | (导航 去 (DESTINATION 清华 科技 馆)) intent=NAVIGATION |

Evaluation—Intent Determination

Setting	F1
1. Baseline (plain XGBoost using lexical features only)	91.36%
2. (1) + slot features	95.26%
3. (2) + query length features	95.46%
4. (3) + mined training samples	95.97%

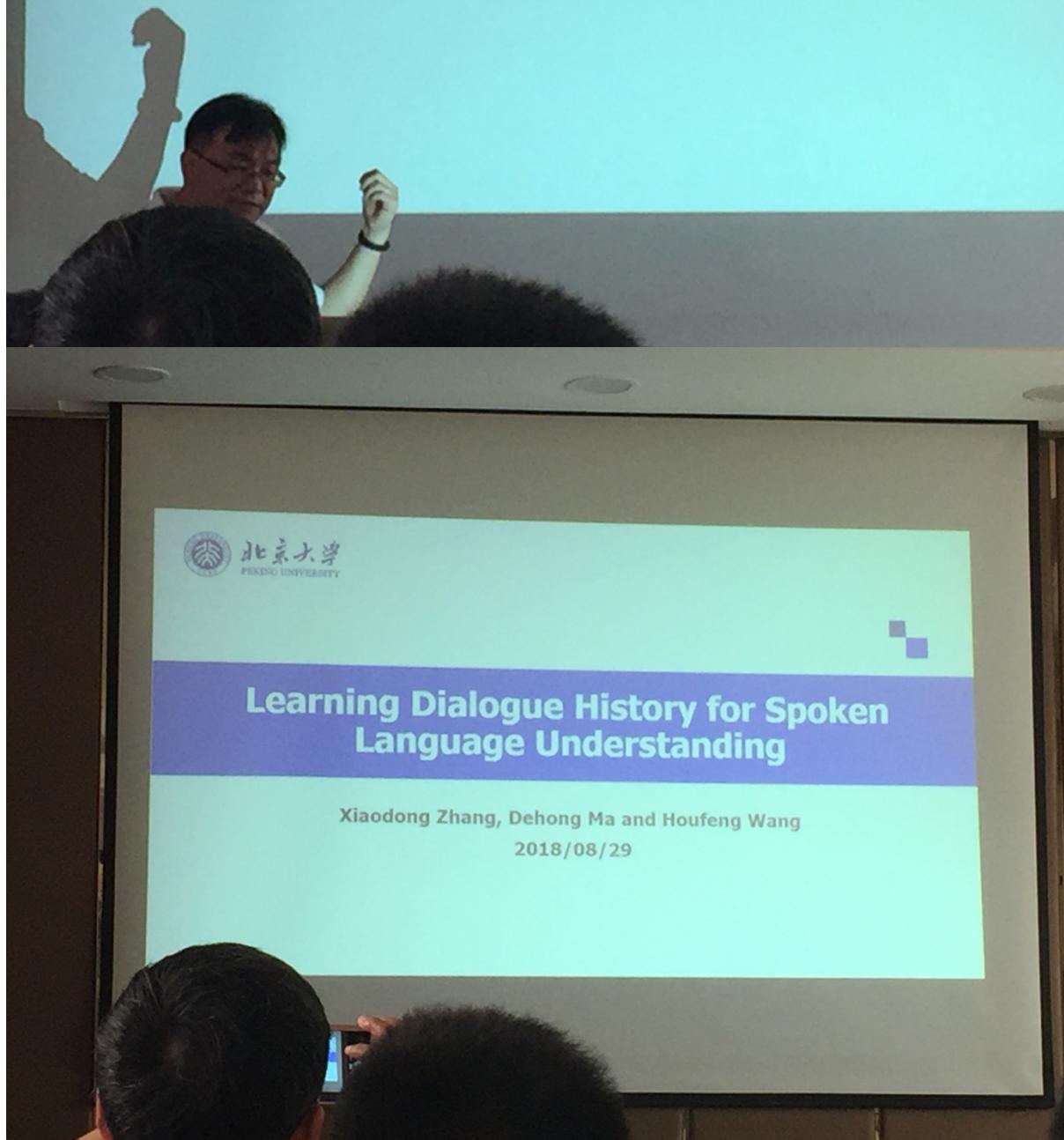
Is deep learning no use at all?

	Intent Determination			
	Lexical Features Only		All Features	
	F1	Precision	F1	Precision
XGBoost	91.36%	91.41%	95.97%	96.01%
CNN	87.96%	88.08%	94.82%	94.93%
CNN based on word embeddings trained on 5M SOGOU queries	91.94%	92.09%	95.13%	95.26%

CNN with pretrained embeddings achieves best performance in our own SLU task!

Evaluation—Slot Filling

Setting	Precision (before slot correction)	Precision (after slot correction)
1. Baseline (plain CRF + LR using lexical features only)	90.58%	91.70%
2. (1) + word/char combo strategy	91.31%	92.43%
3. (2) + dictionary features	92.39%	93.51%
4. (3) + mined training samples	92.73%	93.85%



■ Introduction

- Spoken language understanding
 - Converting users' queries expressed by natural language to structured representations.
- Two subtasks:
 - Intent Identification (classification)
 - Slot filling (sequence labeling)

Utterance	Go to Peking University from the Forbidden City
Slots	O O B-dest I-dest O B-orig I-orig I-orig
Intent	Navigation

2

■ Motivation

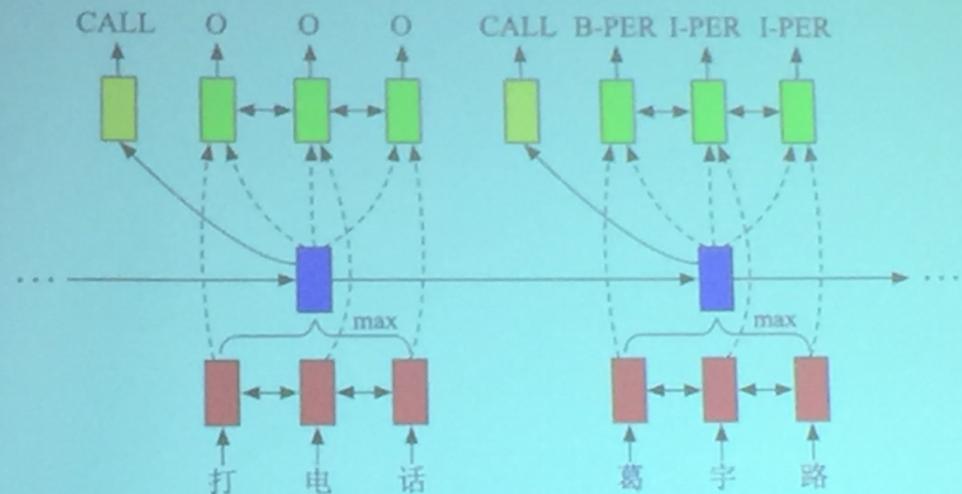
- Dialogue history is important to the understanding of the current utterance.

ID	Dialogue
1	A: Go to Peking University. B: Planning route to Peking University A: Cancel.
2	A: Call Tom. B: Calling Tom for you. A: Cancel.
3	A: Where are you going? B: Washington.
4	A: Who would you like to call? B: Washington.

- How to use dialogue history for the two subtasks?

3

HLSTM-SLU



4

Input

- HLSTM-SLU processes Chinese text at character level.
- The input contains two parts:
 - Characters
 - Additional features: part-of-speech and domain lexicons
- An example:

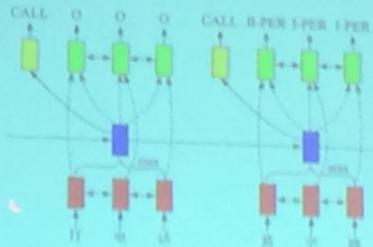
	B-A	L-A	B-NLR	(B-song)	I-song	
打	1	0	0	0	1	0
电	0	0	1	0	0	1
话	0	0	0	1	0	1

- The sparse feature vector is converted to a low-dimensional dense vector via a linear transformation.

5

Model Details

- A character-level LSTM to learn representation of a turn:
$$h_{i,j}^c = BLSTM(x_{i,j}, h_{i,j-1}^c, h_{i,j+1}^c)$$
- A fixed-length representation via max-pooling:
$$r_i^T = \max_{k=1}^n h_{i,k}^c$$
- A turn-level LSTM for a dialog:
$$r_i^D = LSTM(r_i^T, r_{i-1}^D)$$



6

Intent Tagger

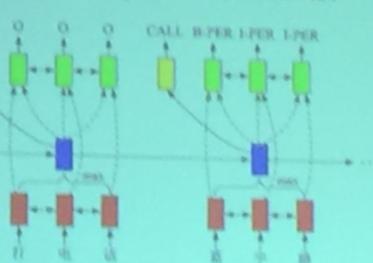
- A fully-connected layer for calculating scores at each turn:
$$s_i^T = W_s^T r_i^D + b_s^T$$
- A CRF layer for modeling label dependency between turns:

$$S(y^T) = \sum_{i=1}^m s_{i,y_i^T}^T + \sum_{i=1}^{m-1} A_{y_i^T, y_{i+1}^T}$$

- Training and testing:

$$p(y^T) = \frac{e^{S(y^T)}}{\sum_{\hat{y}^T \in Y^T} e^{S(\hat{y}^T)}}$$

$$\hat{y}^T = \underset{\hat{y}^T \in Y^T}{\operatorname{argmax}} S(\hat{y}^T)$$



7

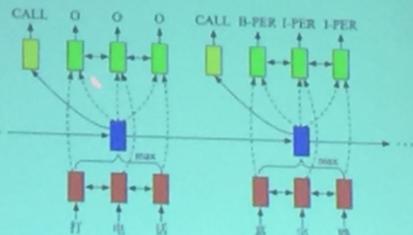
Slot Tagger

- For slot filling, the model rereads a turn with the contextual turn representation:

$$h_{i,j}^d = \text{BLSTM}(x_{i,j}^d, h_{i,j-1}^d, h_{i,j+1}^d)$$
$$x_{i,j}^d = [h_{i,j}^e; r_i^D]$$

- A CRF layer is also used for predict slot tags for each turn.
- The loss function of is the sum of loss for intent and slot:

$$L = -\frac{1}{m} \left(\log(p(y^T)) + \sum_{i=1}^m \log(p(y_i^e)) \right)$$



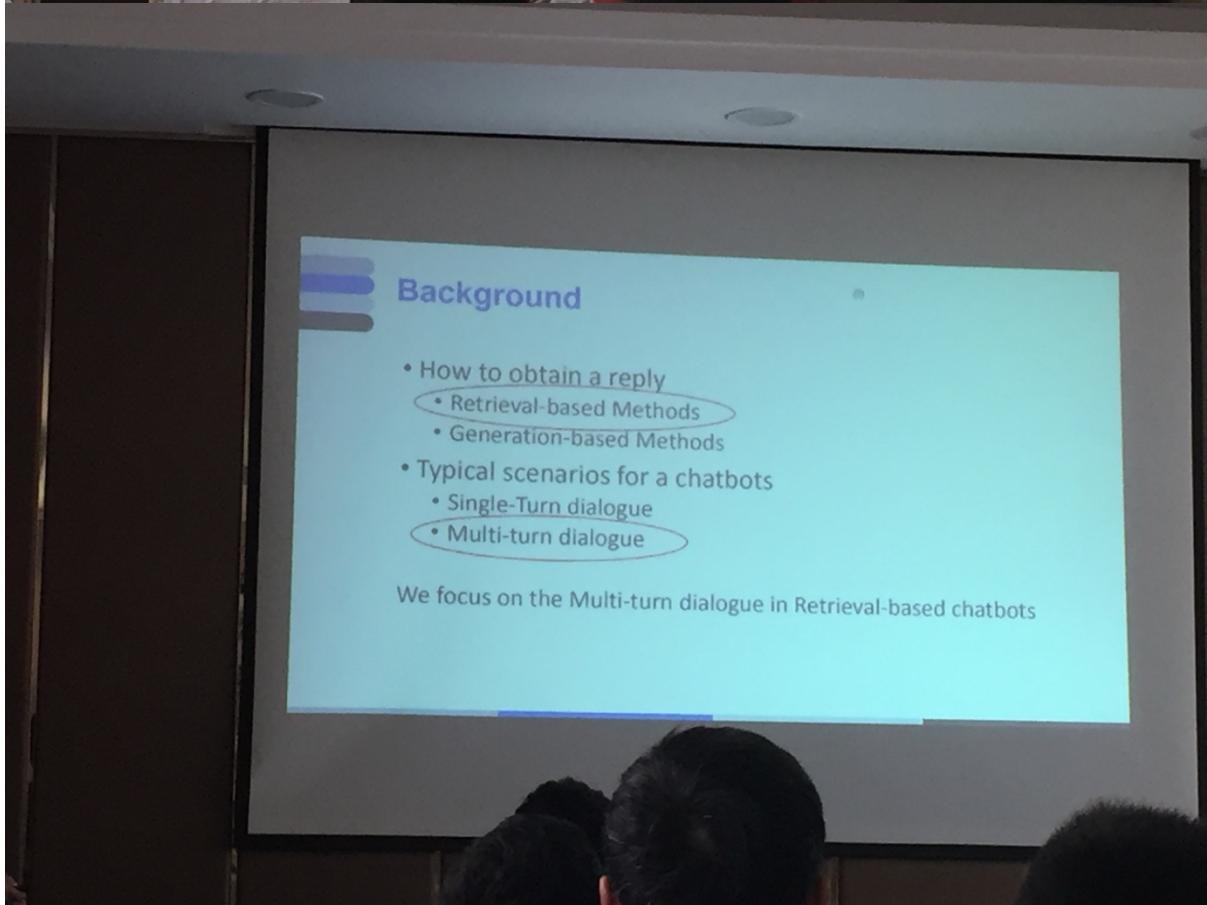
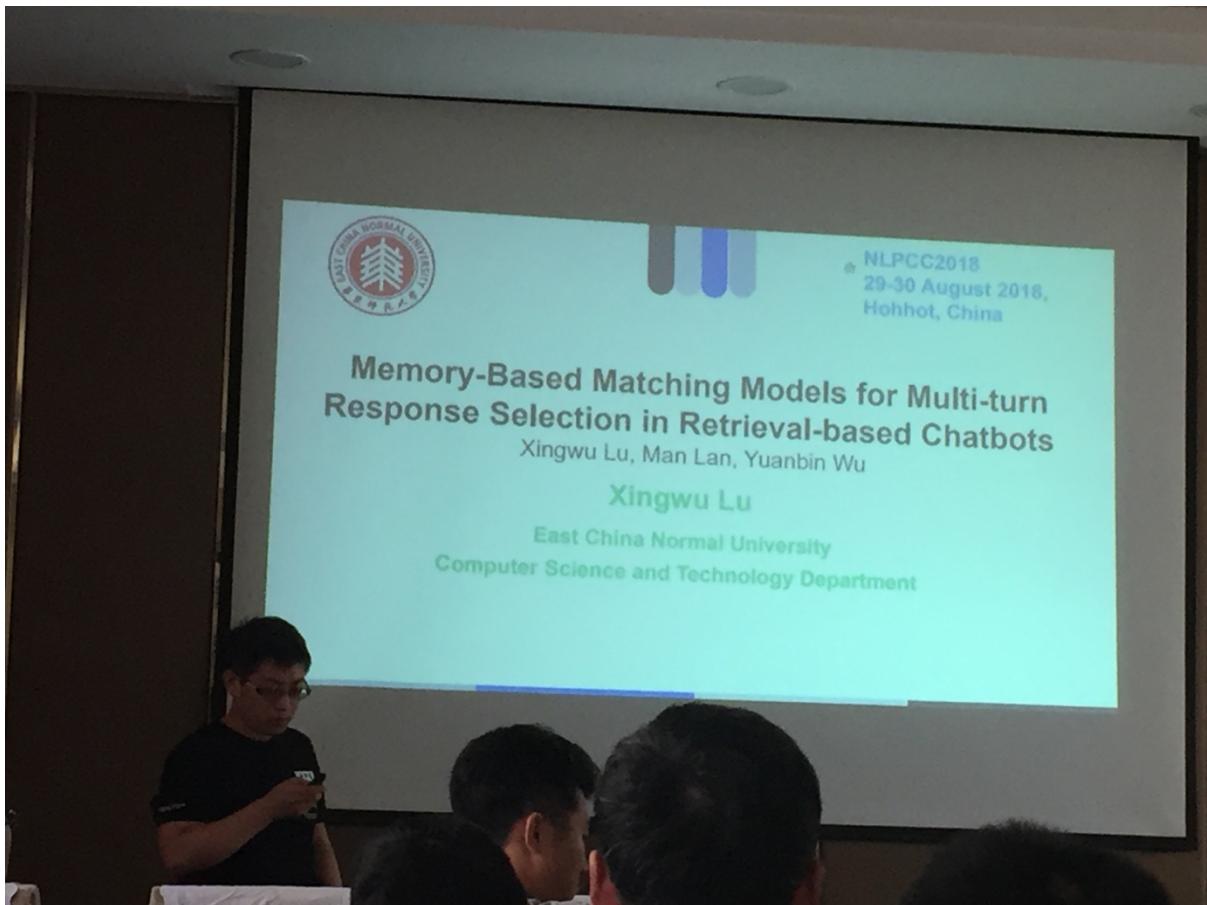
8

Experiments

- NLPCC 2018 Shared Task 4 dataset:
 - 4705 dialogues with 21352 turns in the training set and 1177 dialogues with 5350 turns in the testing set.
 - We do not use additional data.

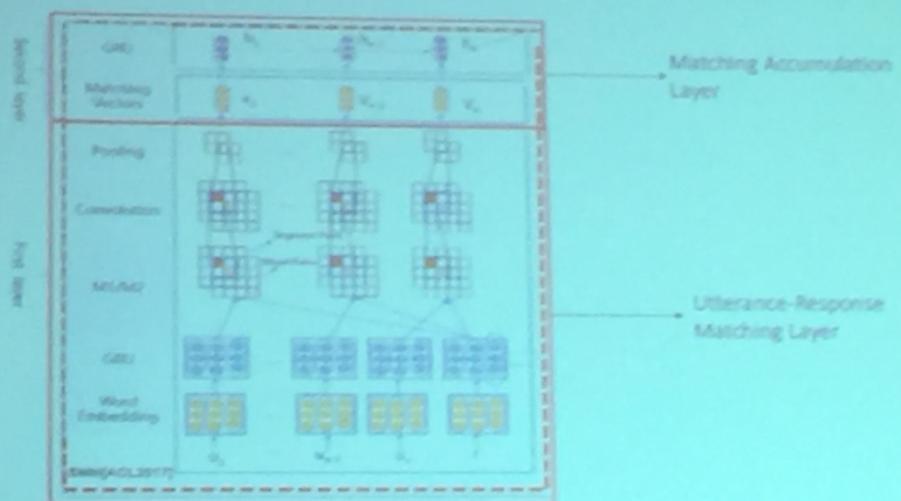
Method	F1 for Intent	Precision for Intent & Slot
BLSTM	88.56	83.38
BLSTM-FT	89.24	86.95
BLSTM-CRF	89.31	88.62
HLSTM-Intent	93.61	90.21
HLSTM-SLU	94.19	90.84

9

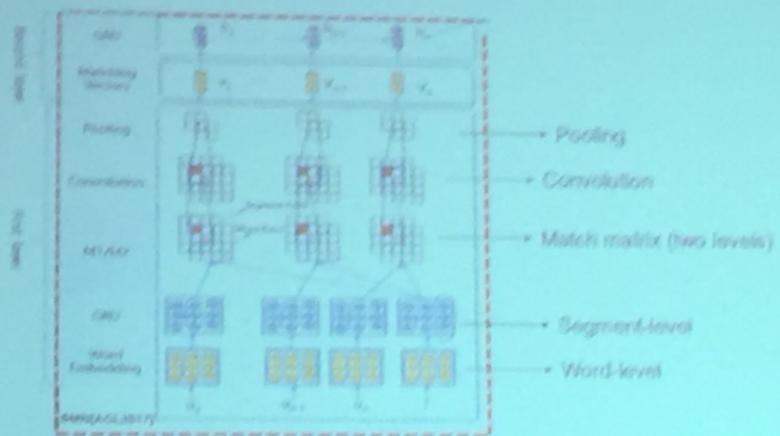


Methodology: Overview

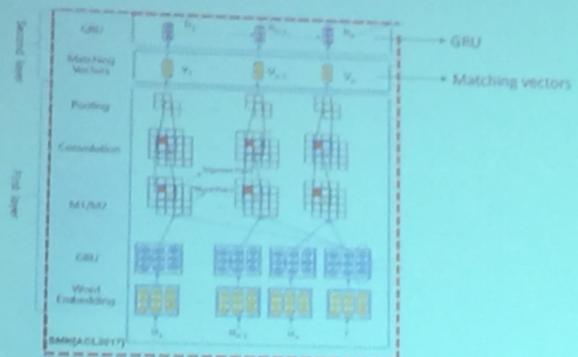
Sequential Matching Network (SMN)



SMN Module: First Layer



SMN Module: Second Layer

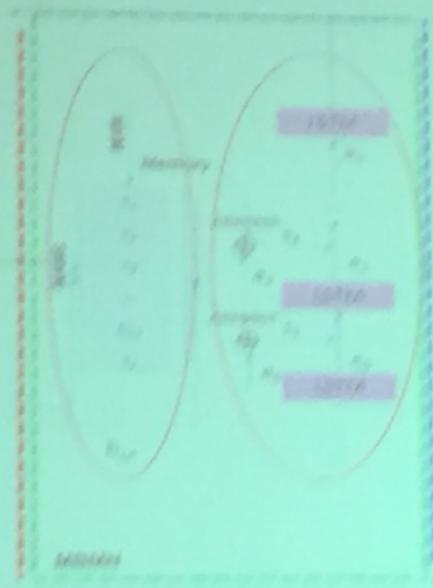


Memory-Based Matching Network (MBMN)

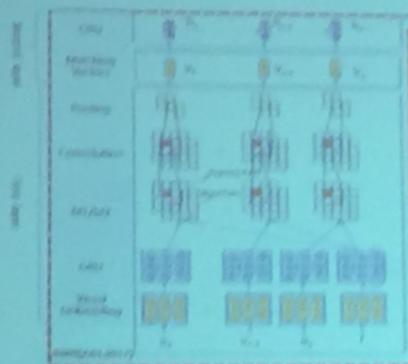
two steps

Memory Building

Apply Multiple
Attentions on
the Memory



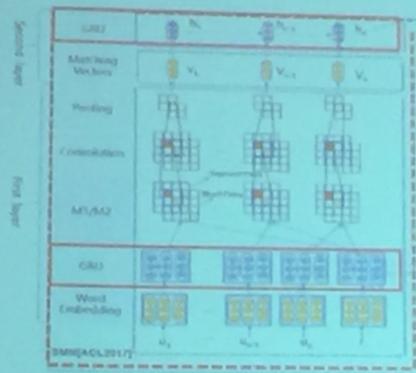
MBMN Module: Memory Building



Matching Vectors Memory (MVM)

y_1
 y_2
 y_3
 \vdots
 y_{n-1}
 y_n

MBMN Module: Memory Building



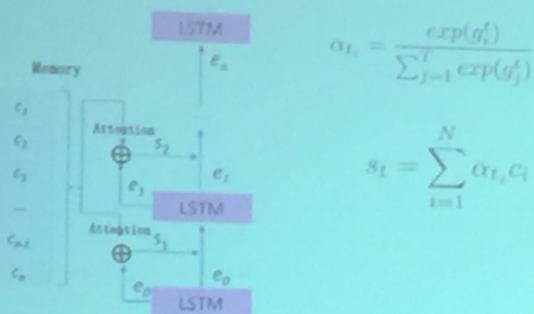
Sequential Matching Vectors Memory (SMVM)

$c_1, c_2, c_3, \dots, c_{n-1}, c_n$

$$c_t = \tanh(W_{1,1}h_{u_t, n_u} + W_{1,2}h_t + b_t)$$

MBMN Module: Multiple Attentions on Memory

$$g_i^t = v^T \tanh(W_e c_i + W_e e_{t-1} + W_r h_{r,n} + b_{attn})$$



NLP Features Module

- Word and Character Matching Feature:
 - following ten measure functions: $|A|, |B|, |A \cap B|, |A \cap B|/|A|, |A \cap B|/|B|,$ $|A - B|/|A|, |A - B|/|B|, |A \cap B|/|A \cup B|,$ $|A \cup B| - |A \cap B|/|A \cup B|, |||A| - |B|||$
- Unigram Feature:
 - linear functions: Cosine, Chebyshev, Manhattan and Euclidean
 - non-linear functions: Polynomial, Sigmoid and Laplacian

Matching Prediction

- Concatenate representations of three modules to calculate the final matching score with *softmax* function

$$g(u, r) = softmax(W_2[p_1, p_2, p_3] + b_2)$$

- Learn $g(u, r)$ by minimizing *cross entropy* with dataset D

$$-\sum_{i=1}^{|D|} [y_i \log(g(u_i, r_i)) + (1 - y_i)(1 - \log(g(u_i, r_i)))]$$

Conclusion

- We proposed an effective Memory-Based Matching model to address Multi-turn Response Selection
- The model consists of a SMN module, a MBMN module, and a NLP features module
- The model performance ranks the 1st among all the participants.

Response Selection of Multi-turn Conversation with Deep Neural Networks

Yunli Wang¹, Zhao Yan², Zhoujun Li¹, and Wenhan Chao¹

¹ Beihang University, Beijing, China

² Tencent, Beijing, China



Problem Formalization

- Assume that we have a data set $D = \{c_i, r_i, y_i\}_1^N$, where c_i is a context, r_i is a response. When $y_i=1$, r_i is a proper response for c_i , the situation is opposite if $y_i=0$.
- So we should use D to train a model M , $M(c, r)$ gives the matching degree between c and r .
- We regard this task as a classification problem, when r is the proper response for c , the label is 1, otherwise the label is 0.

Method Overview

- We proposed a matching network named RCFMFI, which focus on modeling relevance consistency of conversations.
- We adapted a matching network named SMFI, which has excellent performance for multi-turn matching problem.
- We proposed an ensemble strategy for the two models.

Relevance Consistency Matching Network

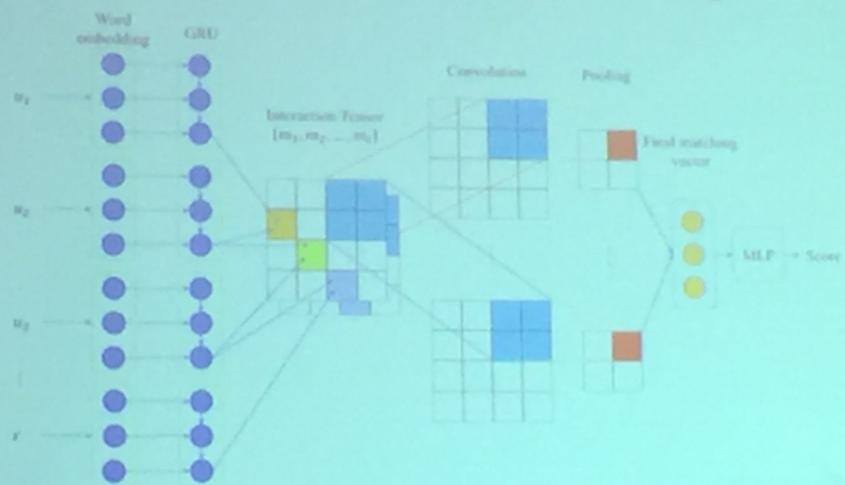


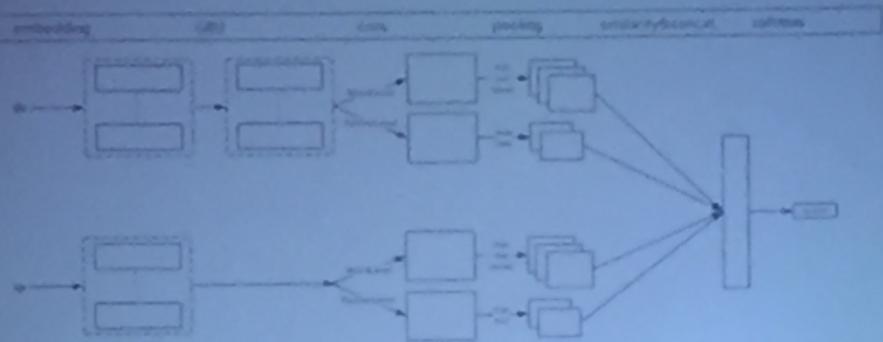
Fig. 1. The architecture of RCMN.

Semantic Matching

- Problem Formalization
 - Identify the $Q_i A_i$ from n candidate QA pairs $\{(Q_1, A_1), \dots, (Q_n, A_n)\}$ of Q-A KB that best matches Q_u .

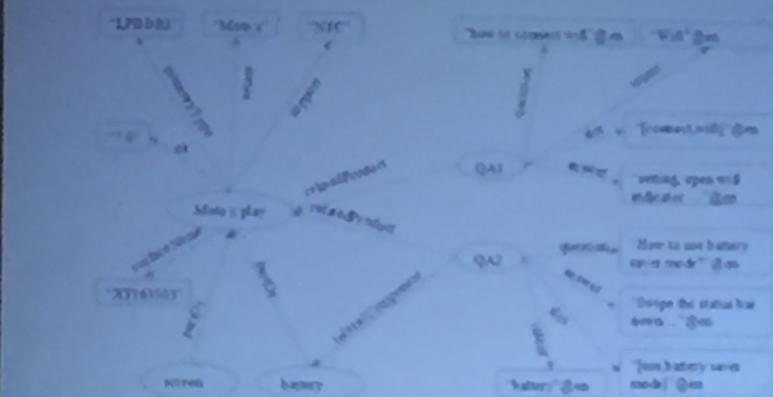
- Algorithm
 - TFIDF
 - WDM
 - SMN
 - MPCNN
 - Product Normalization
 - Matching

◦ Semantic Matching



Lenovo

◦ Data Set, QA-KB and Product-KB



Lenovo

Intent Category Classification

Model	Precision	Recall	F1
GBDT	0.67	0.68	0.67
BiLSTM	0.68	0.70	0.69
SVM	0.74	0.76	0.75

Semantic Matching

Model	Without Answer	With Answer	Average Response Time
TF-IDF	0.60	0.62	null
WMD	0.58	0.60	null
SMN	0.62	0.68	200ms(CPU)
MPCNN	0.72	0.84	50ms(CPU)
MPCNN_GRU	0.72	0.85	55ms(CPU)