

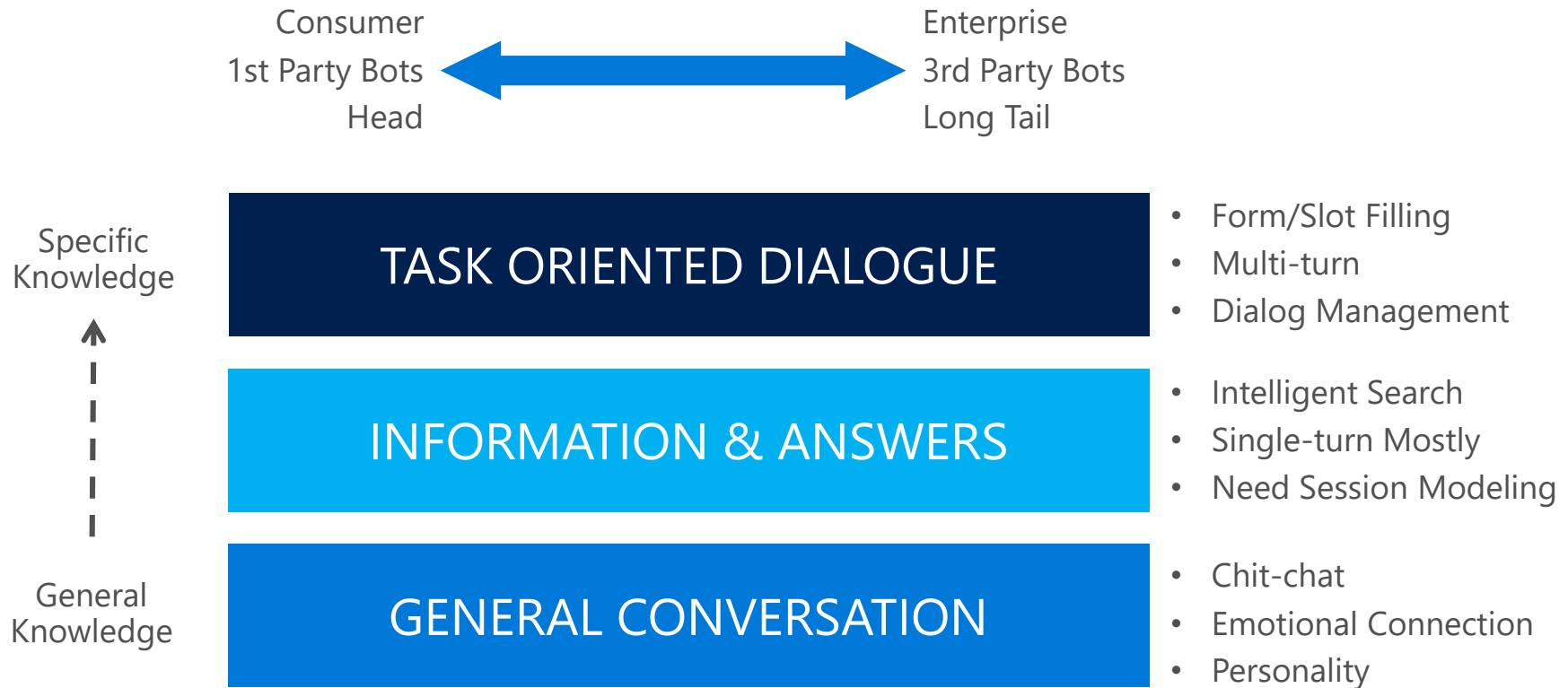
Natural Language Conversation Engine: Chit-chat, QnA and Dialogue

Ming Zhou

Microsoft Research Asia

编程之美2017 @Tsinghua University

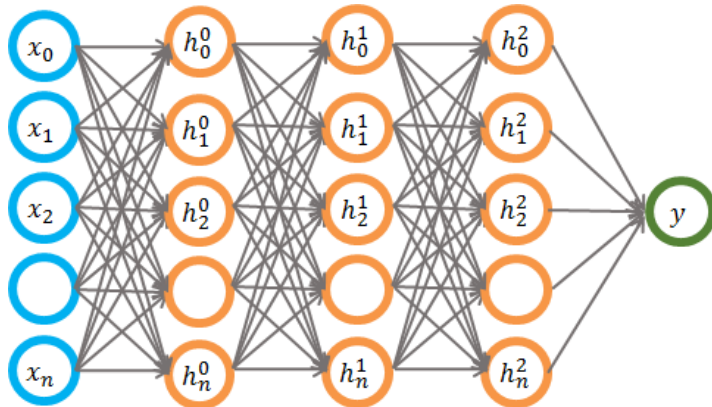
Conversation Engine Architecture



DNN4NLP Fundamentals

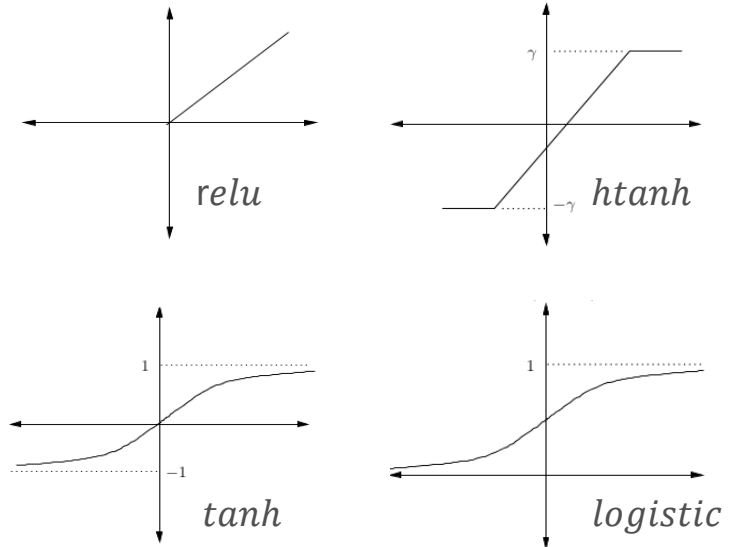
Deep Neural Network

- Deep Neural Network :
 - Involve multiple level neural networks
 - Non-Linear Learner



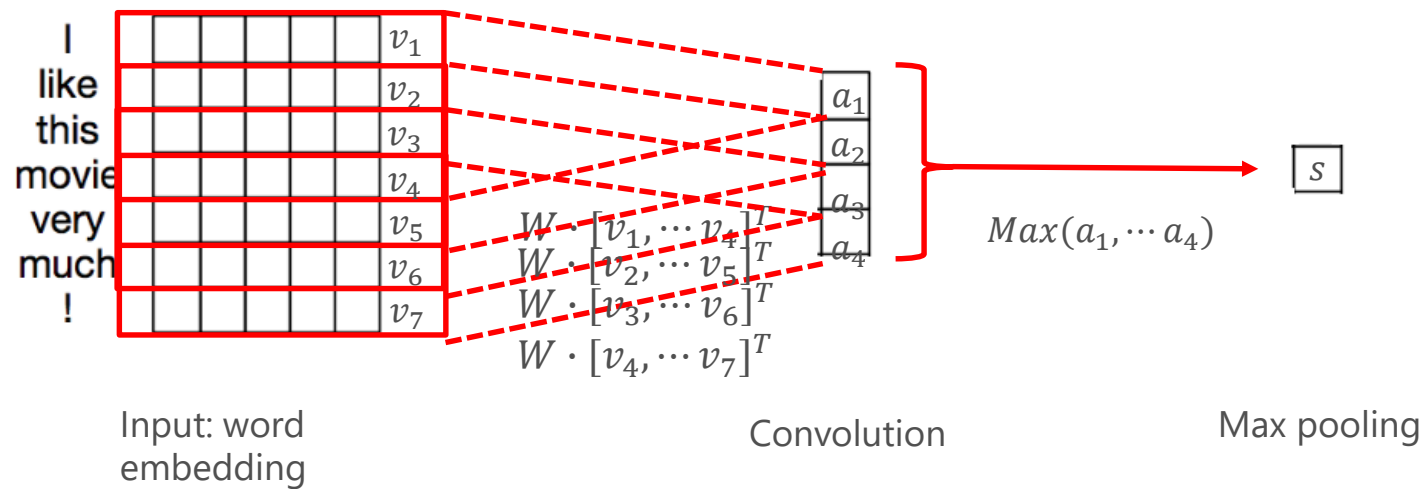
$$h^0 = f(w^0 x) \quad h^1 = f(w^1 h^0)$$

$$h^2 = f(w^2 h^1) \quad y = f(w^3 h^2)$$

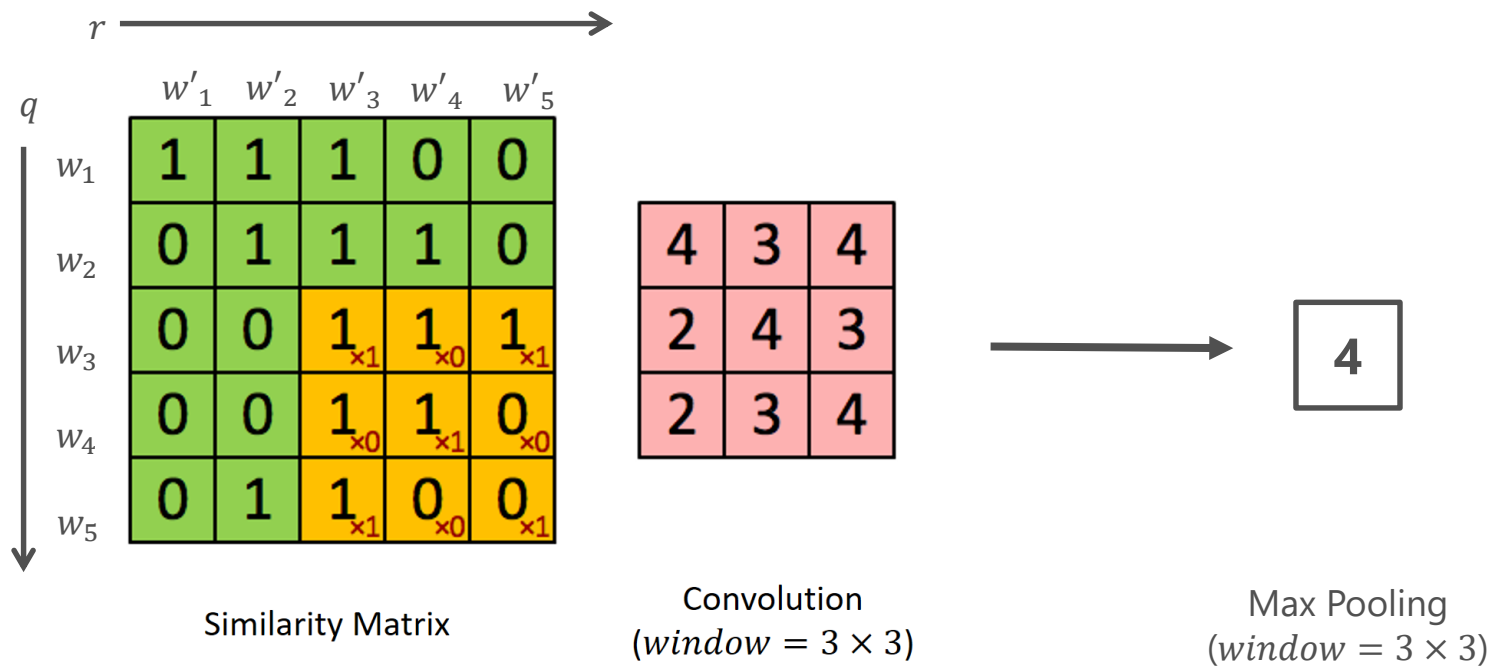


Active functions: $y = f(x)$

CNN (1D)

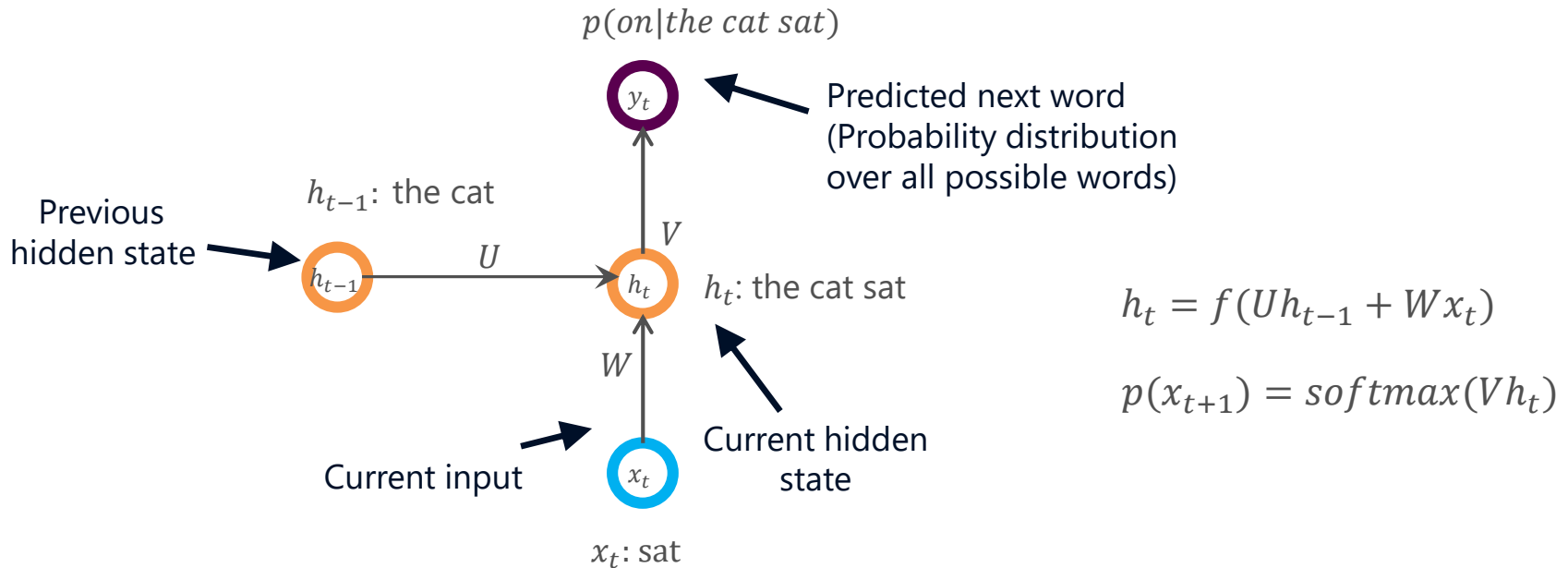


CNN (2D)

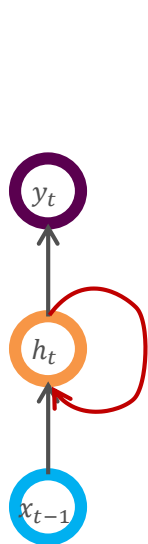


RNN(Recurrent Neural Network)

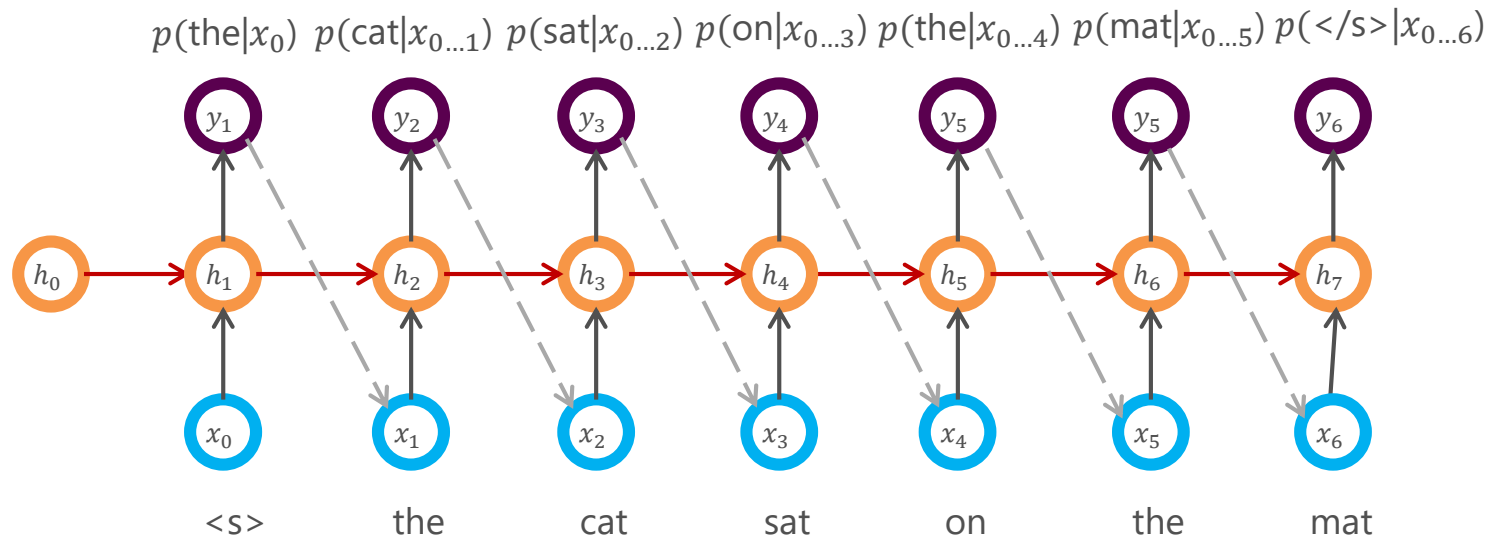
- Inputs: History h_{t-1} at time $t - 1$ and input x_t at time t
- Output: History h_t at time t and probability of next word y_t



RNN(Recurrent Neural Network)



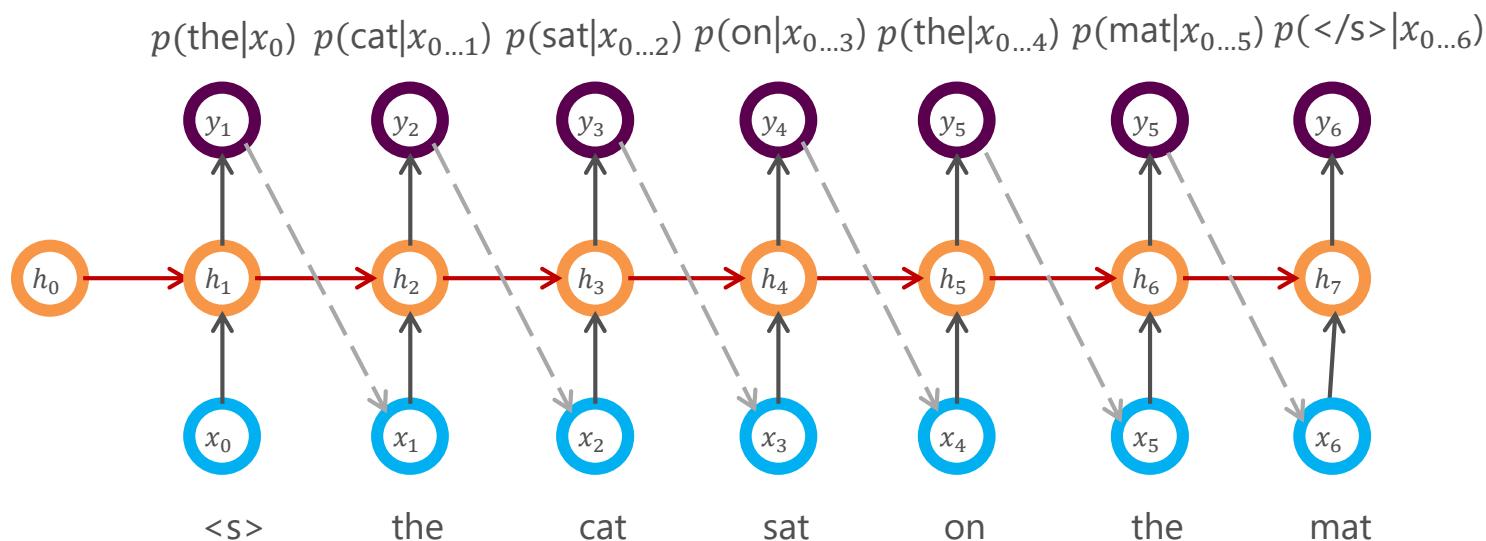
RNN



Unfolded RNN for Language Modeling

Sentence: The cat sat on the mat

Training of RNN



$$p(\text{the cat sat on the mat } \langle s \rangle) = \\ p(\text{the}|x_0) * p(\text{cat}|x_{0...1}) * p(\text{sat}|x_{0...2}) * p(\text{on}|x_{0...3}) \\ * p(\text{the}|x_{0...4}) * p(\text{mat}|x_{0...5}) * p(\text{</s>}|x_{0...6})$$

$$\text{Loss} = -\log(p(\text{the cat sat on the mat } \langle s \rangle)) \\ = -\sum_i \log(p(x_i|x_{0..i-1}))$$

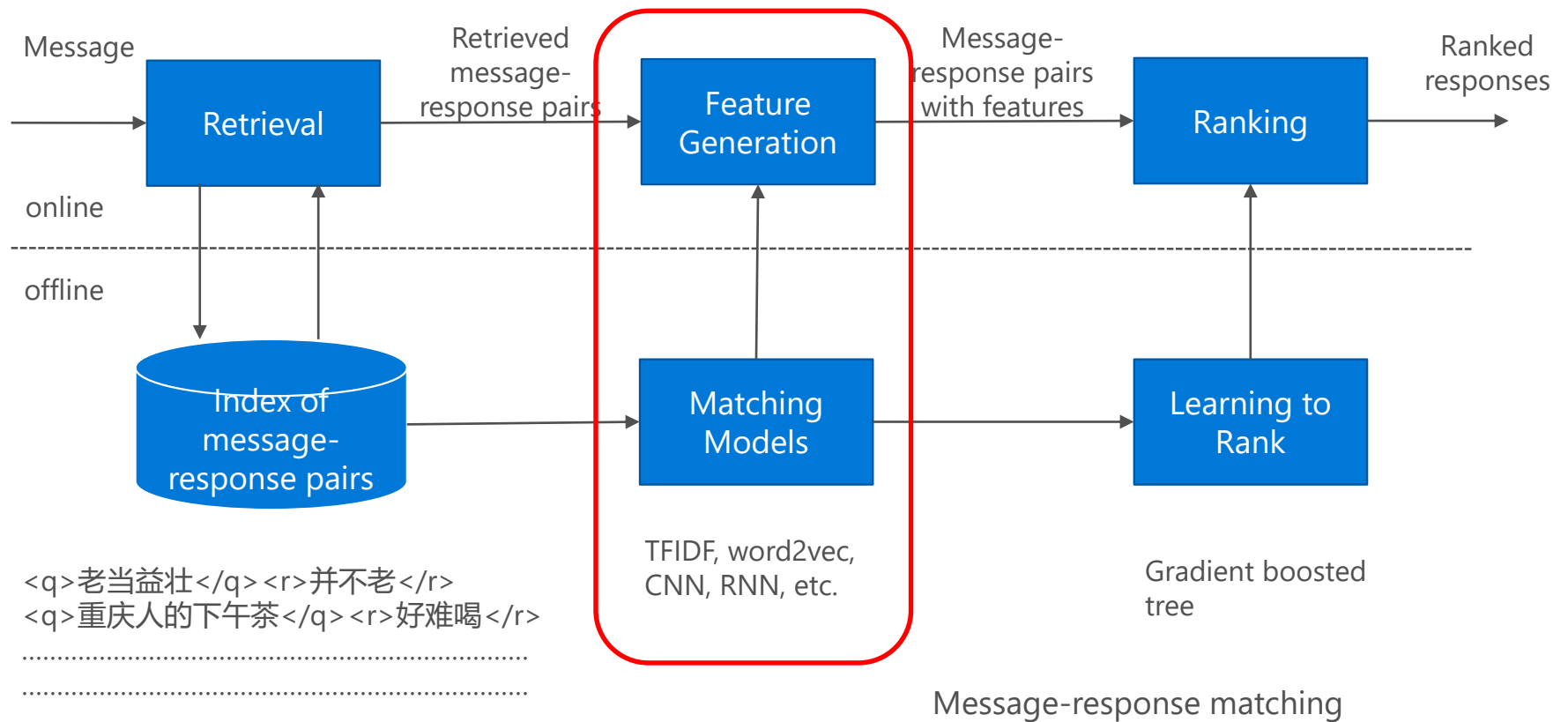
SGD Training

1. Sample a sentence S from corpus;
2. Build unfolded RNN for S ;
3. Run forward to compute **Loss**;
4. Run backward to compute gradient;
5. Update parameters;
6. Repeat until **perplexity** $< \epsilon$ or epoch $> n$;

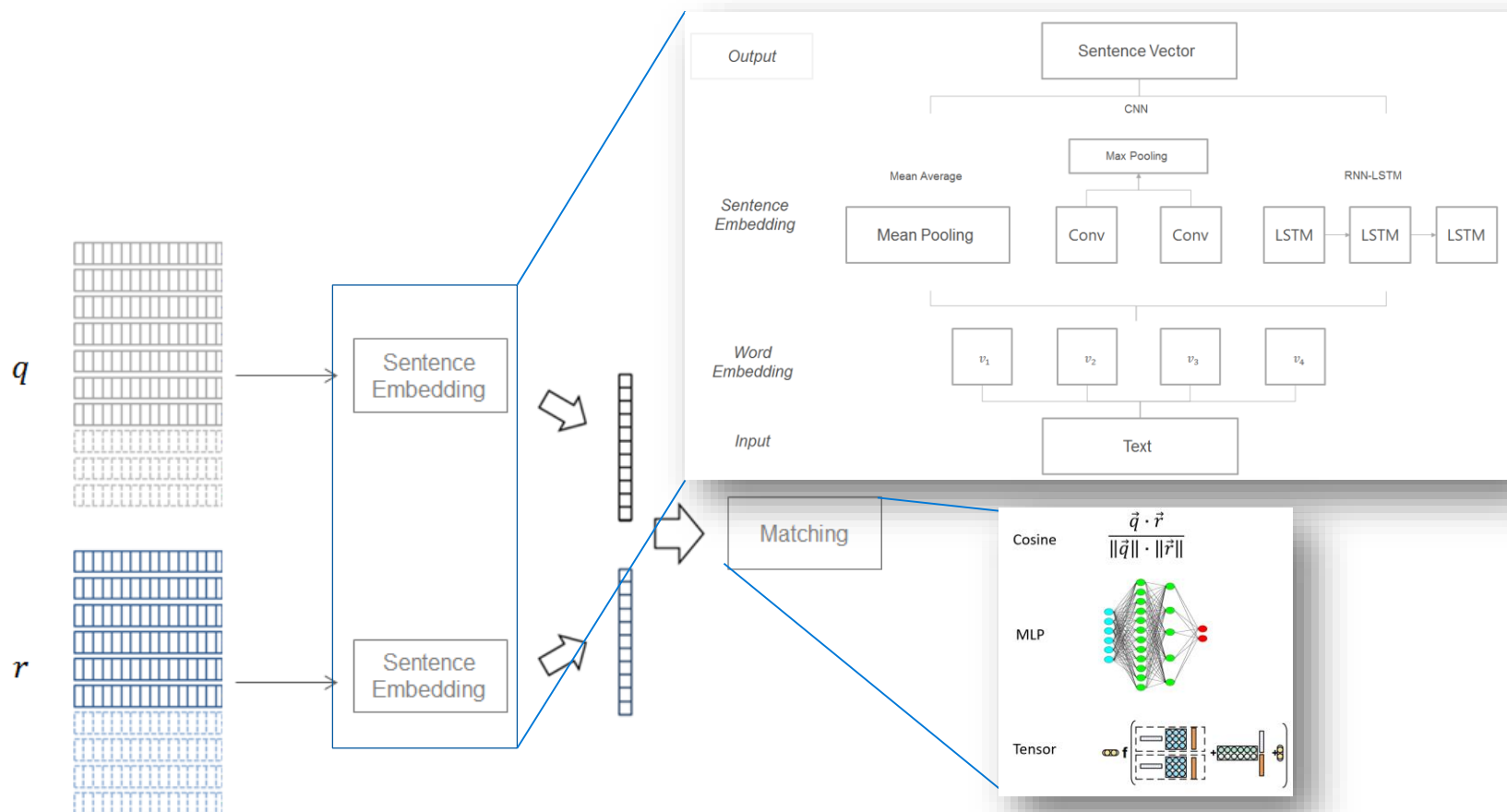
Skip the details of LSTM/GRU here

General Chat Engine

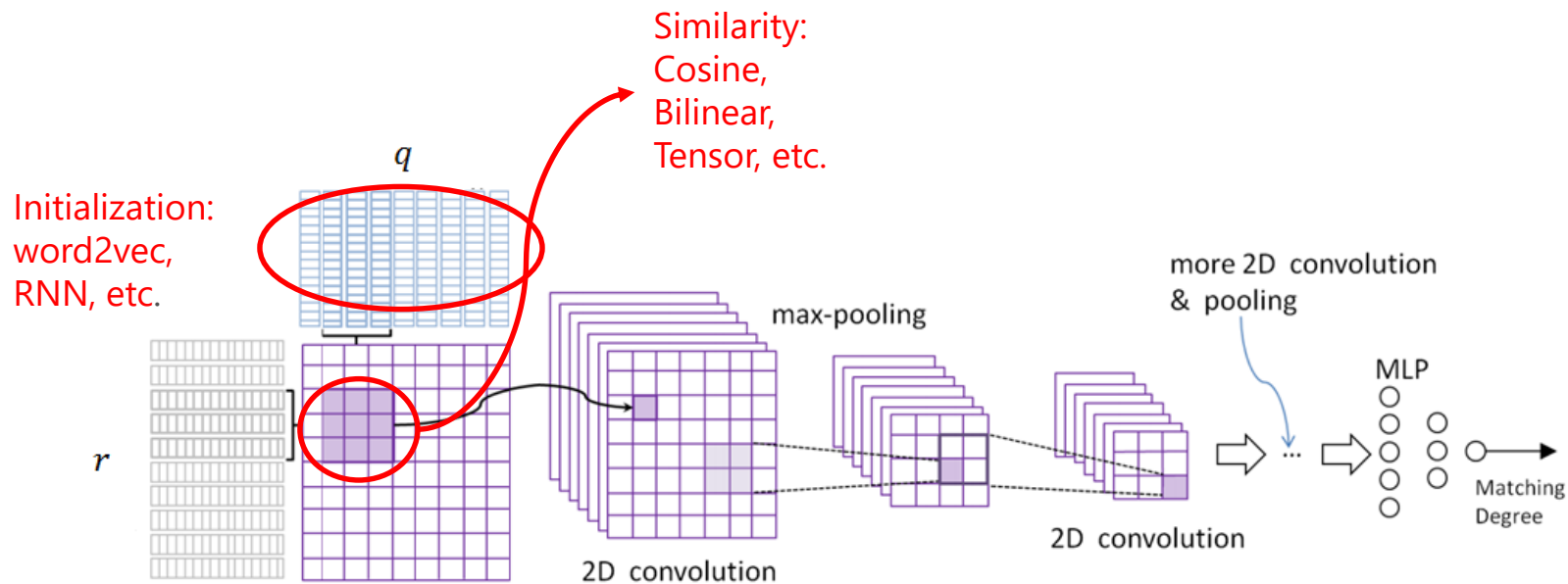
Architecture of Retrieval-based Chatbot (Single-Turn)



Basic Models for Message-Response Matching : Architecture I



Basic Models for Message-Response Matching : Architecture II



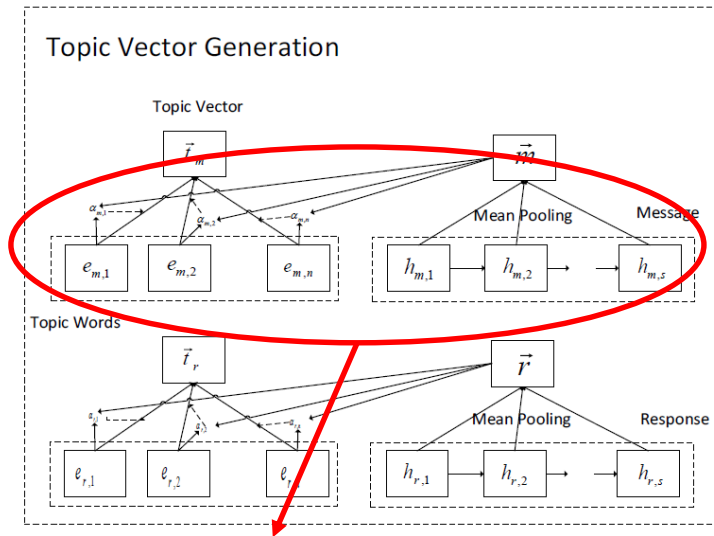
Baotian Hu et al. *Convolutional Neural Network Architectures for Matching Natural Language Sentences*, In NIPS'14

Liang Pang et al. *Text Matching as Image Recognition*, In AAAI'16

Shengxian Wan et al. *A Deep Architecture for Semantic Matching with Multiple Positional Sentence Representations*, In AAAI'16

Fusing with External Knowledge I

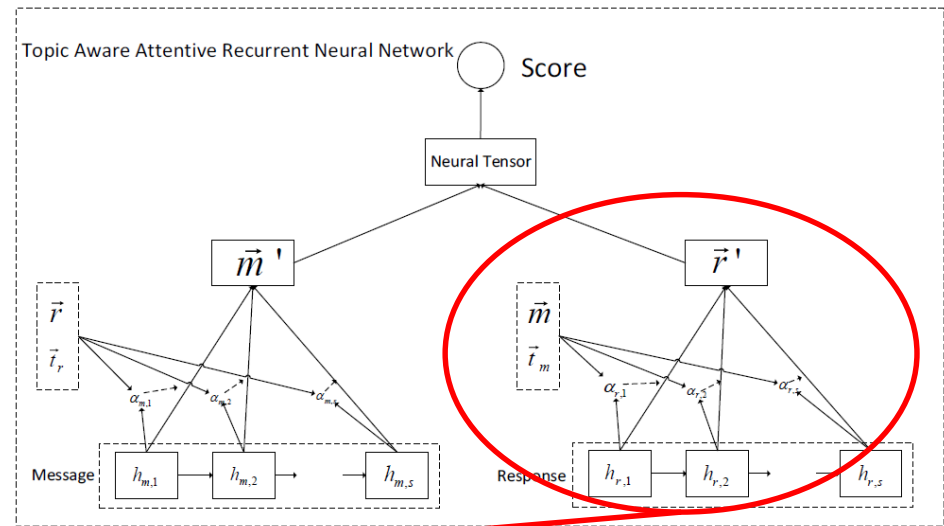
- Topic Aware Attentive Recurrent Neural Network (TAARNN)



Let message/response attend to important parts in external knowledge (topics)

$$\vec{t}_m = \vec{\alpha}_m \cdot T_m$$

$$\vec{\alpha}_m \propto T_m \cdot A \cdot \vec{m}$$



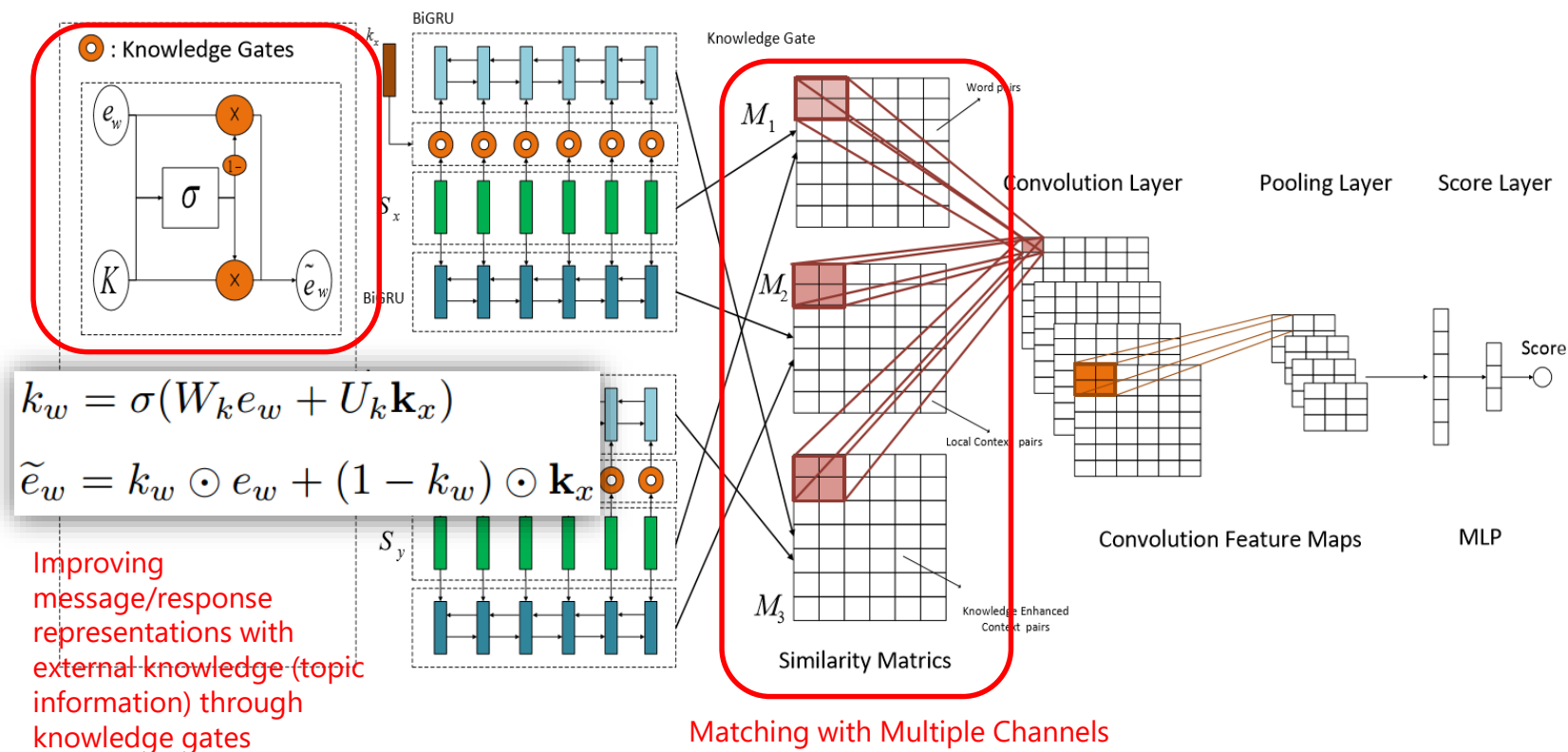
Joint attention with message/response and knowledge (topics)

$$\vec{r}' = \vec{\alpha}_r \cdot \vec{r}$$

$$\alpha_{r,i} \propto \tanh\left(\sum_j h_{r,i} \cdot A_2 \cdot h_{m,j} + h_{r,i} \cdot A_3 \cdot \vec{t}_m\right)$$

Fusing with External Knowledge II

- Knowledge Enhanced Hybrid Neural Network (KEHNN)



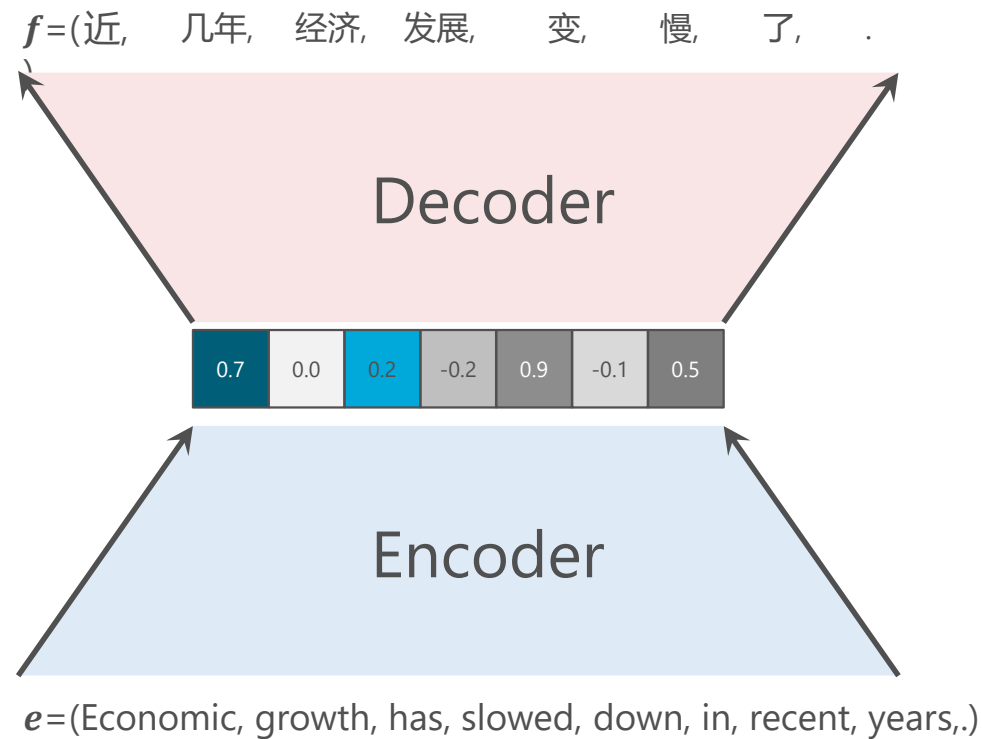
Evaluation on the Largest Public Data - Ubuntu Corpus

	$R_2@1$	$R_{10}@1$	$R_{10}@2$	$R_{10}@5$
Cosine (I)	68.1	38.3	48.2	68.6
MLP (I)	65.1	25.6	38	70.3
CNN (I)	66.5	22.1	36	68.4
CNN+Tensor (I)	74.3	34.9	51.2	79.7
LSTM (I)	72.5	36.1	49.4	80.1
CNN (II)	73.6	38	53.4	77.7
MatchPyramid (II)	74.3	42	55.4	78.6
MV-LSTM (II)	76.7	41.0	56.5	81.0
TAARNN (I)	77.0	40.4	56.0	81.7
KEHNN (II)	78.6	46	59.1	81.9

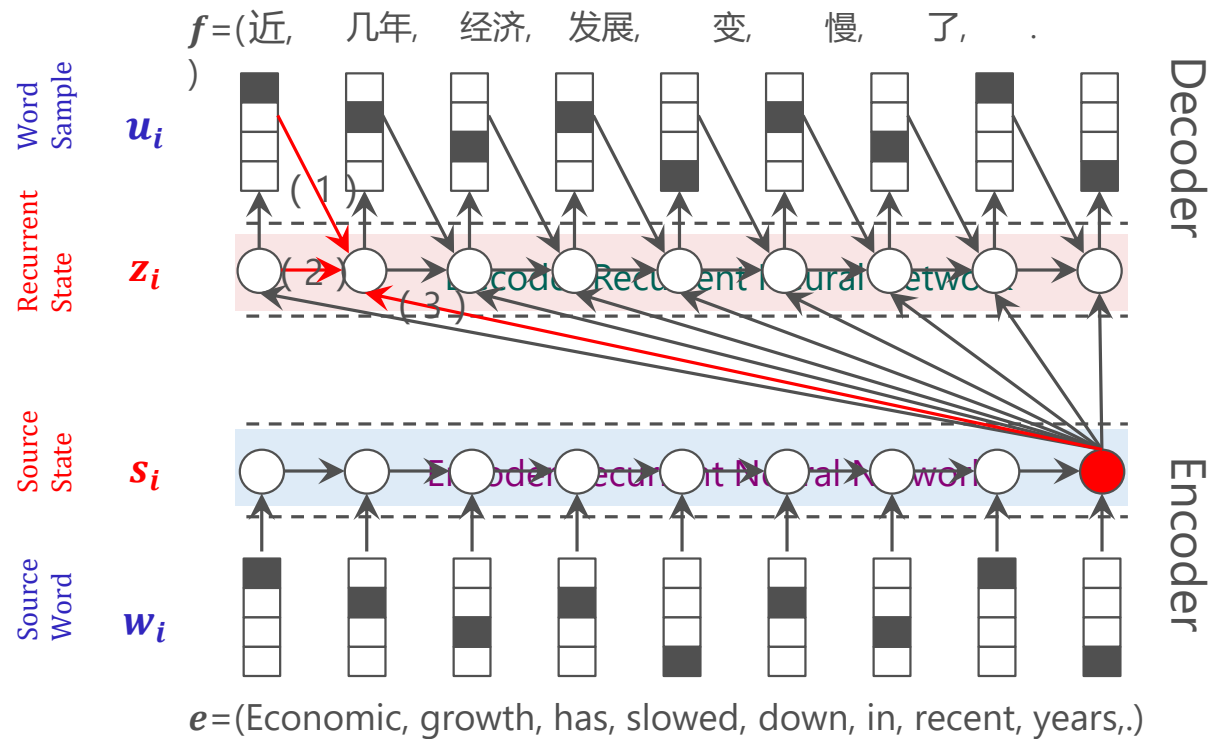
- Train : Validation : Test = 1M : 0.5M : 0.5M
- Negative examples are **randomly sampled**
- $R_n@k$ means recall at position k in n candidates
- (I) and (II) mean “in architecture I or II” respectively

Neural Generation Model

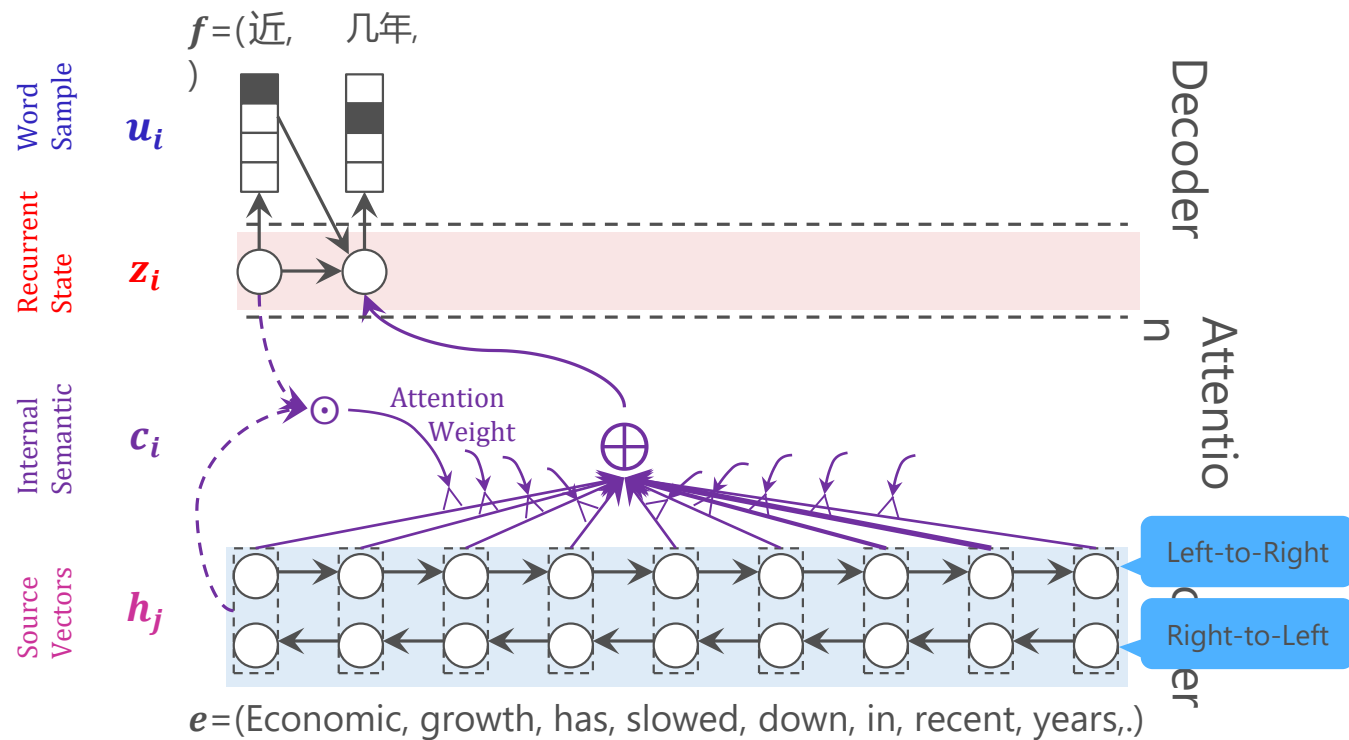
Encoder-Decoder for Sentence Generation



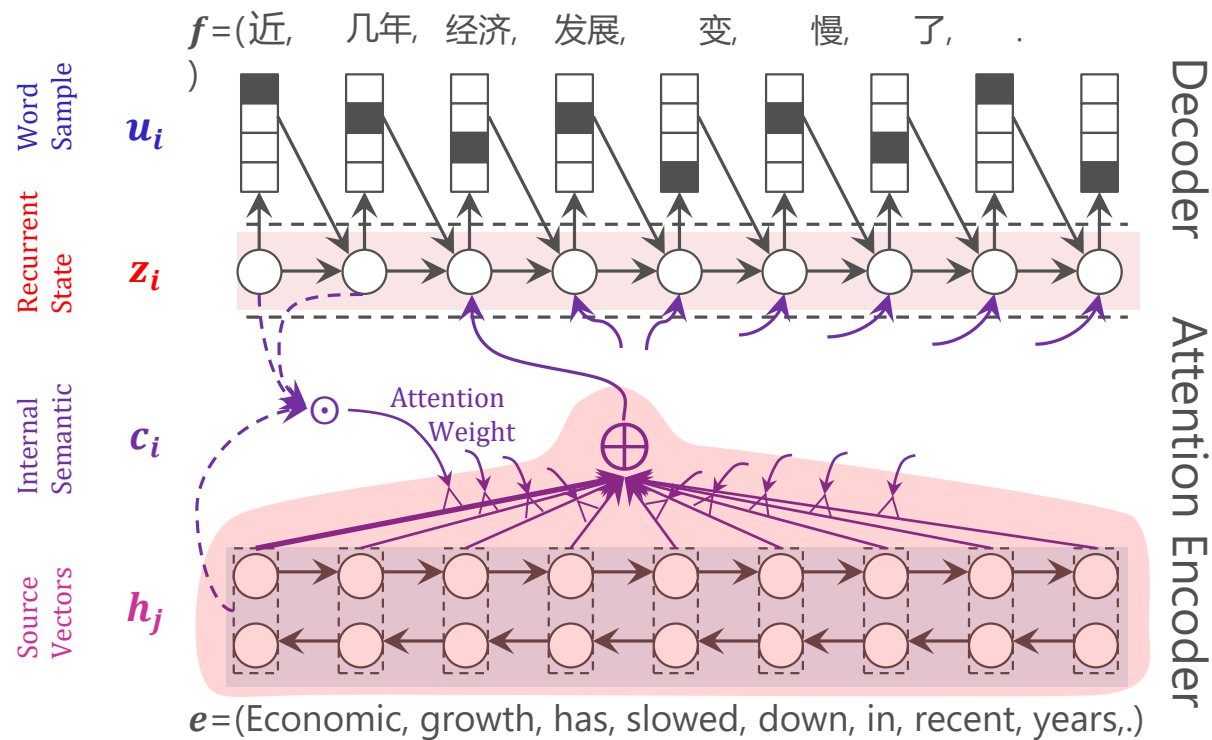
Encoder-Decoder for Sentence Generation



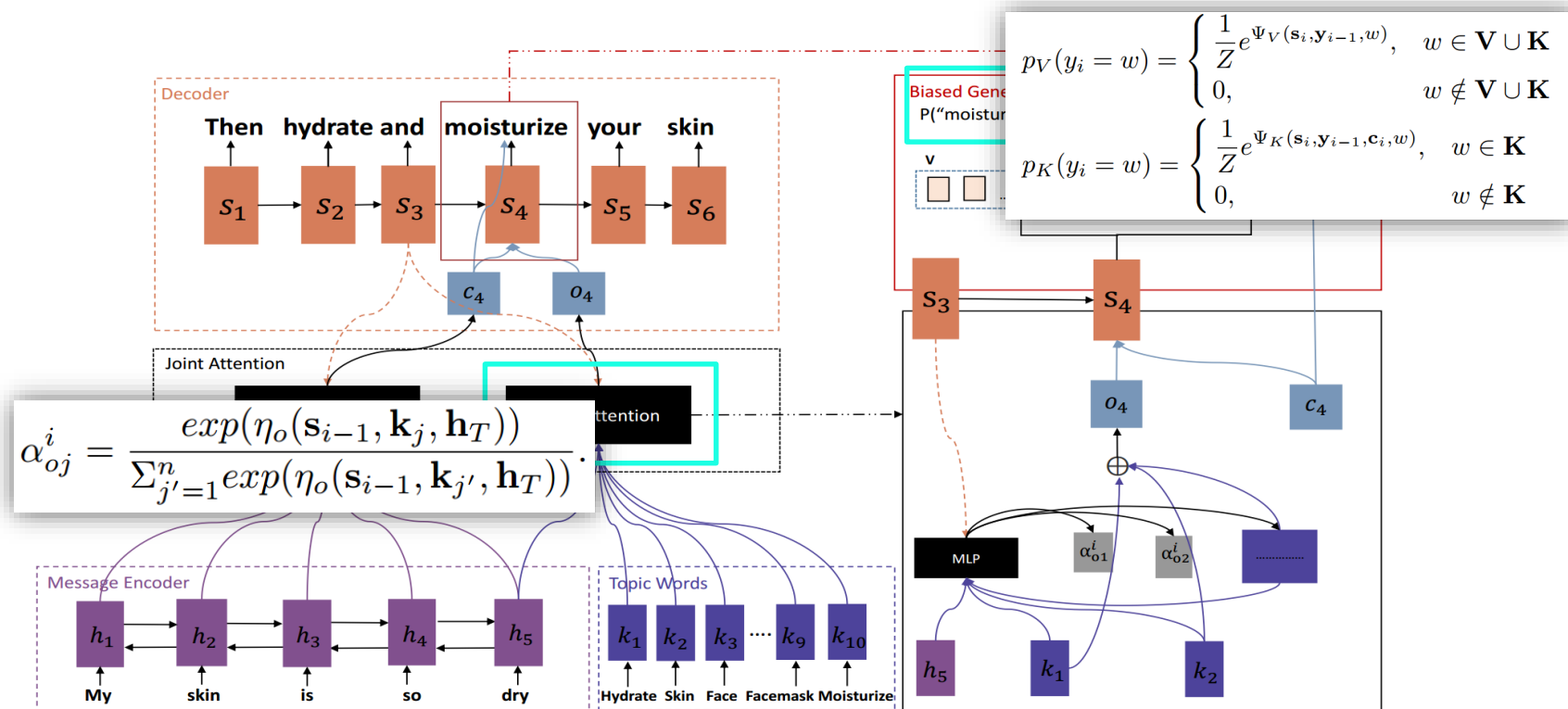
Attention based Encoder-Decoder



Attention based Encoder-Decoder



Topic-aware Neural Response Generation (TA-Seq2Seq)



Evaluation Results

Models	+2	+1	0	Kappa
S2SA	32.3%	36.7%	31.0%	0.8116
S2SA-MMI	33.1%	34.8%	32.1%	0.7848
S2SA-TopicConcat	35.9%	29.3%	34.8%	0.6633
S2SA-TopicAttention	42.3%	27.6%	30.0%	0.8299
TA-Seq2Seq	44.7%	24.9%	30.4%	0.8417

Table 1: Human annotation results

Models	PPL-D	PPL-T	distinct-1	distinct-2
S2SA	147.04	133.11	604/.091	1168/.207
S2SA-MMI	147.04	133.11	603/.151	1073/.378
S2SA-TopicConcat	150.45	132.12	898/.116	2197/.327
S2SA-TopicAttention	133.81	119.55	894/.106	2057/.277
TA-Seq2Seq	134.63	122.82	1355/.161	2970/.401

Table 2: Results on automatic metrics

- Data are crawled from Baidu Tieba
- Train : Validation :
Test = 5M : 10K : 1K

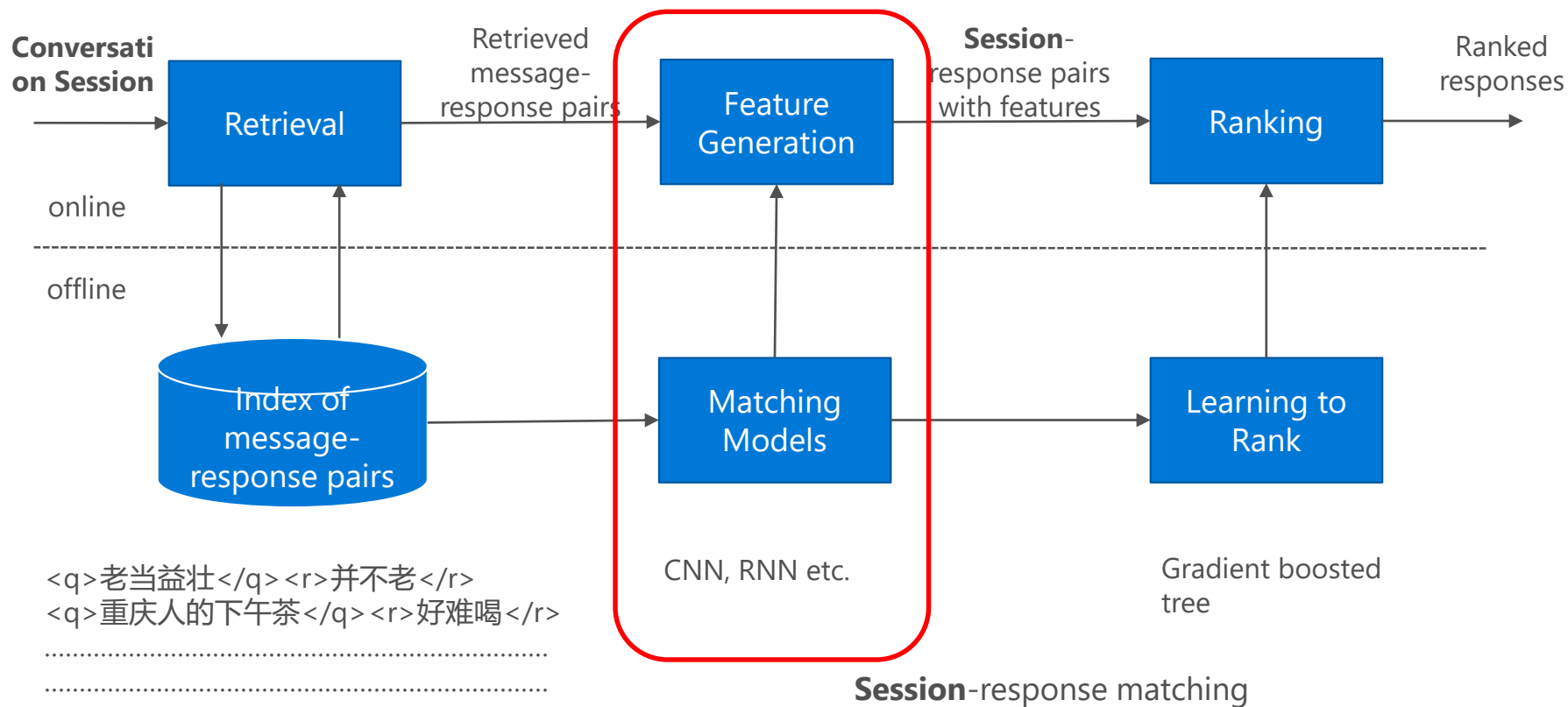
Examples

Message	TA-Seq2Seq	Baseline ¹	Baseline ²
		S2SA-MMI	S2SA
你也玩摄影 You also like photography	我不玩摄影，只是想拍个照片而已 I don't work on photography . I just want to take a few pictures .	额额额... Well well well....	额... Well...
我的左脸颊上有一些红血丝 There is some redness on my left cheek.	我也是敏感肌 My skin is also very sensitive .	我也是啊 Me too.	啊啊啊啊 Oh my god oh my god
大学生的话可以去实习吗 Can a college student apply for an internship	可以的，如果你着急的话可以先填个申请表 Yes, you can fill an application form first if you are in a hurry.	可以的可以的 Yes, you can.	可以的 Yes.
我的皮肤好干 My skin is so dry.	那就补水保湿吧 Then hydrate and moisturize your skin.	我也是啊 Me too.	啊啊啊啊 Oh my god oh my god.

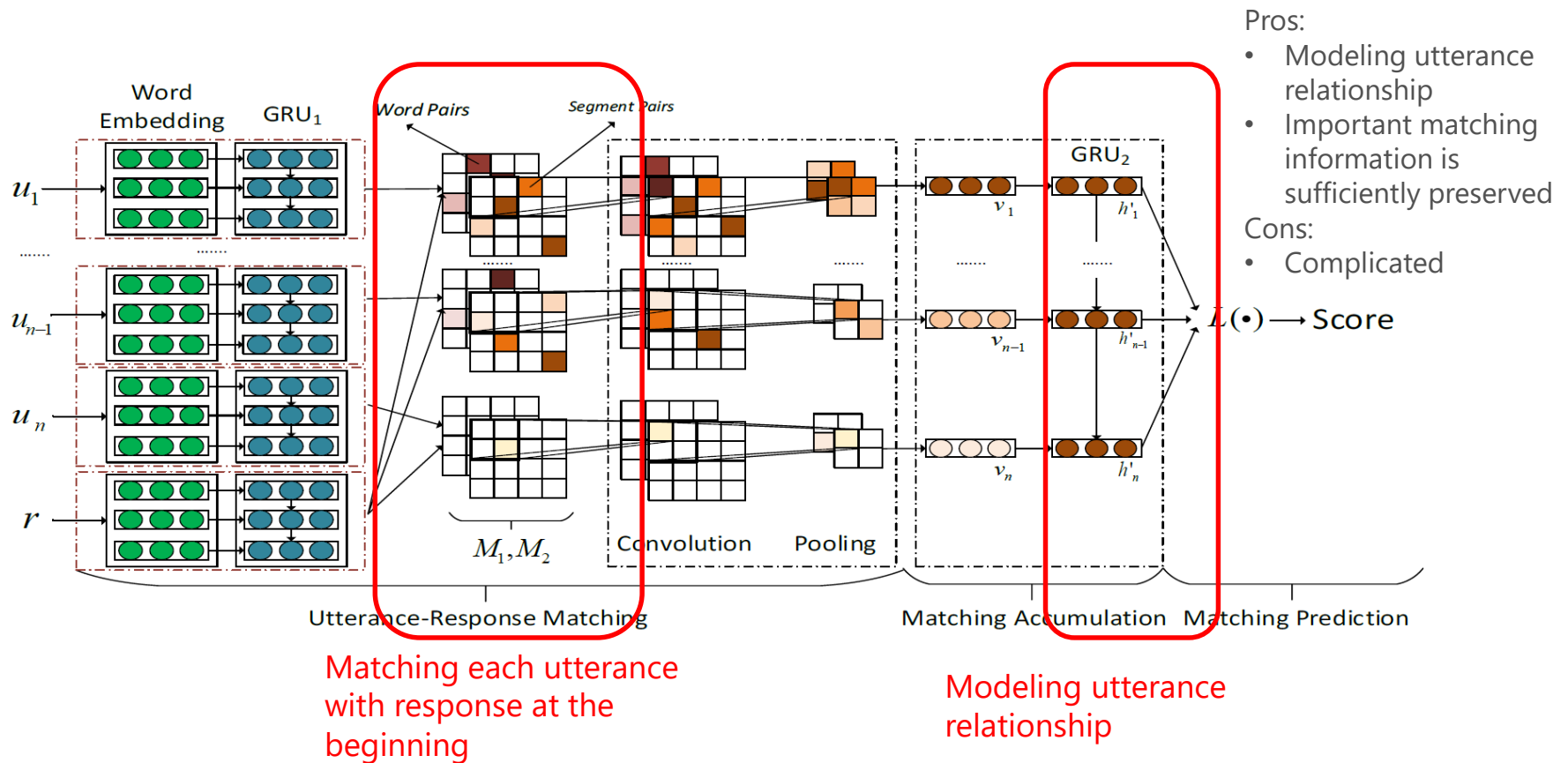
* Words in bold are topic words

Multi-Turn Conversation

System Overview



Session-Response Matching : Sequential Matching Network

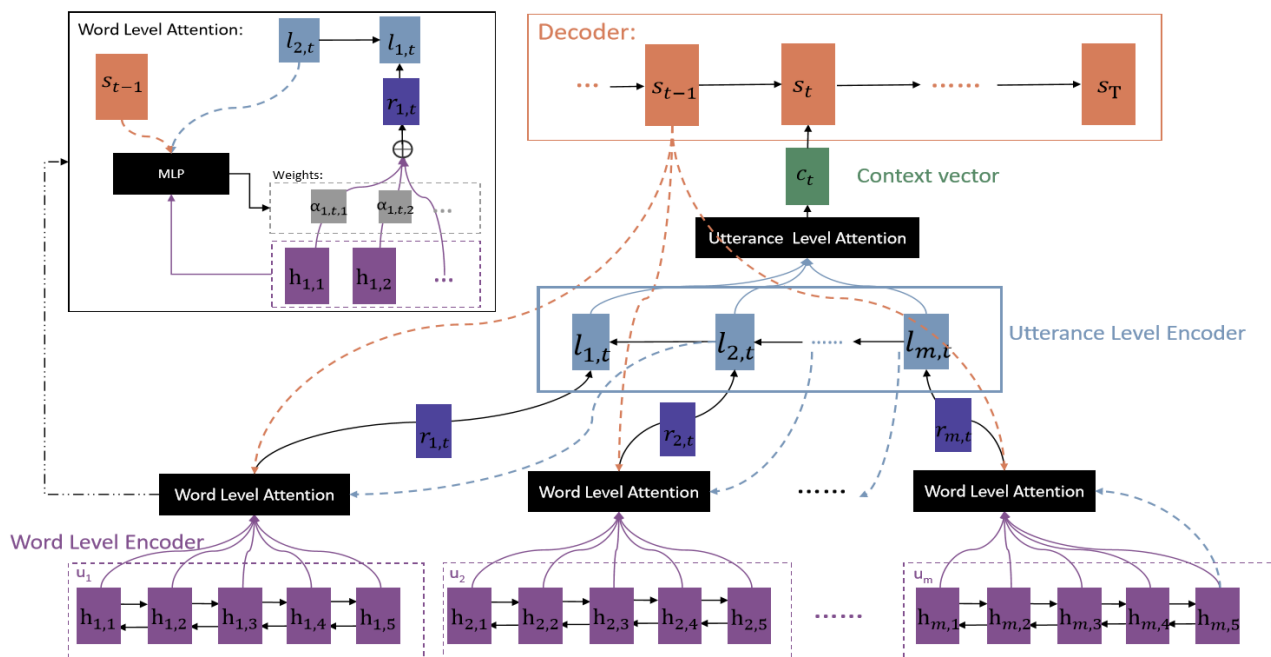


Evaluation on the Ubuntu Corpus

	$R_2@1$	$R_{10}@1$	$R_{10}@2$	$R_{10}@5$
Cosine	65.9	41.0	54.5	70.8
CNN	84.8	54.9	68.4	89.6
LSTM	90.1	63.8	78.4	94.9
BiLSTM	89.5	63.0	78.0	94.4
Multi-View	90.8	66.2	80.1	95.1
DL2R	89.9	62.6	78.3	94.4
SMN	92.6	72.6	84.7	96.1

- Train : Validation : Test = 1M : 0.5M : 0.5M
- Negative examples are **randomly sampled**
- $R_n@k$ means recall at position k in n candidates

Multi-turn Response Generation: Hierarchical Recurrent Attention Network



Chen Xing et al., *Hierarchical Recurrent Attention Network for Response Generation*, In arxiv

Evaluation Results

Model	Validation Perplexity	Test Perplexity
S2SA	43.679	44.508
HRED	46.279	47.467
VHRED	44.548	45.484
HRAN	40.257	41.138

Table 1: Perplexity results

Models	Win	Loss	Tie	Kappa
HRAN v.s. S2SA	27.3	20.6	52.1	0.37
HRAN v.s. HRED	27.2	21.2	51.6	0.35
HRAN v.s. VHRED	25.2	20.4	54.4	0.34

Table 2: Human annotation results (in %)

- Data are crawled from Douban Group
- Train : Validation : Test = 1M : 10K : 1K
- Kappa measures the agreement of human labeling

HRED: Iulian V Serban et al., *Building End-to-End Dialogue Systems Using Generative Hierarchical Neural Network Models*, In AAAI'16
VHRED: Iulian V Serban et al., *A Hierarchical Latent Variable Encoder-Decoder Model for Generating Dialogues*, In AAAI'17

Case Study

ID	Context	Response
1	<p>求唇部死皮怎么去⇒用化妆水一遍遍的擦嘴唇就能擦下来⇒好的回去试试⇒我是能擦下来</p> <p>How to exfoliate my lips?⇒Wipe your lips with lotion again and again. ⇒OK, let me try it. ⇒It works for me.</p>	<p>HRAN: 你用的什么牌子的化妆水</p> <p>Which brand of lotion are you using?</p> <p>VHRED: 嗯我知道了</p> <p>OK I see.</p>
2	<p>我不能去，你可以找丽丽陪你吃饭⇒她住松江，离市区太远了⇒有好吃的多远都要去啊⇒你为什么不能来呢</p> <p>I can't have dinner with you, but you can ask Lily instead.⇒She lives in Songjiang, too far away from downtown.⇒I think she'd like to go because of the delicious food.⇒ Why can't you come?</p>	<p>HRAN: 吃多了上火了</p> <p>Because I've eaten too much and gotten inflamed.</p> <p>VHRED: 因为我是男的</p> <p>Because I'm male.</p>
3	<p>征男友，160cm的妹子真的找不到男友吗⇒你找不到一定不是因为160⇒我知道脸也是硬伤嘛⇒是你非要175以上</p> <p>I want a boyfriend. Why can't a 160cm girl find a boyfriend?⇒It's definitely not because you are 160cm.⇒ Well I know I'm not good-looking.⇒No, it's because you always hit on someone higher than 175cm.</p>	<p>HRAN: 身高不是硬性要求</p> <p>No, I don't care much about height.</p> <p>VHRED: 你是男的还是女的啊</p> <p>Are you male or female?</p>
4	<p>啦啦啦销售助理在哪里呀在哪里⇒学历不达标⇒对做销售有兴趣么⇒没做过不过我不挑工作</p> <p>I'm looking for a sales assistant.⇒I might not be qualified because of my education.⇒Are you interested in sales? ⇒I have no experience before but I'm open to any jobs.</p>	<p>HRAN: 那你做过设计么</p> <p>Then have you done any designing jobs before?</p> <p>VHRED: 好吧好吧我也觉得</p> <p>OK, I think so.</p>

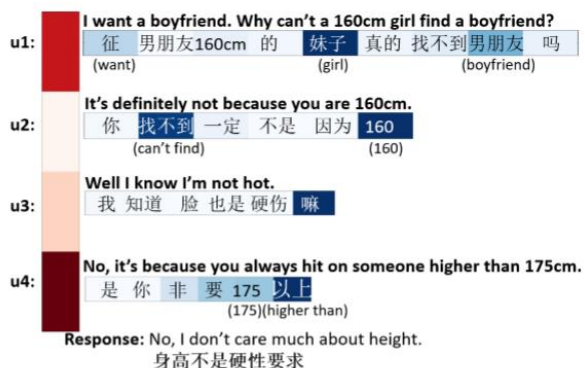
Visualization



(a) Visualization of case 1



(b) Visualization of case 2



(c) Visualization of case 3

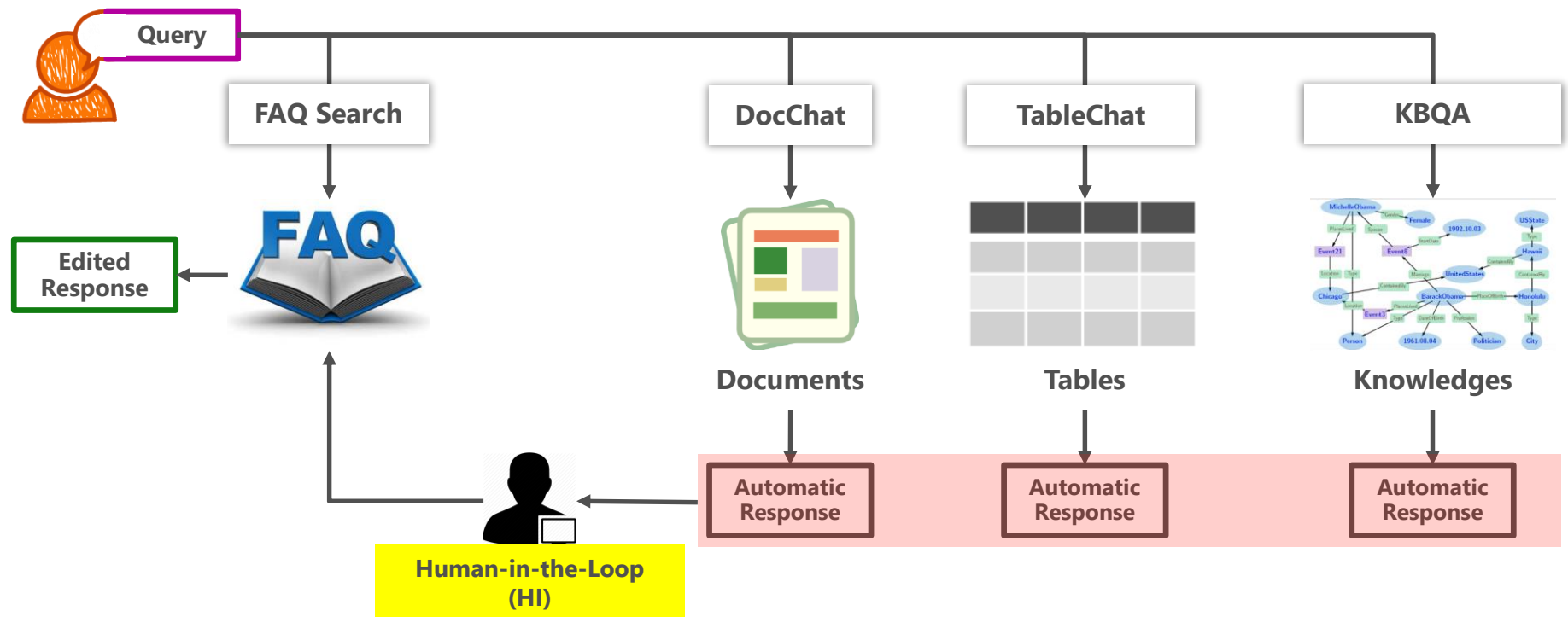


(d) Visualization of case 4

Information and Answer(InfoBot)

InfoBot Overview

Build conversational bots for information queries based on various genre of content and knowledge



Knowledge-based QA (KB-QA)

CCG: Combinatory Categorical Grammar
DCS: Dependency-based Compositional Semantics
SMT: Statistical Machine Translation

Semantic Parsing-based KB-QA(SP-QA)

Query: where was Barack Obama born ?



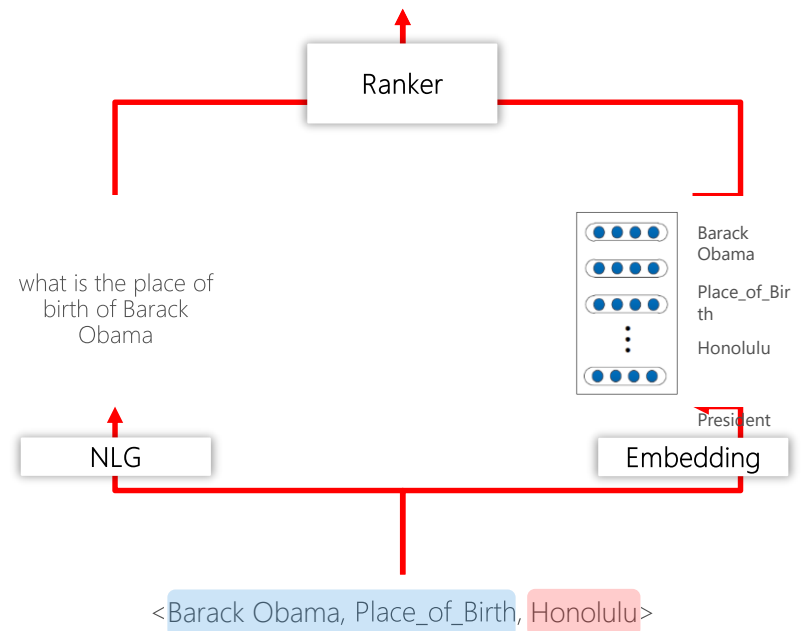
Logical Form: $\lambda x. Place_of_Birth(Barack\ Obama, x)$

KB Lookup

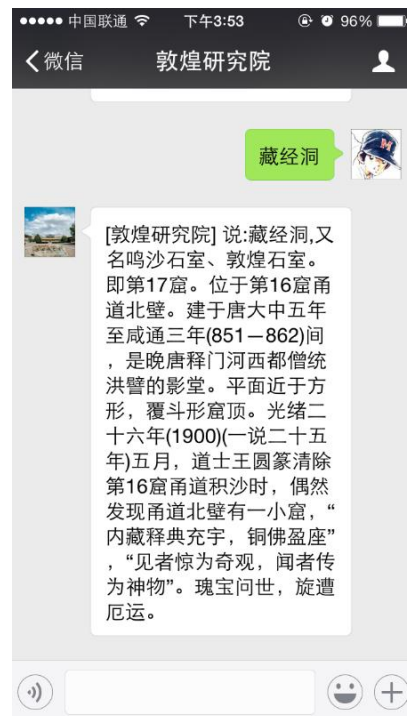
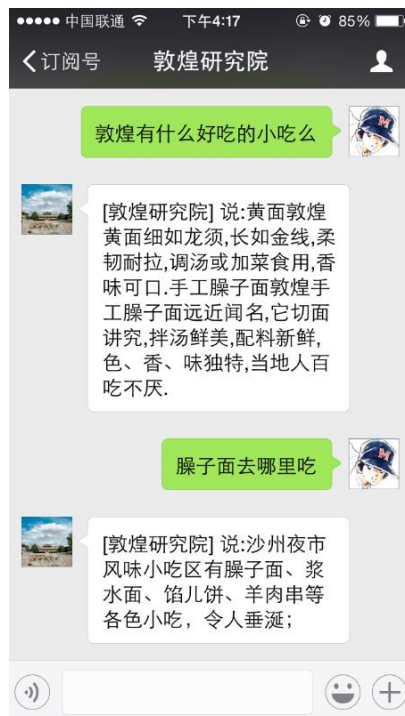
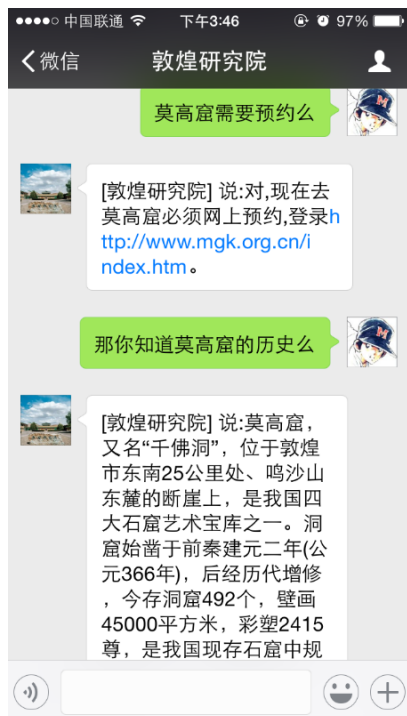
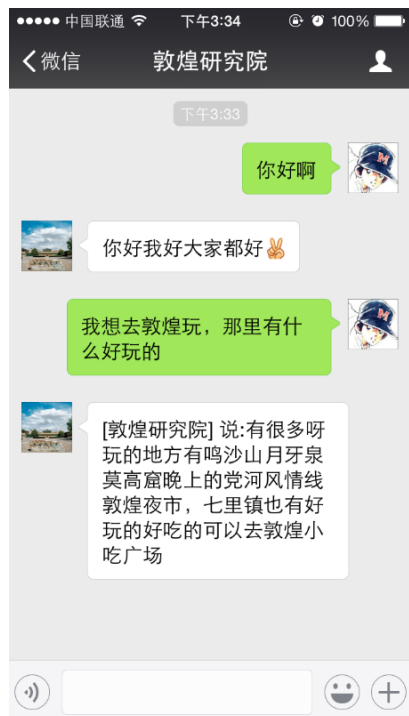
<Barack Obama, Place_of_Birth, Honolulu>

Information Retrieval-based KB-QA(IR-QA)

Query: where was Barack Obama born ?

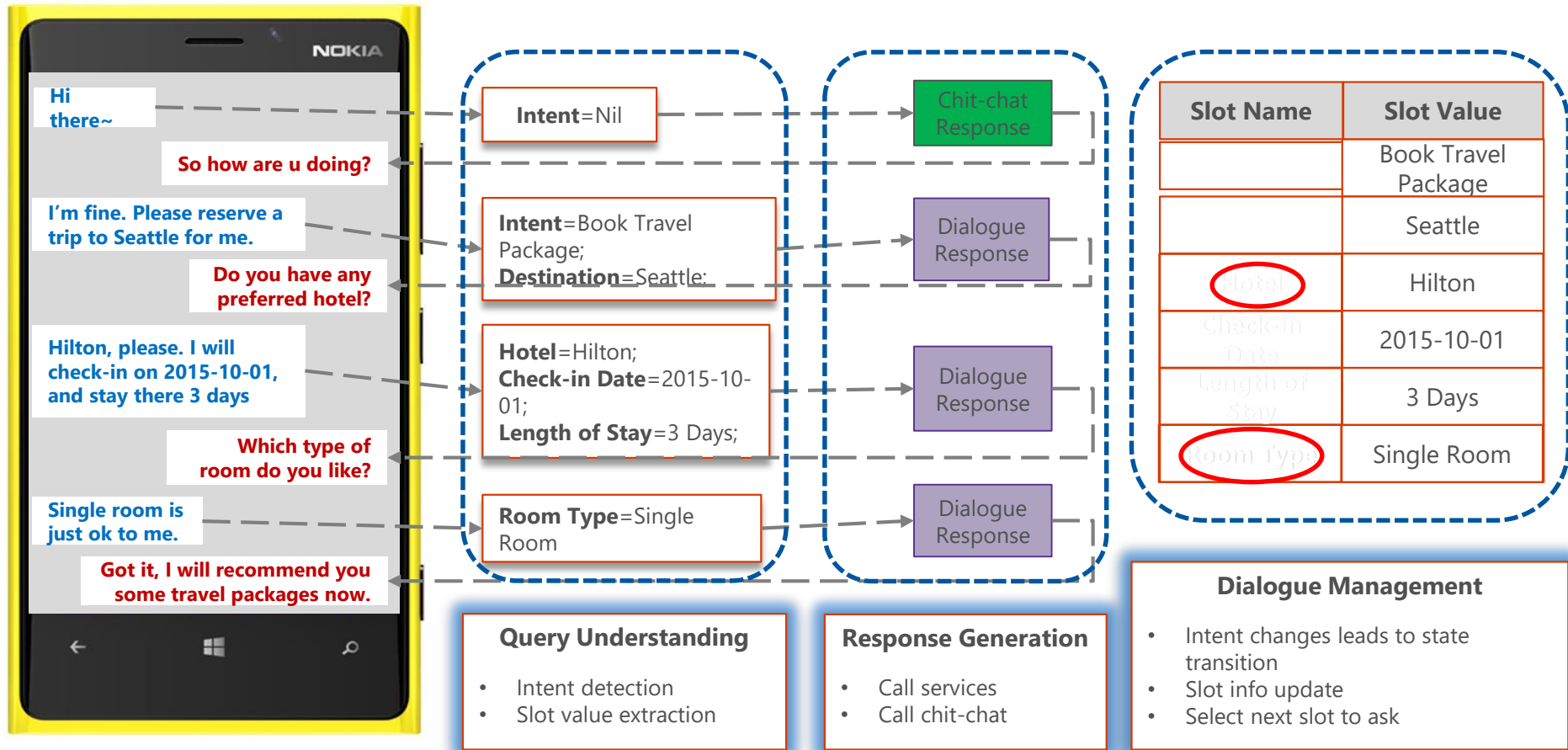


敦煌小冰



Task-Oriented Dialogue System

Dialogue Process



Conclusion and Future Work

CAAP

- Chit-chat: Retrieval-based and generation-based, single turn and multi-turn, fusing with knowledge
- Infobot based on FAQ, Doc-Chat, KB-QA
- Task-Oriented dialogue

Future work

- Better modelling the context and multi-turn
- Fusing user (type) profile into response generation
- Response style transformation