

Scalable Neural Language Generation for Open Domain Dialogue Systems

Speaker: Tsung-Hsien Wen

Supervisor: Professor Steve Young

Problem Definition

- Given a meaning representation, map it into a natural language representation
 - `inform(type=Seven_days,food=Chinese)`
 - Seven_days serves good Chinese food.
- What we care about?
 - adequacy, fluency, readability, variation (Stent et al 2005)

Motivation

- Traditionally, NLG is not scalable because :
 - Embrace a rule-based regime
 - Highly specialised for in-domain applications
- Talking to NLG is not enjoyable because of :
 - Frequent repetition of certain output forms
 - Awkward responses that are not colloquial

Why RNN for NLG?

- Elegant structure for modeling sequences.
- Flexible architecture for adding auxiliary information.
- Collecting data is convenient and quick (crowdsourcing).
- More human-like and colloquial.
- No expert knowledge is required.
- Extensible, adaptation techniques exist.
- Distributed representation
- Less cost, quicker development cycle
- End-to-End trainable

Challenges

- How to render the exact information we want (with the existence of language variation)?
- Adopted methods:
 - Overgeneration – Reranking paradigm (Oh and Rudnicky 2000)
 - Sample words from a Recurrent Generation Model output.
 - Select top candidates based on some scoring criteria.

Part 1

Heuristically Gated RNN Generator

Outline

- Recurrent Generation Model
- Convolutional Semantic Reranker
- Backward RNN Reranker
- Experiments
 - Setup
 - Automatic Evaluation
 - Human Evaluation

Outline

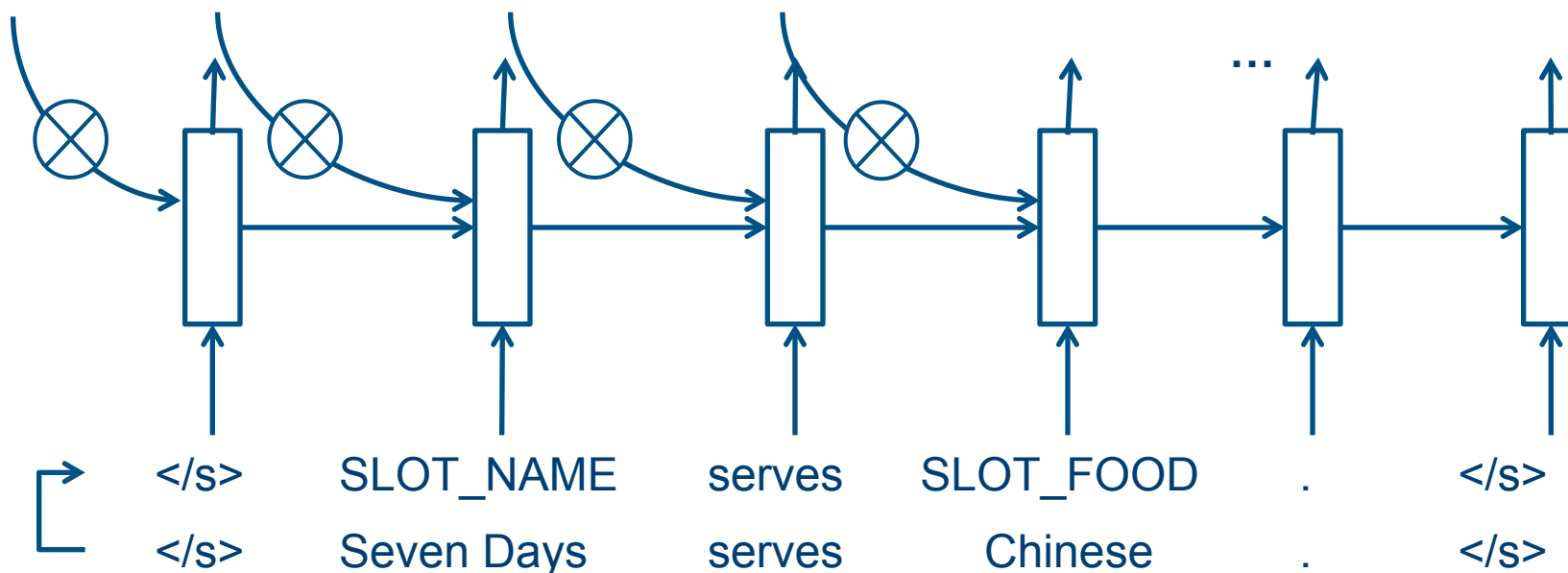
- **Recurrent Generation Model**
- Convolutional Semantic Reranker
- Backward RNN Reranker
- Experiments
 - Setup
 - Automatic Evaluation
 - Human Evaluation

Recurrent Generation Model (1/2)

Inform(name=Seven_Days, food=Chinese)

{ 0, 0, 1, 0, 0, ..., 1, 0, 0, ..., 1, 0, 0, 0, 0, 0... }

dialog act 1-hot representation



delexicalisation

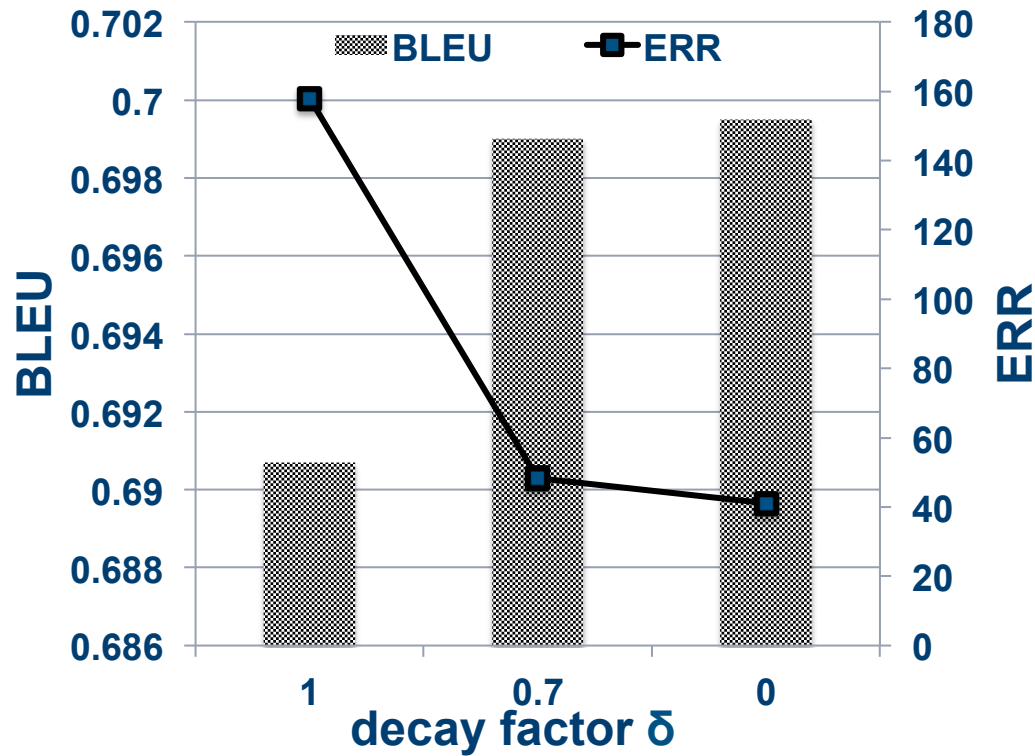
(Mikolov et al 2010)

Recurrent Generation Model (2/2)

- Heuristically check (**exact match**) whether a given slot token has been generated.
- Apply a decay factor $\delta < 1$ on generated feature values.
- Use features to configure the network NOT to re-generate slots that have already generated.
- Binary slots and don't care values cannot be handled.

Feature value	</s>	SLOT_NAME	serves	SLOT_FOOD	.	</s>
NAME	1	1	δ	δ^2	δ^3	δ^4
FOOD	1	1	1	1	δ	δ^2

Recurrent Generation Model (3/3)



- ERR: # of missing/redundant slots
- BLEU: BLEU-4 against multiple references

Outline

- Recurrent Generation Model
- **Convolutional Semantic Reranker**
- Backward RNN Reranker
- Experiments
 - Setup
 - Automatic Evaluation
 - Human Evaluation

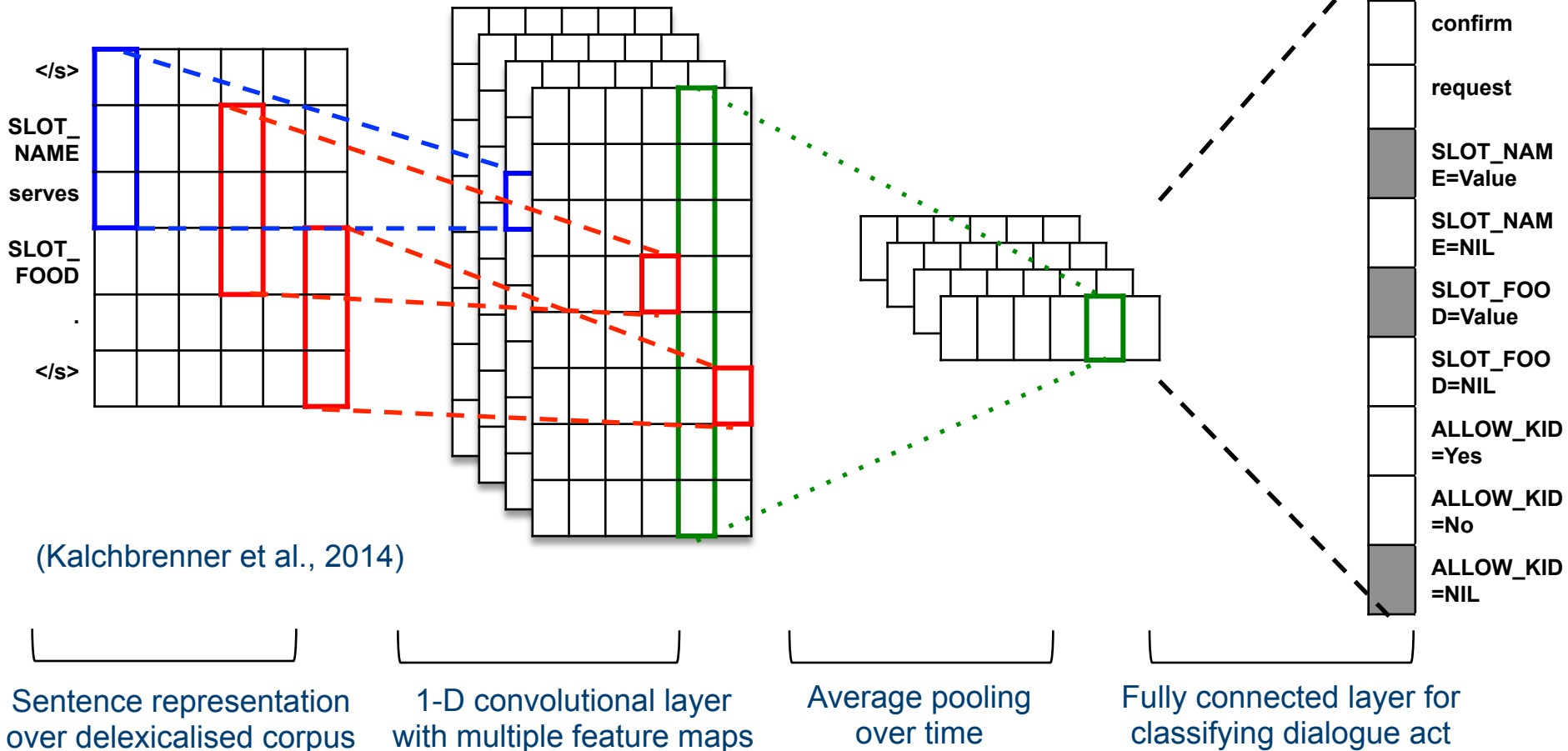
Convolutional Semantic Reranker (1/2)

- Designed to handle :
 - Binary slots: ALLOW_KID=yes/no
 - “don’t care” values: AREA=dont_care
- Use CNN for semantic validation

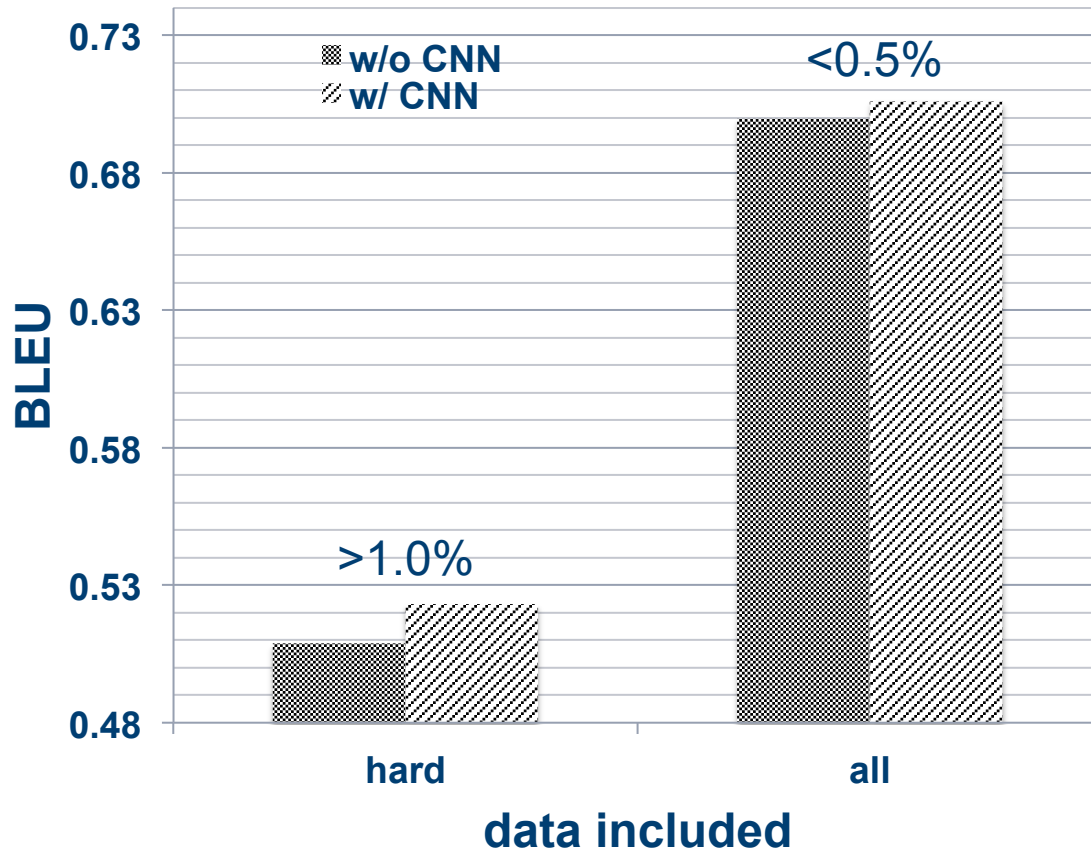
Convolutional Semantic Reranker (2/2)

Target dialogue act: inform(name=Seven_days, food=Chinese)

Generated candidate: </s> SLOT_NAME serves SLOT_FOOD . </s>



Convolutional Semantic Reranker (3/3)



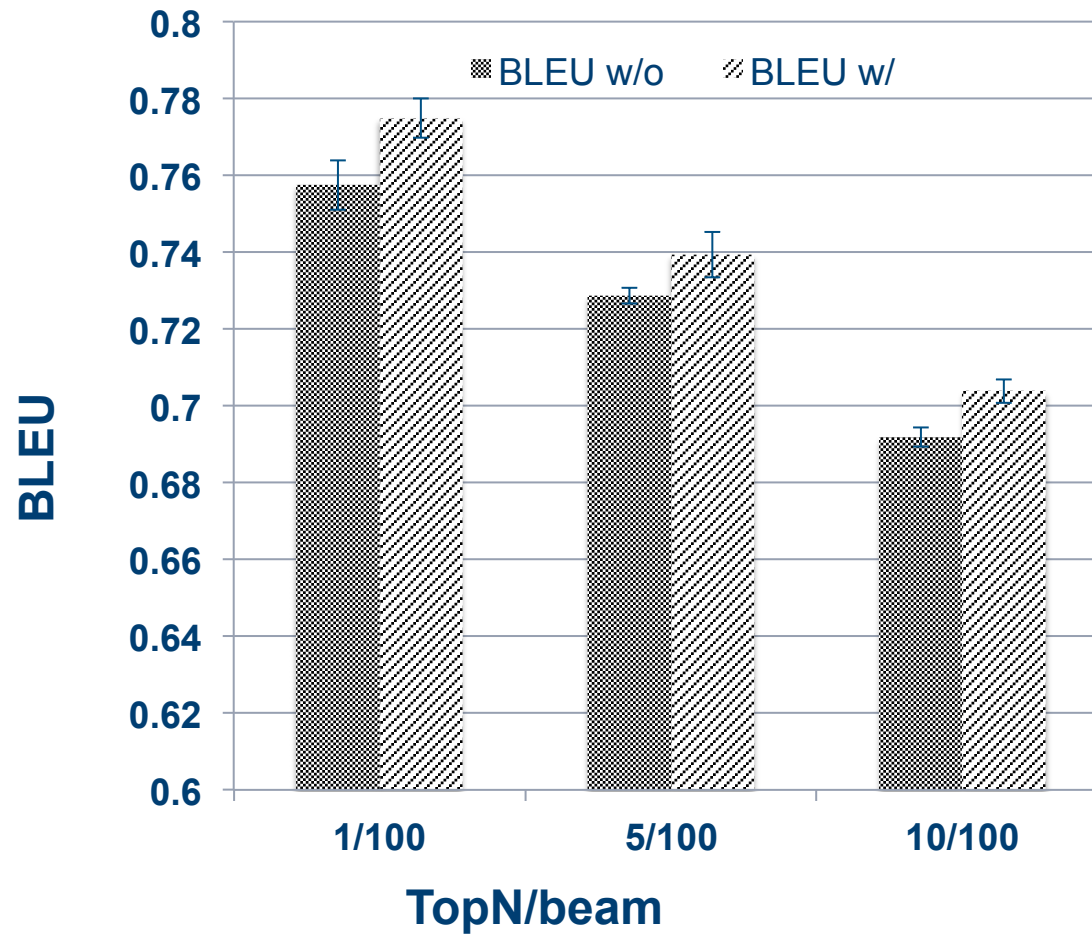
Outline

- Recurrent Generation Model
- Convolutional Semantic Reranker
- **Backward RNN Reranker**
- Experiments
 - Setup
 - Automatic Evaluation
 - Human Evaluation

Backward RNN Reranker

- Motivation:
 - Considering backward context can reduce grammatical errors.
 - Ex. “*Seven Days is **an** exceptional restaurant.*”
- Integrating information from both directions is tricky.
 - The generation procedure is sequential in one direction only.
- Alternative => train an RNN in reverse direction and use it for rescoring.

Backward RNN Reranker (3/3)



Outline

- Recurrent Generation Model
- Convolutional Semantic Reranker
- Backward RNN Reranker
- **Experiments**
 - **Setup**
 - Automatic Evaluation
 - Human Evaluation

Setup

- Data collection:
 - SFX Restaurant domain: 8 system act types, 12 slots (1 is binary).
 - Workers recruited from Amazon MT
 - Asked to generate system responses given a dialogue act.
 - Result in ~5.1K utterances, 228 distinct acts
- Training: BPTT, L2 regularisation, SGD w/ early stopping.
train/valid/test: 3/1/1, data up-sampling

Outline

- Recurrent Generation Model
- Convolutional Semantic Reranker
- Backward RNN Reranker
- Experiments
 - Setup
 - **Automatic Evaluation**
 - Human Evaluation

Automatic Evaluation (1/2)

- Test set: 1039 utterances, 1848 required slots.
- Metrics: BLEU-4 (against multiple references), ERR(slot errors)
- Results averaged over 10 random initialised networks
- Compared with class-based LM (classlm), handcrafted generator (hdc), and kNN based model.

Automatic Evaluation (2/2)

	BLEU	hdc	knn	classlm	rnn
Selection Beam	1/20	0.440	0.591	0.757	0.777
	5/20	-	-	0.678	0.712

	ERR	hdc	knn	classlm	rnn
Selection Beam	1/20	0	17.2	47.8	0
	5/20	-	-	104.6	3.1

Outline

- Recurrent Generation Model
- Convolutional Semantic Reranker
- Backward RNN Reranker
- Experiments
 - Setup
 - Automatic Evaluation
 - **Human Evaluation**

Human Evaluation (1/3)

- Setup
 - Judges (~60) recruited from Amazon MT.
 - Asked to evaluate two system responses pairwise.
 - Comparing handcrafted (hdc), RNN top-1 (rnn₁), RNN sample from top-5 (rnn₅), and class-based LM sampled from top-5 (classlm₅) .
- Metrics:
 - Informativeness, Naturalness (rating out of 5)
 - Preference

Human Evaluation (2/3)

Metrics	hdc	rnn1	hdc	rnn5
Info.	3.75	3.81	3.85	3.93*
Nat.	3.58	3.74**	3.57	3.94**
Pref.	44.8%	55.2%*	37.2%	62.8%**
Metrics	rnn1	rnn5	classlm5	rnn5
Info.	3.75	3.72	4.02	4.15%*
Nat.	3.67	3.58	3.91	4.02
Pref.	47.5%	52.5%	47.1%	52.9%

*=p<.05, **<.005

Human Evaluation (2/3)

Metrics	hdc	rnn1	hdc	rnn5
Info.	3.75	3.81	3.85	3.93*
Nat.	3.58	3.74**	3.57	3.94**
Pref.	44.8%	55.2%*	37.2%	62.8%**
Metrics	rnn1	rnn5	classlm5	rnn5
Info.	3.75	3.72	4.02	4.15%*
Nat.	3.67	3.58	3.91	4.02
Pref.	47.5%	52.5%	47.1%	52.9%

*=p<.05, **<.005

Human Evaluation (2/3)

Metrics	hdc	rnn1	hdc	rnn5
Info.	3.75	3.81	3.85	3.93*
Nat.	3.58	3.74**	3.57	3.94**
Pref.	44.8%	55.2%*	37.2%	62.8%**
Metrics	rnn1	rnn5	classlm5	rnn5
Info.	3.75	3.72	4.02	4.15%*
Nat.	3.67	3.58	3.91	4.02
Pref.	47.5%	52.5%	47.1%	52.9%

*=p<.05, **<.005

Human Evaluation (2/3)

Metrics	hdc	rnn1	hdc	rnn5
Info.	3.75	3.81	3.85	3.93*
Nat.	3.58	3.74**	3.57	3.94**
Pref.	44.8%	55.2%*	37.2%	62.8%**
Metrics	rnn1	rnn5	classlm5	rnn5
Info.	3.75	3.72	4.02	4.15%*
Nat.	3.67	3.58	3.91	4.02
Pref.	47.5%	52.5%	47.1%	52.9%

*=p<.05, **<.005

Human Evaluation (2/3)

Metrics	hdc	rnn1	hdc	rnn5
Info.	3.75	3.81	3.85	3.93*
Nat.	3.58	3.74**	3.57	3.94**
Pref.	44.8%	55.2%*	37.2%	62.8%**
Metrics	rnn1	rnn5	classlm5	rnn5
Info.	3.75	3.72	4.02	4.15%*
Nat.	3.67	3.58	3.91	4.02
Pref.	47.5%	52.5%	47.1%	52.9%

*=p<.05, **<.005

Human Evaluation (3/3)

#	Dialogue act and example realisations of our system, by sampling from top-5 candidates
1	inform(name="stroganoff restaurant",pricerange=cheap,near="fishermans wharf") <i>stroganoff restaurant is a cheap restaurant near fishermans wharf .</i> <i>stroganoff restaurant is in the cheap price range near fishermans wharf .</i>
2	reject(kidsallowed=yes,food="basque") <i>unfortunately there are 0 restaurants that allow kids and serve basque .</i>
3	informonly(name="bund shanghai restaurant", food="shanghainese") <i>i apologize , no other restaurant except bund shanghai restaurant that serves shanghainese .</i> <i>sorry but there is no place other than the restaurant bund shanghai restaurant for shanghainese .</i>
4	confirm(goodformeal=dontcare) <i>i am sorry . just to confirm . you are looking for a restaurant good for any meal ?</i> <i>can i confirm that you do not care about what meal they offer ?</i>
5	request(near) <i>would you like to dine near a particular location ?</i>
6	reqmore() <i>is there anything else i can do for you ?</i>
7	select(kidsallowed=yes, kidsallowed=no) <i>are you looking for a restaurant that allows kids , or does not allow kids ?</i>
8	goodbye() <i>thank you for calling . good bye .</i>

A Brief Summary

- RGM learns generation decisions from corpus.
- No rules, grammars, semantic alignments, or heavy feature engineering are required.
- Can we do better?
 - No heuristic rules for gates.
 - Direct control of generating arbitrary slot-value pairs.
 - Better performance.

Part 2

Semantically Controlled LSTM Generator

Outline

- **SC-LSTM**
- Deep Model
- Experiments
 - Automatic Evaluation
 - Human Evaluation

SC-LSTM (1/4)

- Original LSTM cell

$$\mathbf{i}_t = \sigma(\mathbf{W}_{wi}\mathbf{w}_t + \mathbf{W}_{hi}\mathbf{h}_{t-1}) \quad (1)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_{wf}\mathbf{w}_t + \mathbf{W}_{hf}\mathbf{h}_{t-1}) \quad (2)$$

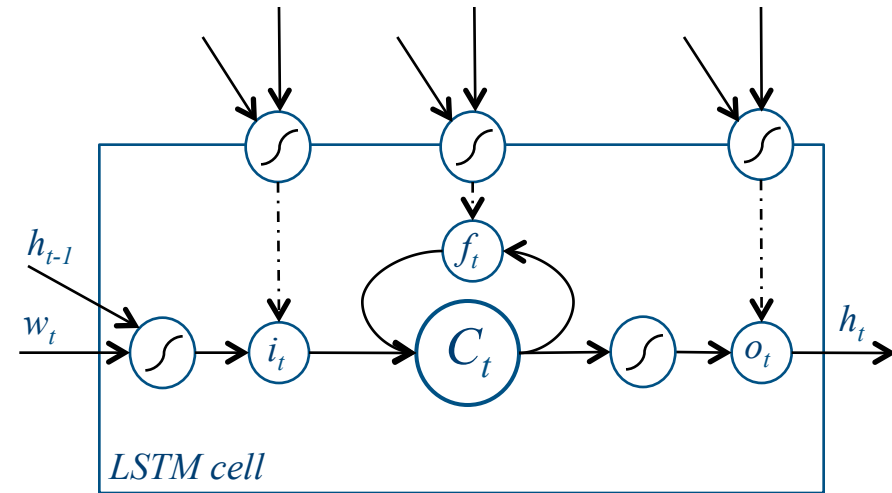
$$\mathbf{o}_t = \sigma(\mathbf{W}_{wo}\mathbf{w}_t + \mathbf{W}_{ho}\mathbf{h}_{t-1}) \quad (3)$$

$$\hat{\mathbf{c}}_t = \tanh(\mathbf{W}_{wc}\mathbf{w}_t + \mathbf{W}_{hc}\mathbf{h}_{t-1}) \quad (4)$$

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \hat{\mathbf{c}}_t \quad (5)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t) \quad (6)$$

(Hochreiter and Schmidhuber, 1997)



SC-LSTM (2/4)

- Original LSTM cell

$$\mathbf{i}_t = \sigma(\mathbf{W}_{wi}\mathbf{w}_t + \mathbf{W}_{hi}\mathbf{h}_{t-1}) \quad (1)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_{wf}\mathbf{w}_t + \mathbf{W}_{hf}\mathbf{h}_{t-1}) \quad (2)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_{wo}\mathbf{w}_t + \mathbf{W}_{ho}\mathbf{h}_{t-1}) \quad (3)$$

$$\hat{\mathbf{c}}_t = \tanh(\mathbf{W}_{wc}\mathbf{w}_t + \mathbf{W}_{hc}\mathbf{h}_{t-1}) \quad (4)$$

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \hat{\mathbf{c}}_t \quad (5)$$

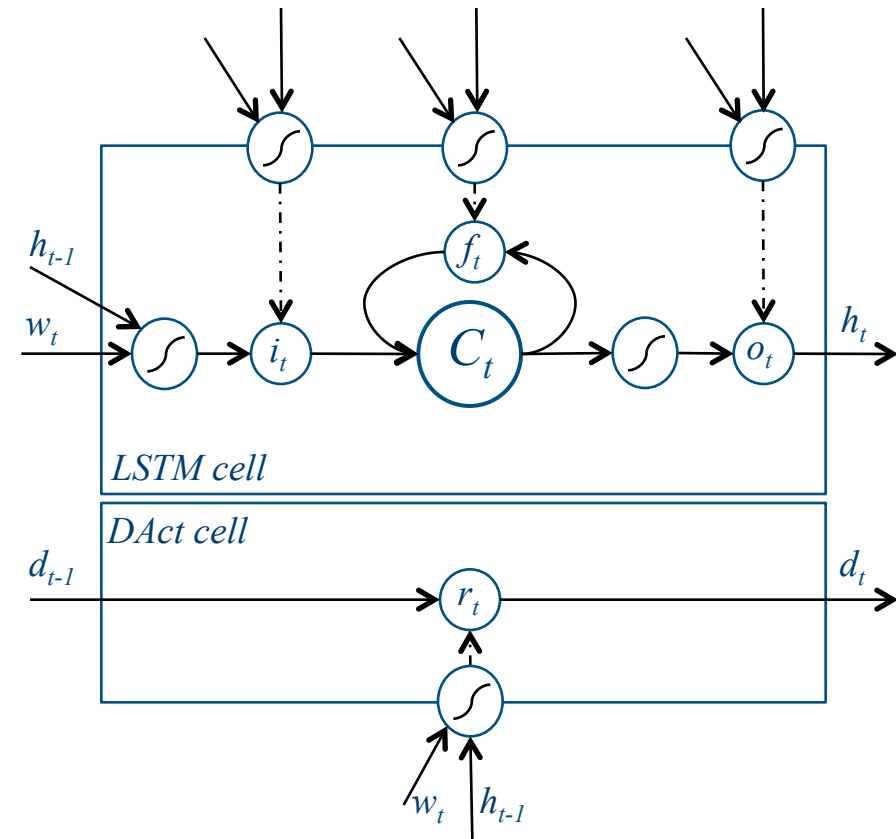
$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t) \quad (6)$$

- DA cell

$$\mathbf{r}_t = \sigma(\mathbf{W}_{wr}\mathbf{w}_t + \alpha\mathbf{W}_{hr}\mathbf{h}_{t-1}) \quad (7)$$

$$\mathbf{d}_t = \mathbf{r}_t \odot \mathbf{d}_{t-1} \quad (8)$$

(Hochreiter and Schmidhuber, 1997)



SC-LSTM (3/4)

- Original LSTM cell

$$\mathbf{i}_t = \sigma(\mathbf{W}_{wi}\mathbf{w}_t + \mathbf{W}_{hi}\mathbf{h}_{t-1}) \quad (1)$$

$$\mathbf{f}_t = \sigma(\mathbf{W}_{wf}\mathbf{w}_t + \mathbf{W}_{hf}\mathbf{h}_{t-1}) \quad (2)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_{wo}\mathbf{w}_t + \mathbf{W}_{ho}\mathbf{h}_{t-1}) \quad (3)$$

$$\hat{\mathbf{c}}_t = \tanh(\mathbf{W}_{wc}\mathbf{w}_t + \mathbf{W}_{hc}\mathbf{h}_{t-1}) \quad (4)$$

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \hat{\mathbf{c}}_t \quad (5)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t) \quad (6)$$

- DA cell

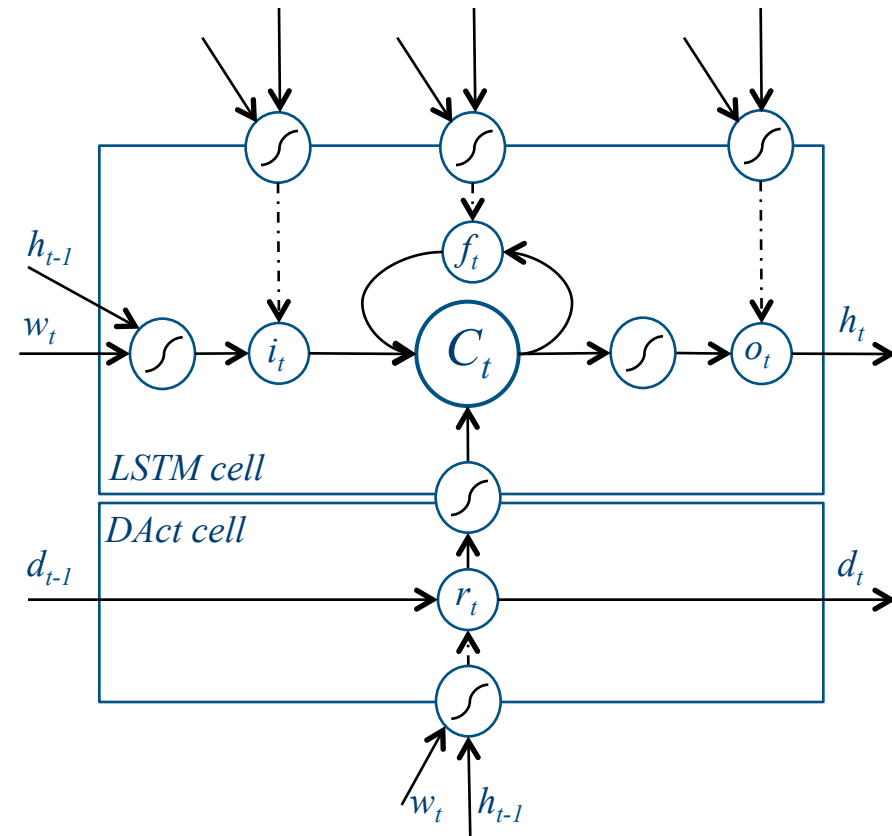
$$\mathbf{r}_t = \sigma(\mathbf{W}_{wr}\mathbf{w}_t + \alpha\mathbf{W}_{hr}\mathbf{h}_{t-1}) \quad (7)$$

$$\mathbf{d}_t = \mathbf{r}_t \odot \mathbf{d}_{t-1} \quad (8)$$

- Modify eq. (6) to

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \hat{\mathbf{c}}_t + \tanh(\mathbf{W}_{dc}\mathbf{d}_t) \quad (9)$$

(Hochreiter and Schmidhuber, 1997)



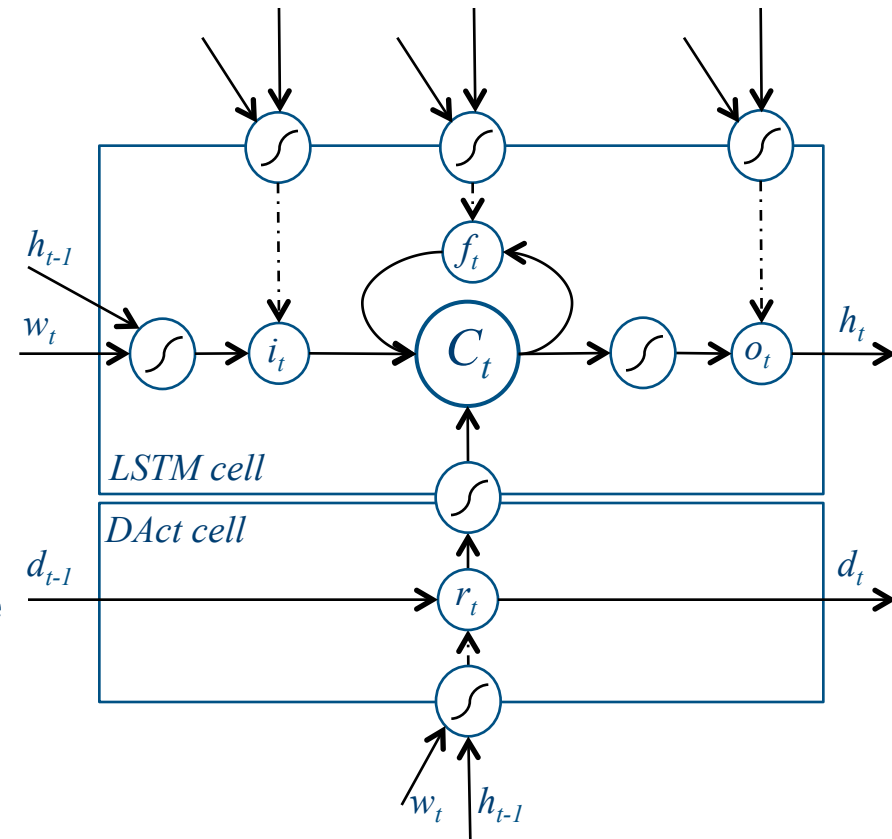
SC-LSTM (4/4)

- Cost function

$$F(\theta) = \sum_t \mathbf{p}_t^\top \log(\mathbf{y}_t) + \|\mathbf{d}_T\| + \sum_{t=0}^{T-1} \eta \xi \|\mathbf{d}_{t+1} - \mathbf{d}_t\|$$

- 1st term : cross entropy error
- 2nd term: make sure rendering all the information needed
- 3rd term: prevent undesirable gating behaviors

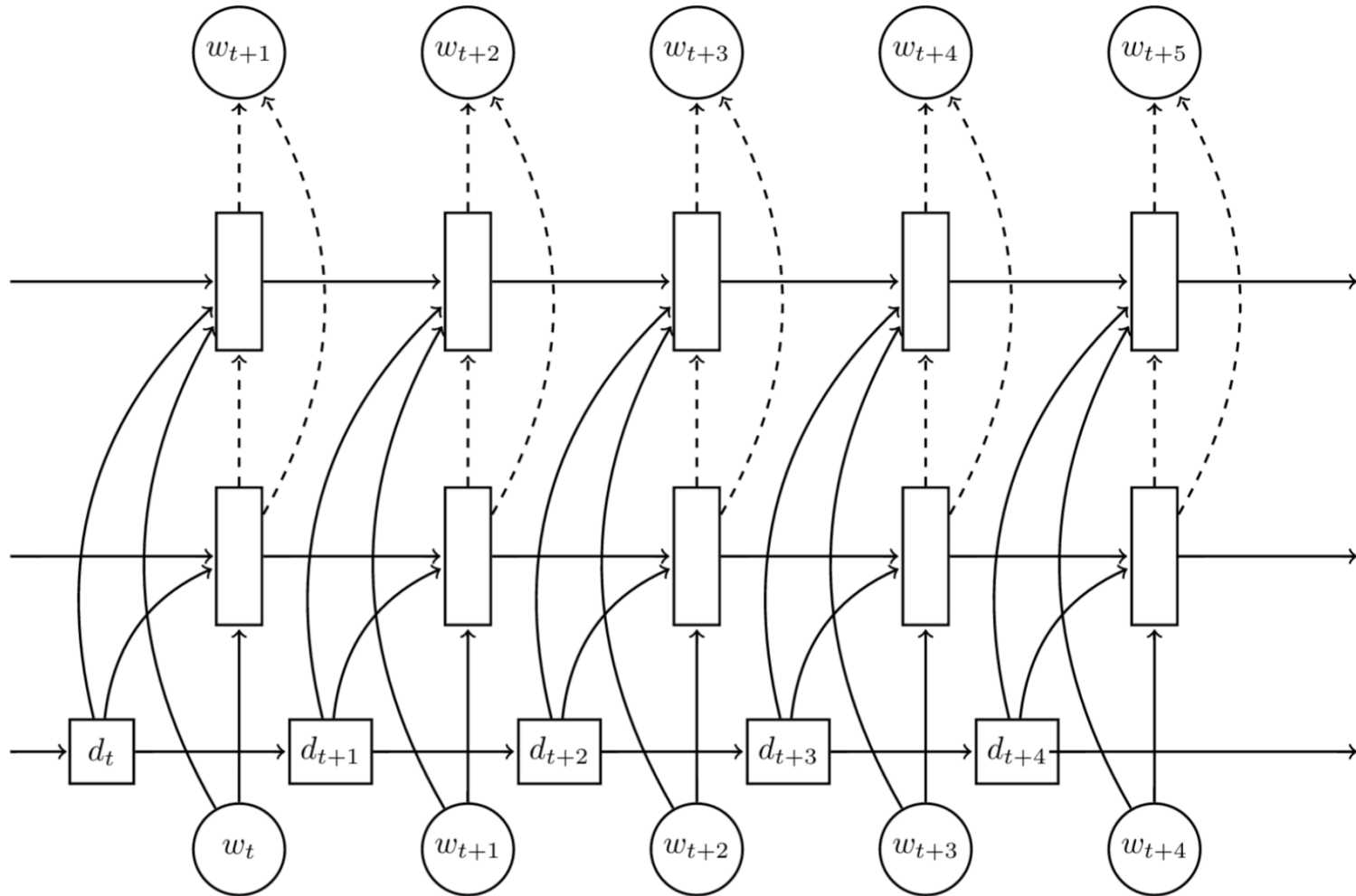
(Hochreiter and Schmidhuber, 1997)



Outline

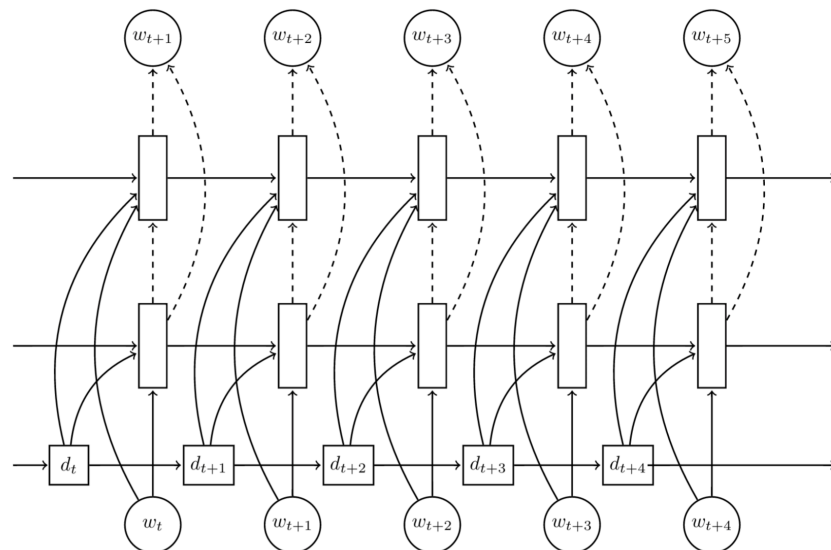
- SC-LSTM
- **Deep Model**
- Experiments
 - Automatic Evaluation
 - Human Evaluation

Deep Model (1/2)



Deep Model (2/2)

- Techniques applied
 - Skip connection (Graves et al 2013)
 - RNN dropout (Srivastava et al 2014)



- Gating Equation is modified

$$\mathbf{r}_t = \sigma(\mathbf{W}_{wr}\mathbf{w}_t + \alpha\mathbf{W}_{hr}\mathbf{h}_{t-1}) \quad (7)$$

- To

$$\mathbf{r}_t = \sigma(\mathbf{W}_{wr}\mathbf{w}_t + \sum_l \alpha_l \mathbf{W}_{hr}^l \mathbf{h}_{t-1}^l) \quad (12)$$

Outline

- SC-LSTM
- Deep Model
- **Experiments**
 - **Automatic Evaluation**
 - Human Evaluation

Automatic Evaluation (1/3)

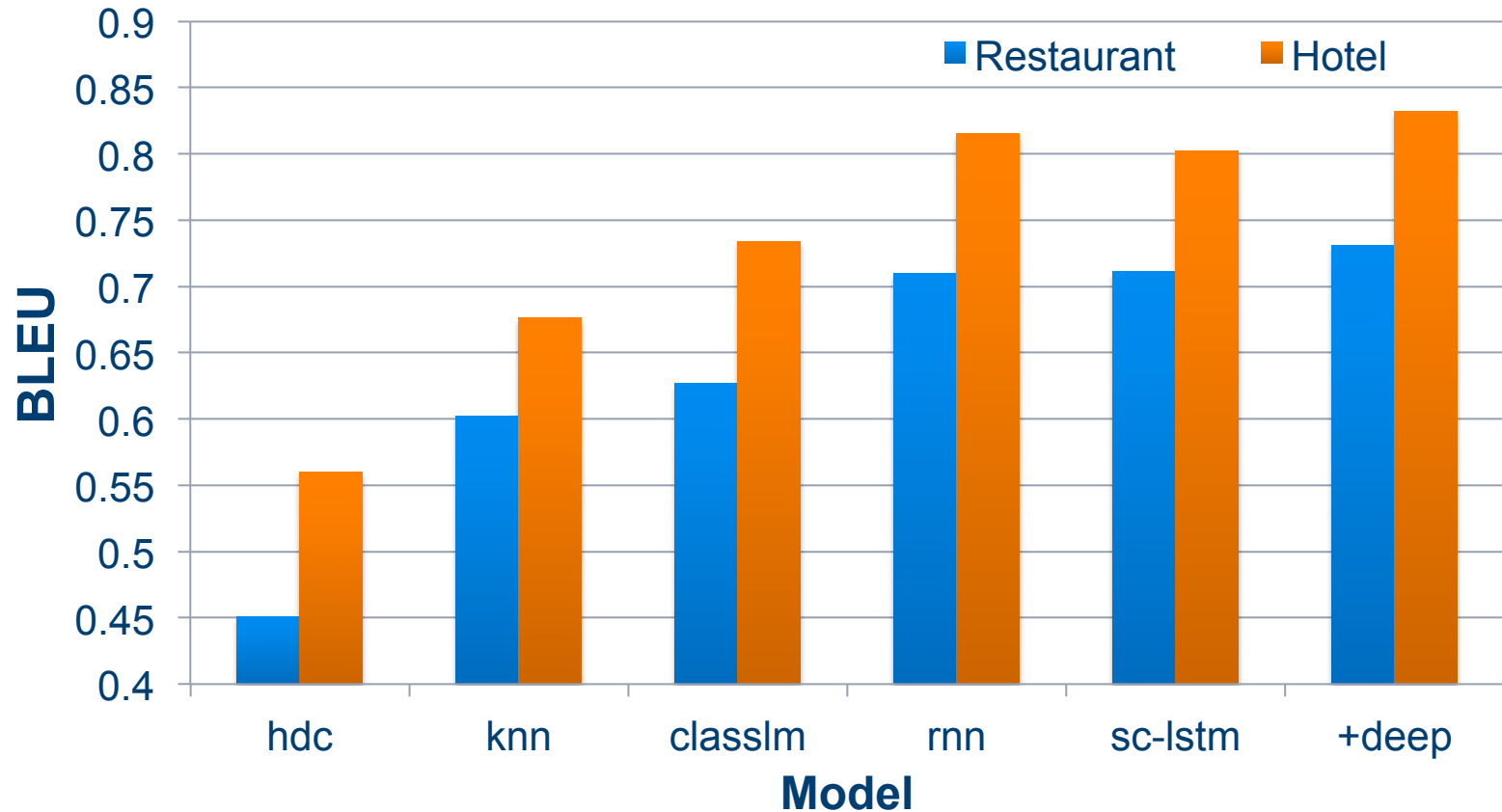
- Dataset: SFX Restaurant & SFX Hotel Domains
 - 5K utterances, 3:1:1 splitting
 - 248/164 distinct acts, 2.25/1.95 # of slot per DA

- Ontologies:

	SF Restaurant	SF Hotel
act type	inform, inform_only, reject, confirm, select, request, reqmore, goodbye	
shared	name, type, *pricerange, price, phone, address, postcode, *area, *near	
specific	*food *goodformeal *kids-allowed	*hasinternet *acceptscards *dogs-allowed

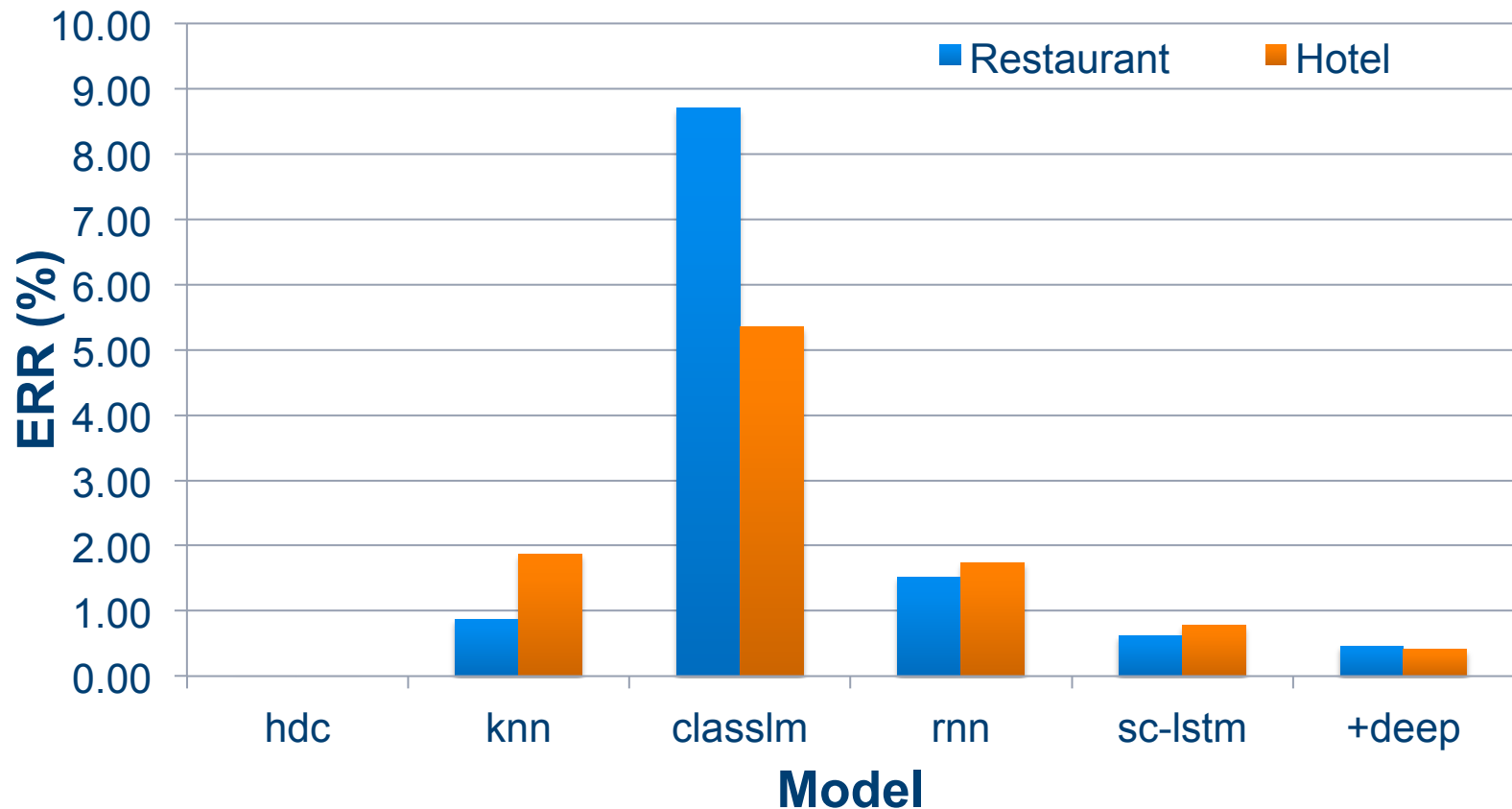
bold=binary slots, *=slots can take “don’t care” value

Automatic Evaluation (2/3)



Selection scheme : 5/20

Automatic Evaluation (3/3)



Selection scheme : 5/20

Outline

- SC-LSTM
- Deep Model
- Experiments
 - Automatic Evaluation
 - **Human Evaluation**

Human Evaluation (1/3)

- Setting
 - Done on SFX Restaurant domain
 - Comparing *classlm*, *rnn w/*, *sc-lstm* and *+deep*
- Metrics
 - Informativeness, Naturalness, Preference

Human Evaluation (2/3)

Method	Informativeness	Naturalness
+deep	2.58	2.51
sc-lstm	2.59	2.50
rnn w/ classlm	2.53	2.42*
	2.46**	2.45

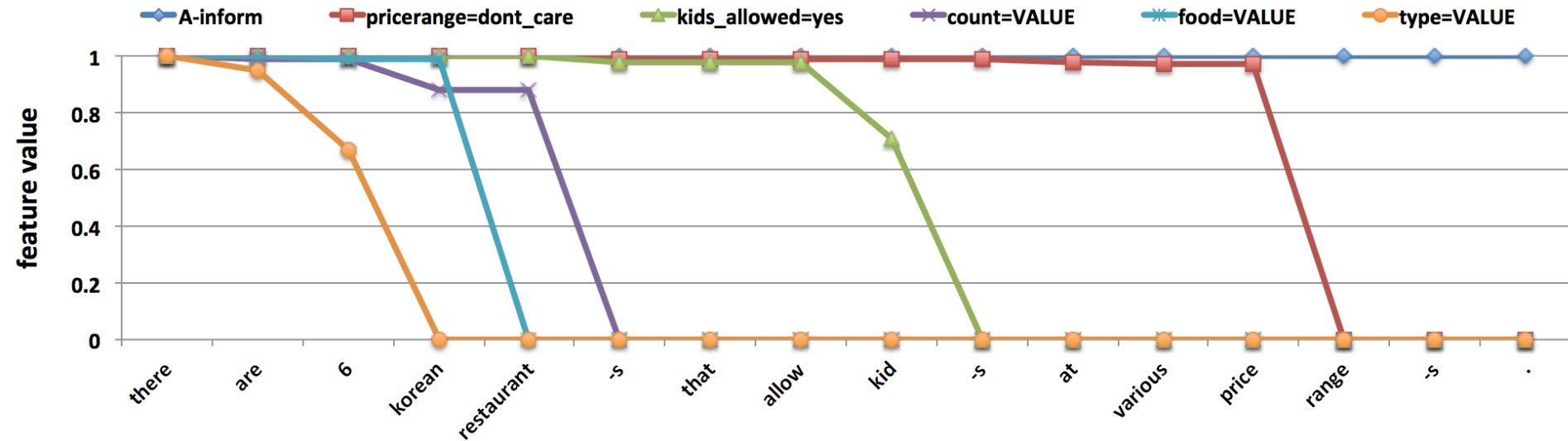
* $p < 0.05$ ** $p < 0.005$

Human Evaluation (3/3)

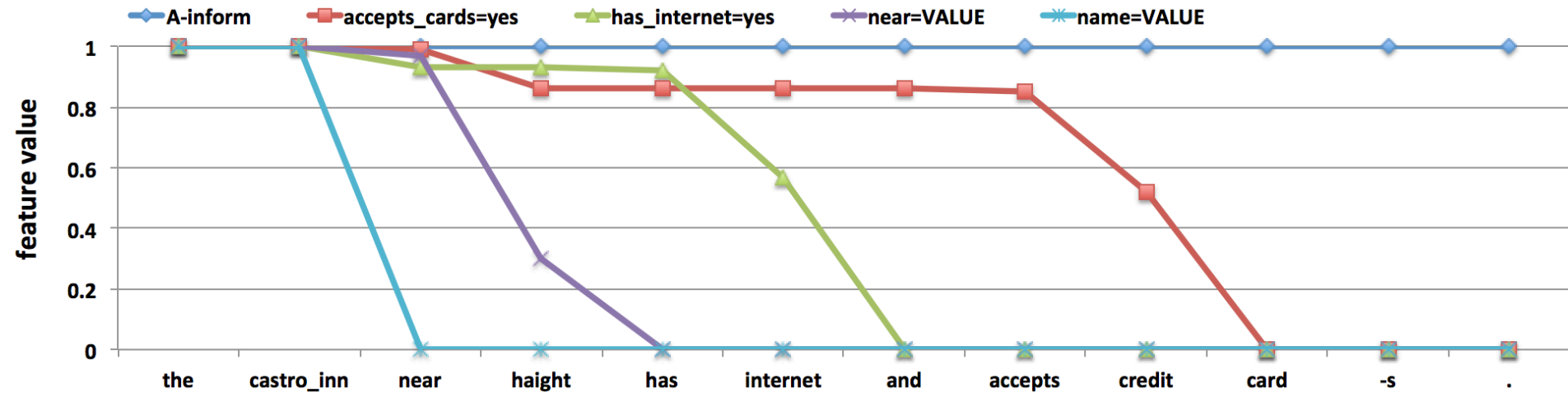
Pref. %	classlm	rnn w/	sc-lstm	+deep
classlm	-	46.0	40.9 ^{**}	37.7 ^{**}
rnn w/	54.0	-	43.0	35.7 [*]
sc-lstm	59.1 [*]	57	-	47.6
+deep	62.3 ^{**}	64.3 ^{**}	52.4	-

^{*} $p < 0.05$ ^{**} $p < 0.005$

Example



Example



Conclusion

Conclusion – Why RNN for NLG?

- ✓ Elegant structure for modeling sequences.
 - Flexible architecture for adding auxiliary information.
- ✓ Collecting data is convenient and quick (crowdsourcing).
- ✓ More human-like and colloquial.
- ✓ No expert knowledge is required.
 - Extensible, adaptation techniques exist.
 - Distributed representation
- ✓ Less cost, quicker development cycle
- ✓ End-to-End trainable

Papers

- Tsung-Hsien Wen, Milica Gasic , Dongho Kim, Nikola Mrksic, Pei-Hao Su, David Vandyke, and Steve Young. 2015. Stochastic language generation in dialogue using recurrent neural networks with convolutional sentence reranking. In *Proceedings of SIGdial*. Association for Computational Linguistics.
- Tsung-Hsien Wen, Milica Gasic , Nikola Mrksic, Pei-Hao Su, David Vandyke, and Steve Young. 2015. Semantically Conditioned LSTM-based Natural Language Generation for Spoken Dialogue Systems. To be appear In *Proceedings of EMNLP*. Association for Computational Linguistics.

Reference

- Amanda Stent, Matthew Marge, and Mohit Singhai. 2005. Evaluating evaluation methods for generation in the presence of variation. In Proceedings of CICLing 2005.
- Alice H. Oh and Alexander I. Rudnicky. 2000. Stochastic language generation for spoken dialogue systems. In Proceedings of the 2000 ANLP/NAACL Workshop on Conversational Systems.
- Tomas Mikolov, Martin Karafit, Lukas Burget, Jan Cernocky, and Sanjeev Khudanpur. 2010. Recurrent neural network based language model. *In Proceedings on InterSpeech*.
- Nal Kalchbrenner, Edward Grefenstette, and Phil Blunsom. 2014. A convolutional neural network for modelling sentences. Proceedings of the 52nd Annual Meeting of ACL.

Reference

- Sepp Hochreiter and Jurgen Schmidhuber. 1997. Long short-term memory. *Neural Computation*.
- Alex Graves. 2013. Generating sequences with recurrent neural networks. *CoRR*, abs/1308.0850.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*.

Thank you! Questions?

*This project is supported by Toshiba Research Europe Ltd,
Cambridge Research Laboratory.*

Dialogue Systems Group