



UNSUPERVISED LEARNING AND MODELING OF KNOWLEDGE AND INTENT FOR SPOKEN DIALOGUE SYSTEMS

Yun-Nung (Vivian) Chen | <http://vivianchen.idv.tw>





OUTLINE



Introduction



Semantic Decoding



- Ontology Induction



- Knowledge Graph Propagation



- Matrix Factorization



- Experiments



Future Work



Conclusions

OUTLINE



Introduction



Semantic Decoding [ACL-IJCNLP'15]



- Ontology Induction



- Knowledge Graph Propagation



- Matrix Factorization



- Experiments



Future Work



Conclusions

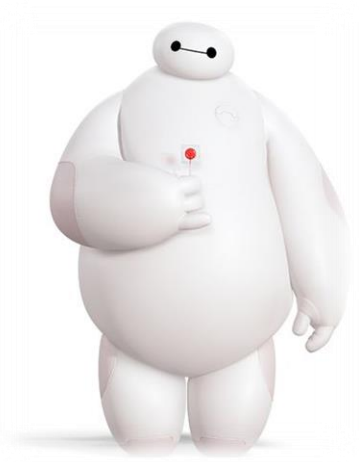
A POPULAR ROBOT - BAYMAX



A POPULAR ROBOT - BAYMAX

Baymax is capable of maintaining a good **spoken dialogue system** and **learning** new knowledge for better **understanding** and **interacting** with people.

The goal is to automate learning and understanding procedures in system development.



SPOKEN DIALOGUE SYSTEM (SDS)

Spoken dialogue systems are the intelligent agents that are able to help users finish tasks more efficiently via speech interactions.

Spoken dialogue systems are being incorporated into various devices (smart-phones, smart TVs, in-car navigating system, etc).



Apple's Siri
Microsoft's Cortana



Microsoft's XBOX Kinect



Amazon's Echo



Samsung's SMART TV



Google Now

<https://www.apple.com/ios/siri/>
<http://www.windowsphone.com/en-us/how-to/wp8/cortana/meet-cortana>
<http://www.xbox.com/en-US/>
<http://www.amazon.com/oc/echo/>
<http://www.samsung.com/us/experience/smart-tv/>
<https://www.google.com/landing/now/>

LARGE SMART DEVICE POPULATION

The number of global smartphone users will surpass **2 billion** in 2016.

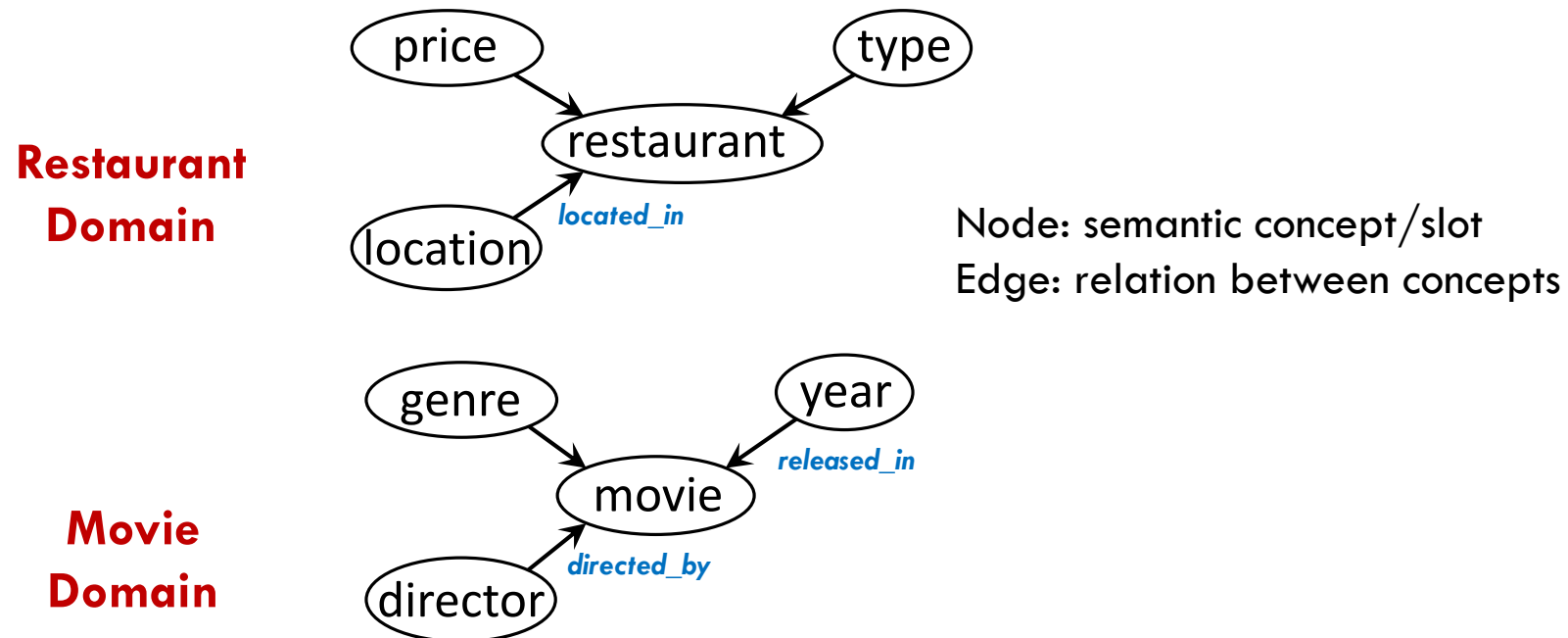
As of 2012, there are **1.1 billion** automobiles on the earth.



The more **natural** and **convenient** input of the devices evolves towards **speech**

KNOWLEDGE REPRESENTATION/ONTOLOGY

Traditional SDSs require **manual annotations** for **specific domains** to represent domain knowledge.



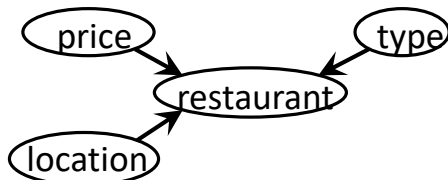
UTTERANCE SEMANTIC REPRESENTATION

A spoken language understanding (SLU) component requires the domain ontology to decode utterances into semantic forms, which contain **core content (a set of slots and slot-fillers)** of the utterance.

Restaurant Domain



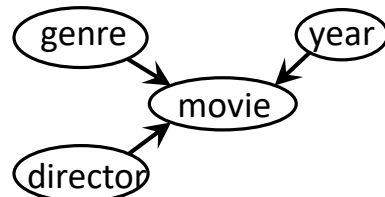
find a cheap taiwanese restaurant in seattle



target="restaurant", price="cheap",
type="taiwanese", location="seattle"

Movie Domain

show me action movies directed by james cameron



target="movie", genre="action",
director="james cameron"

CHALLENGES FOR SDS

An SDS in a new domain requires

- 1) A hand-crafted domain ontology
- 2) Utterances labelled with semantic representations
- 3) An SLU component for mapping utterances into semantic representations

With increasing spoken interactions, building domain ontologies and annotating utterances cost a lot so that the data does not scale up.

The goal is to **enable an SDS to automatically learn this knowledge** so that open domain requests can be handled.

INTERACTION EXAMPLE

User



find an inexpensive eating place for taiwanese food



Intelligent Agent

Inexpensive Taiwanese eating places include Din Tai Fung, Boiling Point, etc. What do you want to choose? I can help you go there.

Q: How does a dialogue system process this request?

SDS PROCESS – AVAILABLE DOMAIN ONTOLOGY

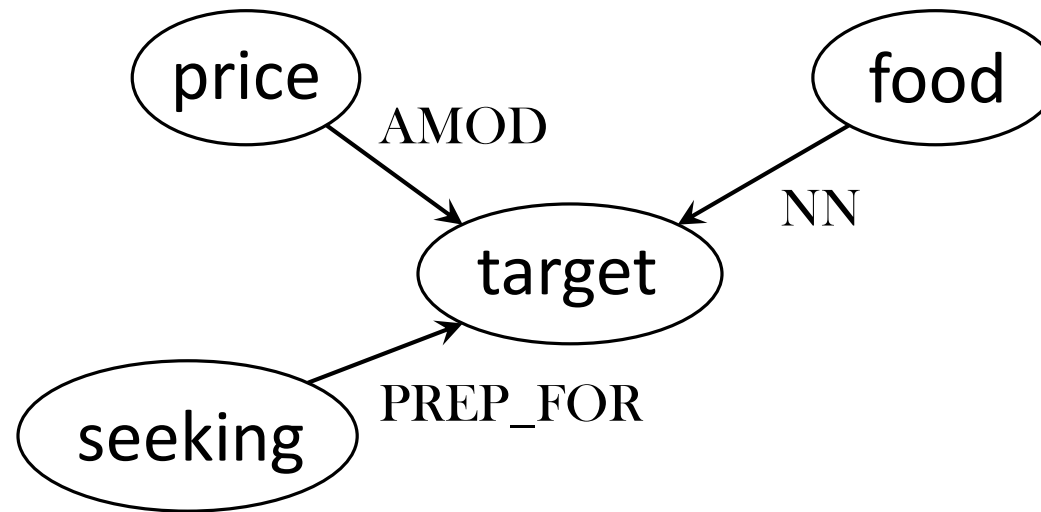
User



find an inexpensive eating place for taiwanese food



Intelligent Agent



Organized Domain Knowledge

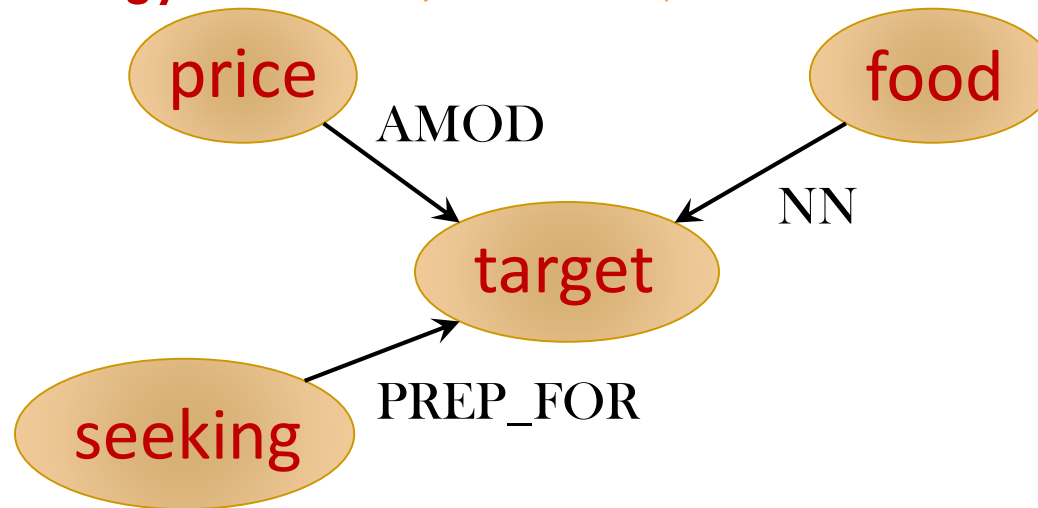
SDS PROCESS – AVAILABLE DOMAIN ONTOLOGY

User



find an inexpensive eating place for taiwanese food

Ontology Induction (*semantic slot*)



Intelligent Agent

Organized Domain Knowledge

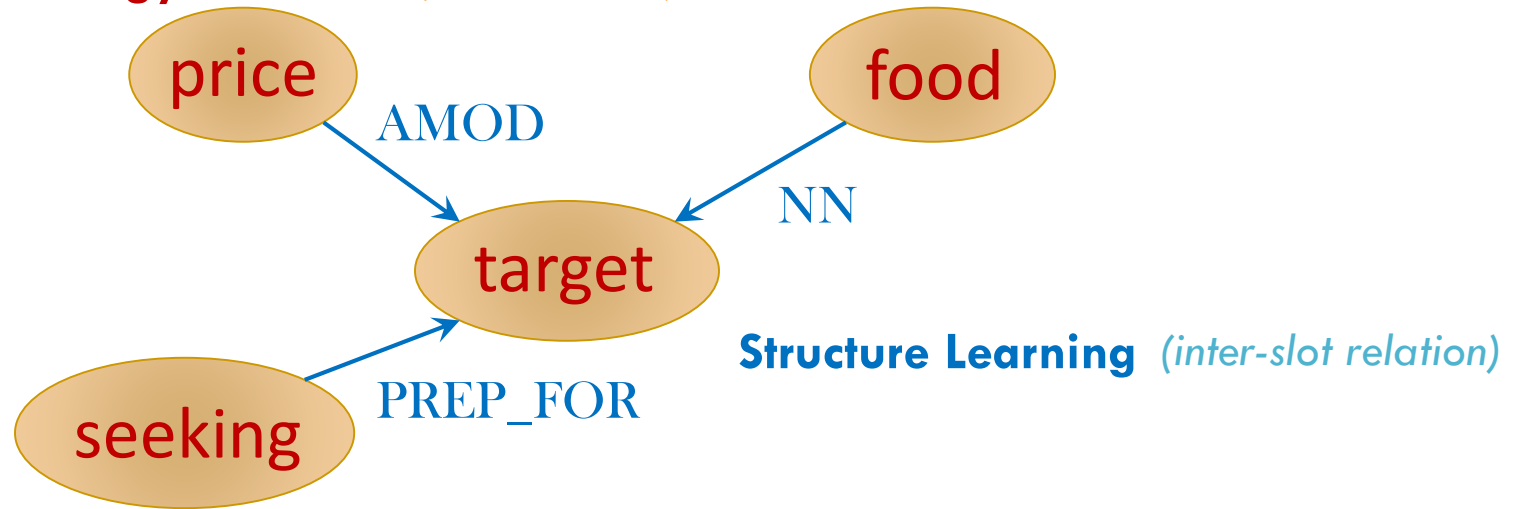
SDS PROCESS – AVAILABLE DOMAIN ONTOLOGY

User



find an inexpensive eating place for taiwanese food

Ontology Induction (*semantic slot*)



Intelligent Agent

Organized Domain Knowledge

SDS PROCESS – SPOKEN LANGUAGE UNDERSTANDING (SLU)

User



find an inexpensive eating place for taiwanese food

price

AMOD

food

NN

target

seeking

PREP_FOR

seeking="find"
target="eating place"
price="inexpensive"
food="taiwanese food"

Intelligent Agent



SDS PROCESS – SPOKEN LANGUAGE UNDERSTANDING (SLU)

User



find an inexpensive eating place for taiwanese food

price

AMOD

food

NN

target

seeking

PREP_FOR

Semantic Decoding

seeking="find"

target="eating place"

price="inexpensive"

food="taiwanese food"

Intelligent Agent



SDS PROCESS – DIALOGUE MANAGEMENT (DM)

User



find an inexpensive eating place for taiwanese food

price

AMOD

food

NN

target

seeking

PREP_FOR

```
SELECT restaurant {  
  restaurant.price="inexpensive"  
  restaurant.food="Taiwanese food"  
}
```

Intelligent Agent



SDS PROCESS – DIALOGUE MANAGEMENT (DM)

User



find an inexpensive eating place for taiwanese food

price

AMOD

food

NN

Surface Form Derivation
(natural language)

target

```
SELECT restaurant {  
  restaurant.price="inexpensive"  
  restaurant.food="Taiwanese food"  
}
```

seeking

PREP_FOR



Intelligent Agent

SDS PROCESS – DIALOGUE MANAGEMENT (DM)

User



find an inexpensive eating place for taiwanese food

```
SELECT restaurant {  
  restaurant.price="inexpensive"  
  restaurant.food="Taiwanese food"  
}
```

Predicted behavior: navigation



Intelligent Agent

Din Tai Fung
Boiling Point
:
:

SDS PROCESS – DIALOGUE MANAGEMENT (DM)

User



find an inexpensive eating place for taiwanese food

```
SELECT restaurant {  
  restaurant.price="inexpensive"  
  restaurant.food="Taiwanese food"  
}
```

Predicted behavior: navigation

Behavior Prediction

Din Tai Fung
Boiling Point
:
:
:



Intelligent Agent

SDS PROCESS – NATURAL LANGUAGE GENERATION (NLG)

User



find an inexpensive eating place for taiwanese food



Intelligent Agent

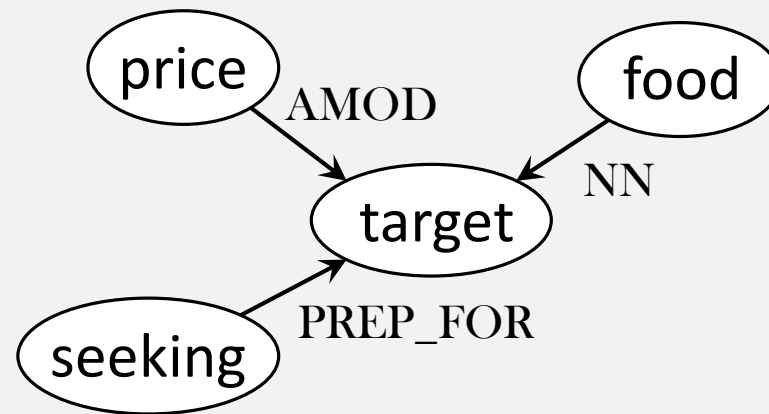
Inexpensive Taiwanese eating places include Din Tai Fung, Boiling Point, etc. What do you want to choose? I can help you go there. **(navigation)**

GOALS

User



find an inexpensive eating place for taiwanese food



```
SELECT restaurant {  
  restaurant.price="inexpensive"  
  restaurant.food="taiwanese food"  
}
```

Predicted behavior: navigation

Required Domain-Specific Information

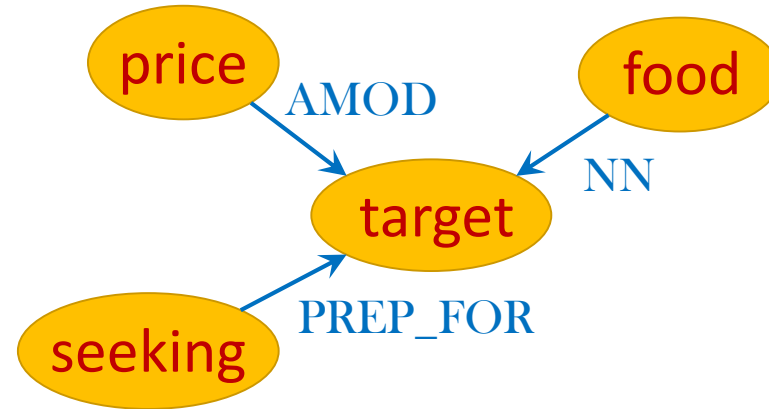
FIVE GOALS

User



find an inexpensive **eating place** for taiwanese food

1. Ontology Induction (semantic slot)



2. Structure Learning (inter-slot relation)

Required Domain-Specific Information

3. Surface Form Derivation

SELECT **restaurant** { (natural language)
restaurant.price="inexpensive"
restaurant.food="taiwanese food"
}

4. Semantic Decoding

Predicted behavior: navigation

5. Behavior Prediction

FIVE GOALS

User



find an inexpensive eating place for taiwanese food

1. Ontology Induction (*semantic slot*)

3. Surface Form Derivation
(*natural language*)

4. Semantic Decoding

2. Structure Learning (*inter-slot relation*)

5. Behavior Prediction

FIVE GOALS

User



find an inexpensive eating place for taiwanese food

- 1. Ontology Induction**
- 2. Structure Learning**
- 3. Surface Form Derivation**

Knowledge Acquisition

- 4. Semantic Decoding**
- 5. Behavior Prediction**

SLU Modeling

OUTLINE



Introduction



Semantic Decoding [ACL-IJCNLP'15]



- Ontology Induction



- Knowledge Graph Propagation



- Matrix Factorization



- Experiments



Future Work

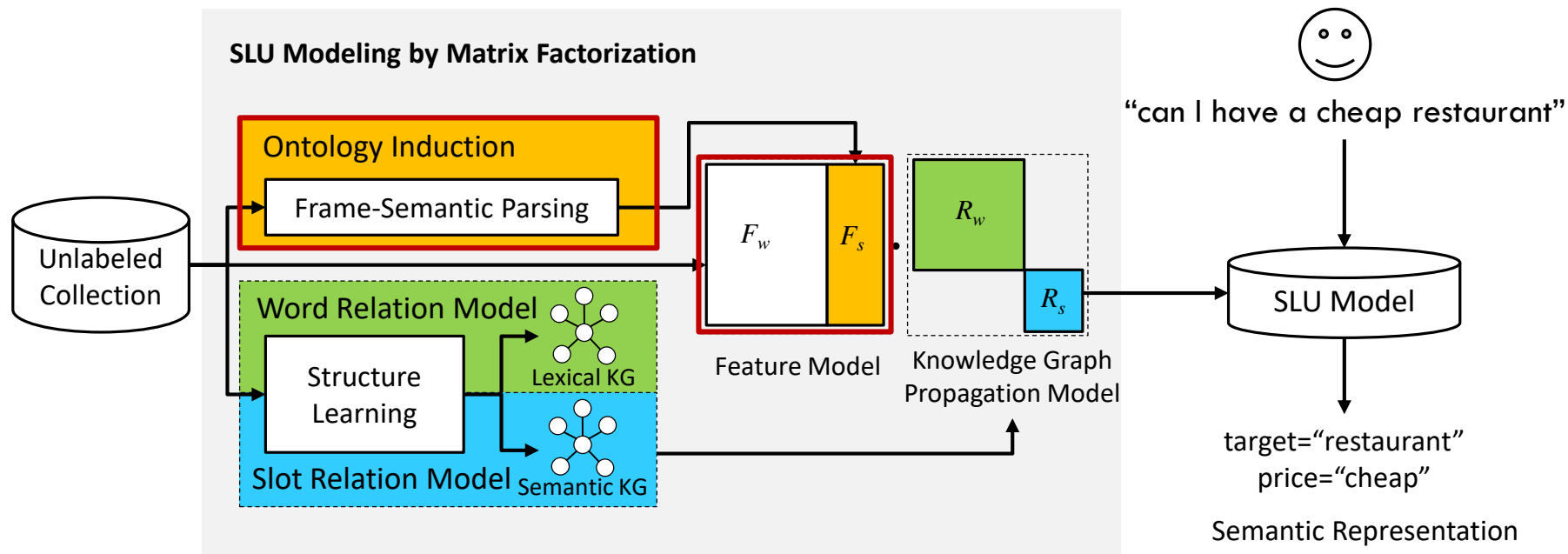


Conclusions

SEMANTIC DECODING

Input: user utterances

Output: the domain-specific semantic concepts included in each individual utterance



OUTLINE



Introduction



Semantic Decoding [ACL-IJCNLP'15]



▪ **Ontology Induction**



▪ Knowledge Graph Propagation



▪ Matrix Factorization



▪ Experiments



Future Work



Conclusions

PROBABILISTIC FRAME-SEMANTIC PARSING

FrameNet [Baker et al., 1998]

- a linguistically semantic resource, based on the frame-semantics theory
- “low fat milk” → “milk” evokes the “food” frame;
“low fat” fills the descriptor frame element

SEMAFOR [Das et al., 2014]

- a state-of-the-art frame-semantics parser, trained on manually annotated FrameNet sentences

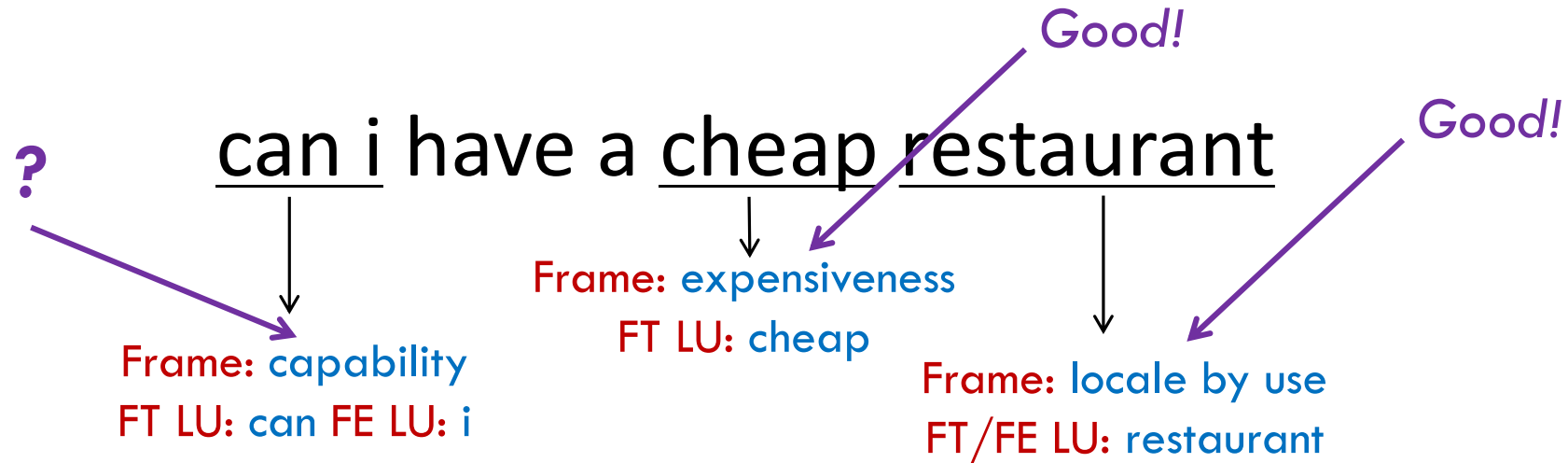


Baker et al., "The Berkeley FrameNet project," in Proc. of International Conference on Computational Linguistics, 1998.

Das et al., "Frame-semantic parsing," in Proc. of Computational Linguistics, 2014.



FRAME-SEMANTIC PARSING FOR UTTERANCES



FT: Frame Target; FE: Frame Element; LU: Lexical Unit

1st Issue: adapting **generic** frames to **domain-specific** settings for SDSs

OUTLINE



Introduction



Semantic Decoding [ACL-IJCNLP'15]



- Ontology Induction



- **Knowledge Graph Propagation (for 1st issue)**



- Matrix Factorization



- Experiments



Future Work

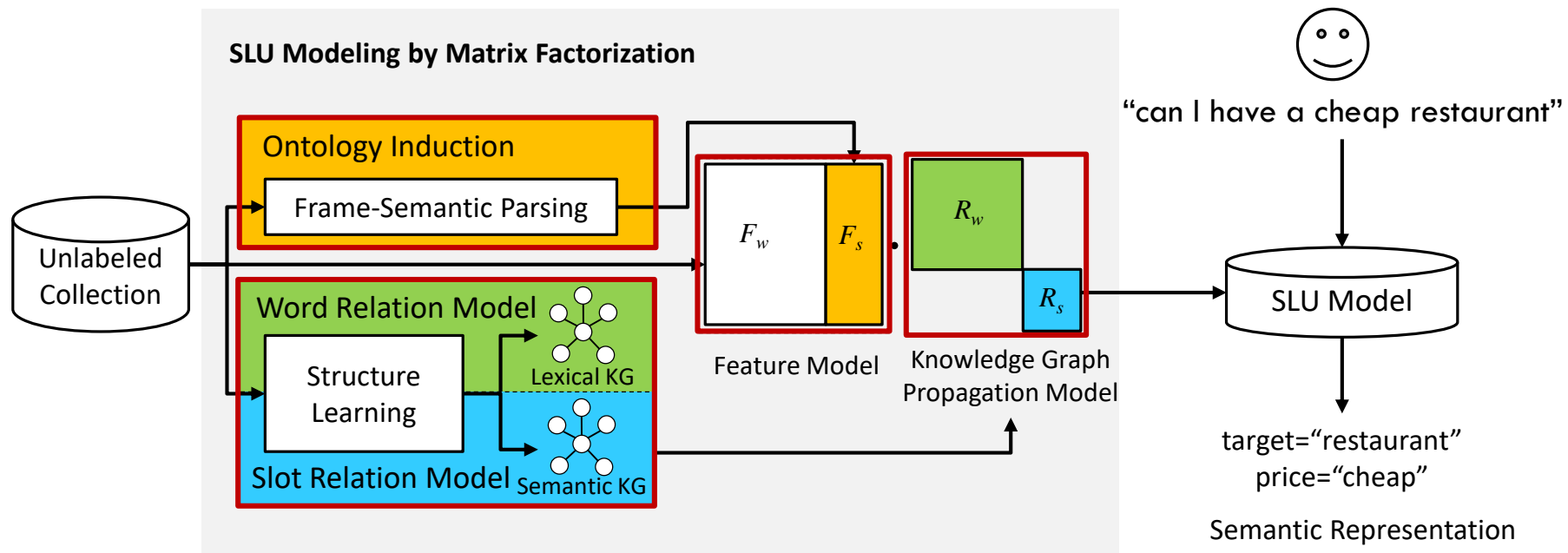


Conclusions

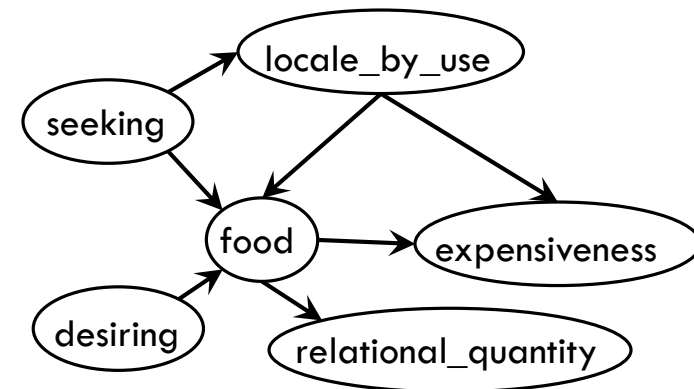
SEMANTIC DECODING

Input: user utterances

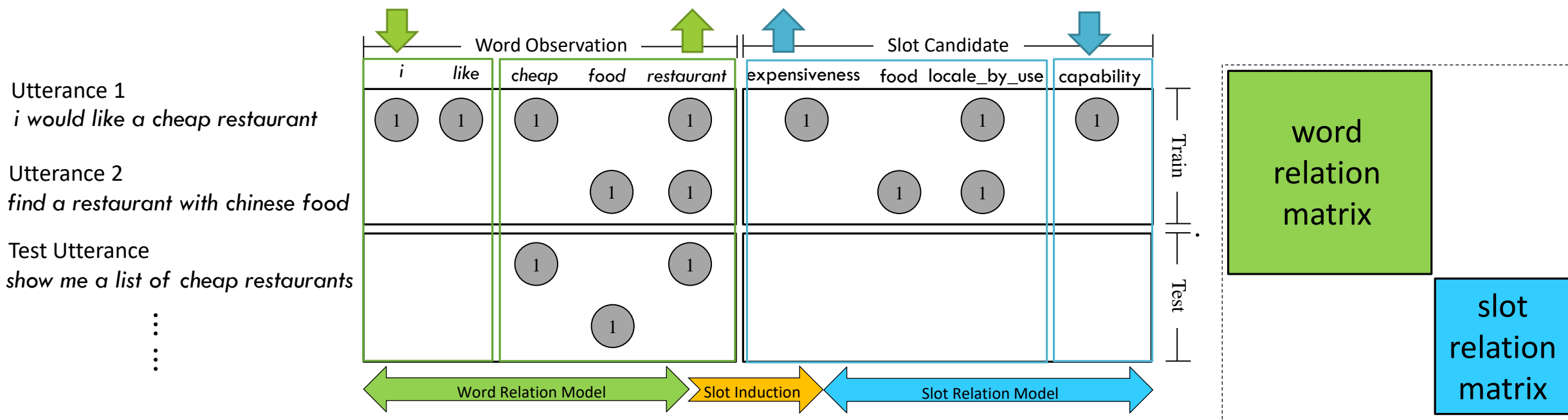
Output: the domain-specific semantic concepts included in each individual utterance



1ST ISSUE: HOW TO ADAPT GENERIC SLOTS TO DOMAIN-SPECIFIC? KNOWLEDGE GRAPH PROPAGATION MODEL



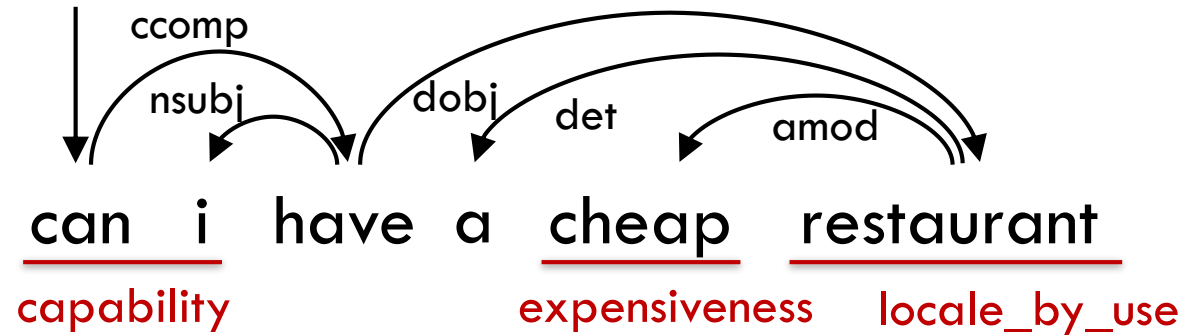
Assumption: The domain-specific words/slots have more dependency to each other.



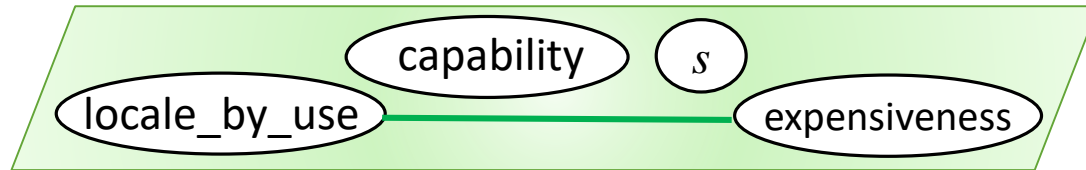
The relation matrices allow each node propagate the scores to its neighbor in the knowledge graph, so that the domain-specific words/slots have higher scores during training.

KNOWLEDGE GRAPH CONSTRUCTION

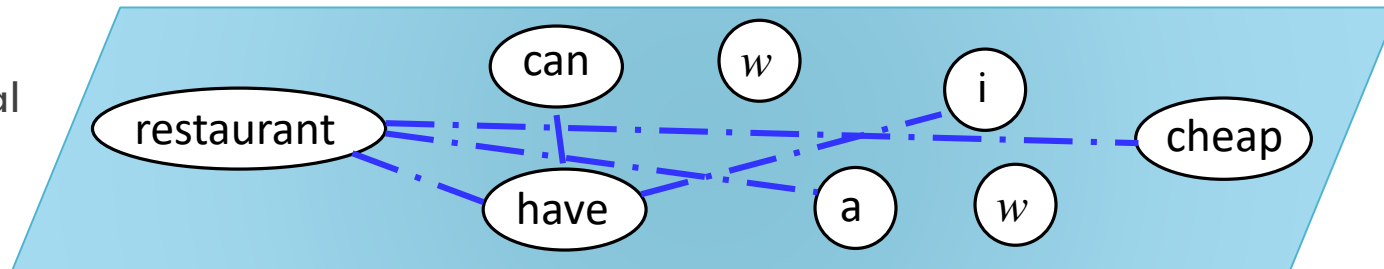
Syntactic dependency parsing on utterances



Slot-based semantic knowledge graph



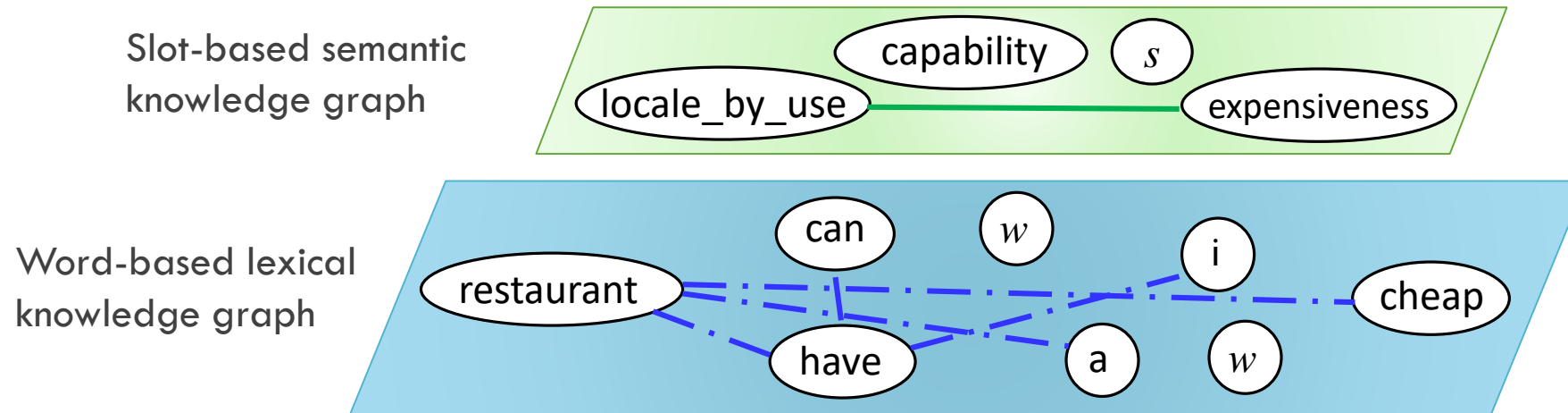
Word-based lexical knowledge graph



KNOWLEDGE GRAPH CONSTRUCTION

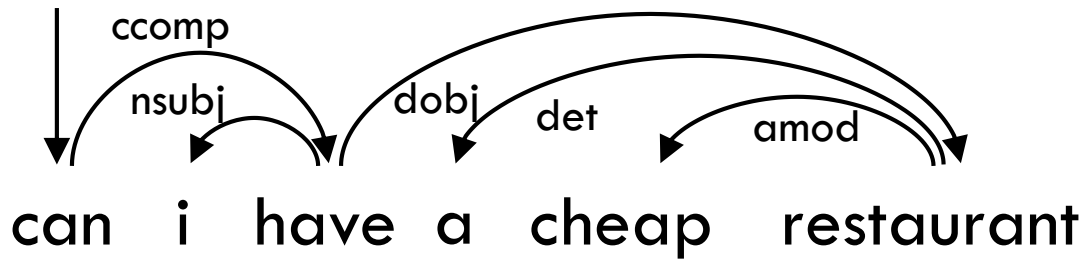
The edge between a node pair is weighted as relation importance for build the matrix

How to decide the weights to represent relation importance?



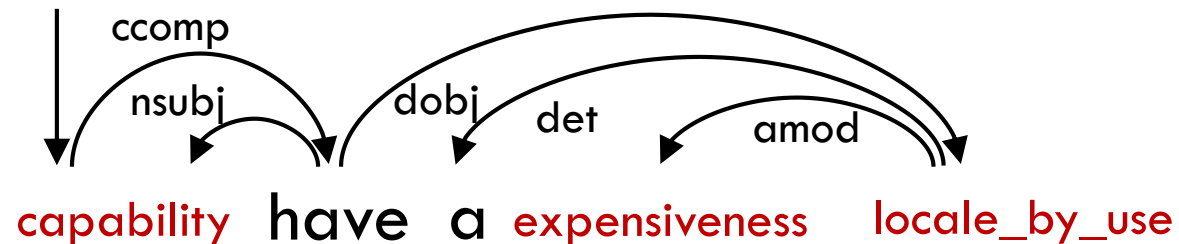
WEIGHT MEASUREMENT BY EMBEDDINGS

Dependency-based word embeddings



can = [0.8 ... 0.24]
have = [0.3 ... 0.21]
:
:

Dependency-based slot embeddings



expensiveness = [0.12 ... 0.7]
capability = [0.3 ... 0.6]
:
:

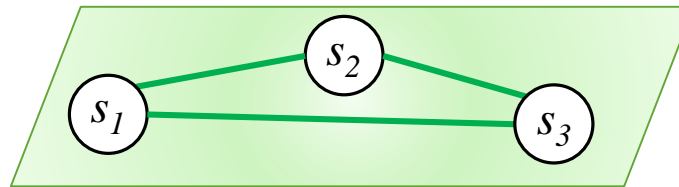


WEIGHT MEASUREMENT BY EMBEDDINGS

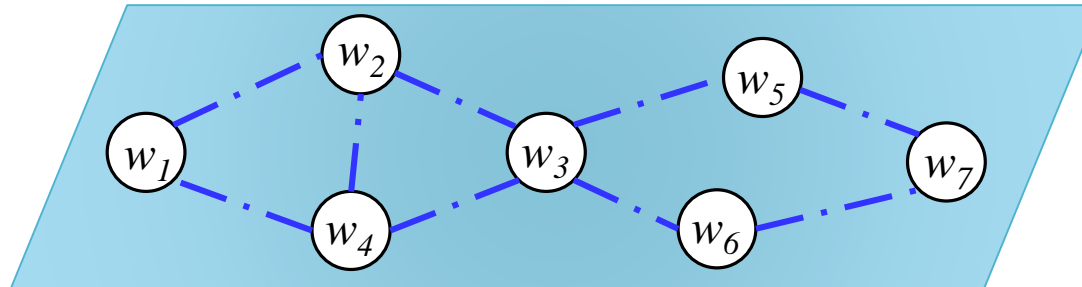
Compute edge weights to represent relation importance

- Slot-to-slot semantic relation R_S^S : similarity between slot embeddings
- Slot-to-slot dependency relation R_S^D : dependency score between slot embeddings
- Word-to-word semantic relation R_W^S : similarity between word embeddings
- Word-to-word dependency relation R_W^D : dependency score between word embeddings

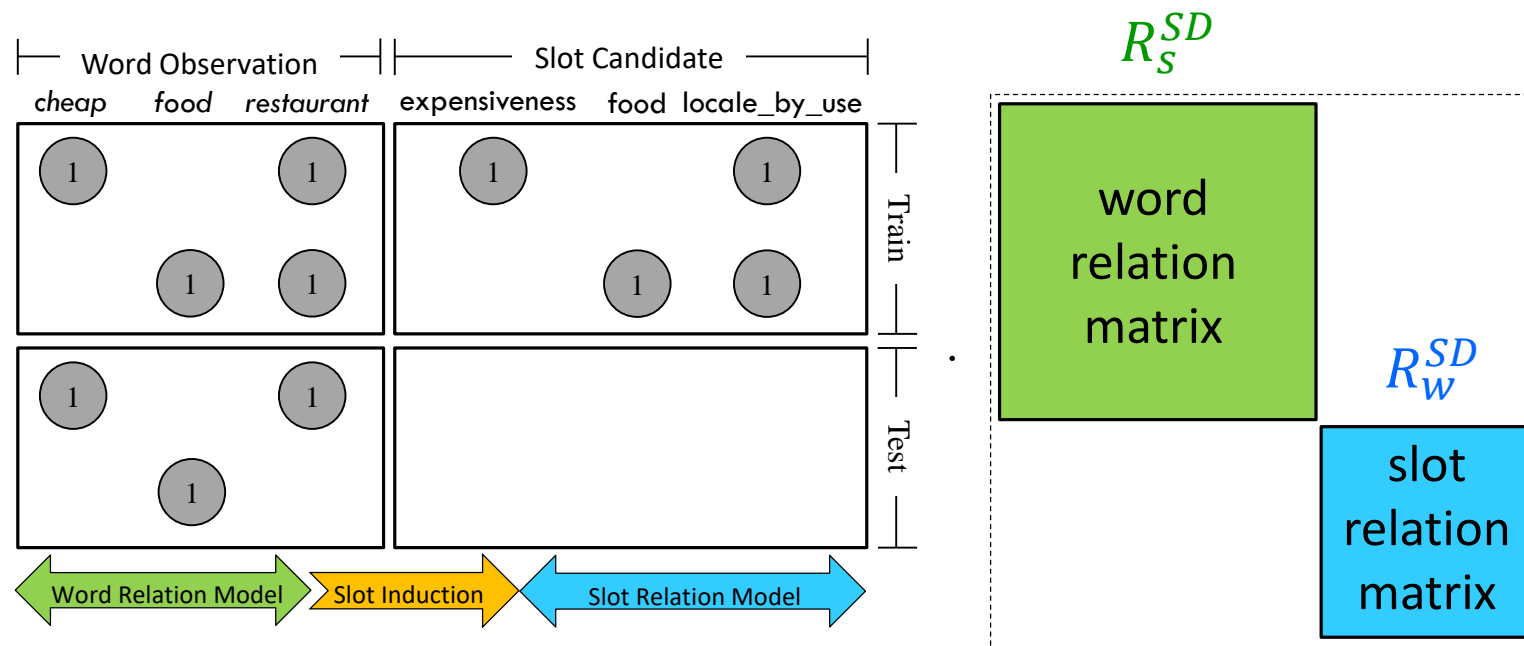
$$R_S^{SD} = R_S^S + R_S^D$$



$$R_W^{SD} = R_W^S + R_W^D$$



KNOWLEDGE GRAPH PROPAGATION MODEL



OUTLINE



Introduction



Semantic Decoding [ACL-IJCNLP'15]



- Ontology Induction



- Knowledge Graph Propagation



- **Matrix Factorization (for 2nd issue)**



- Experiments



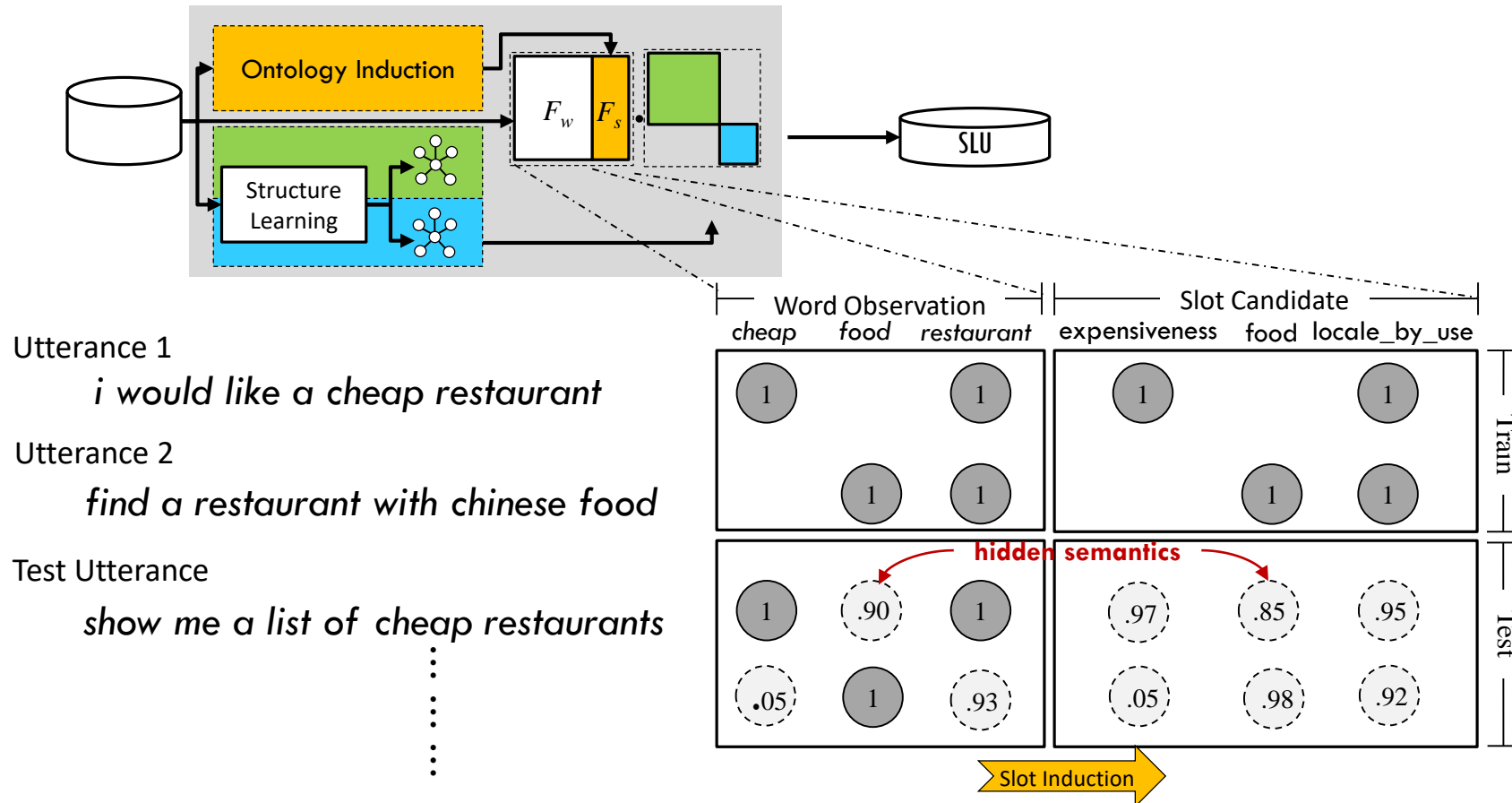
Future Work



Conclusions

MATRIX FACTORIZATION (MF)

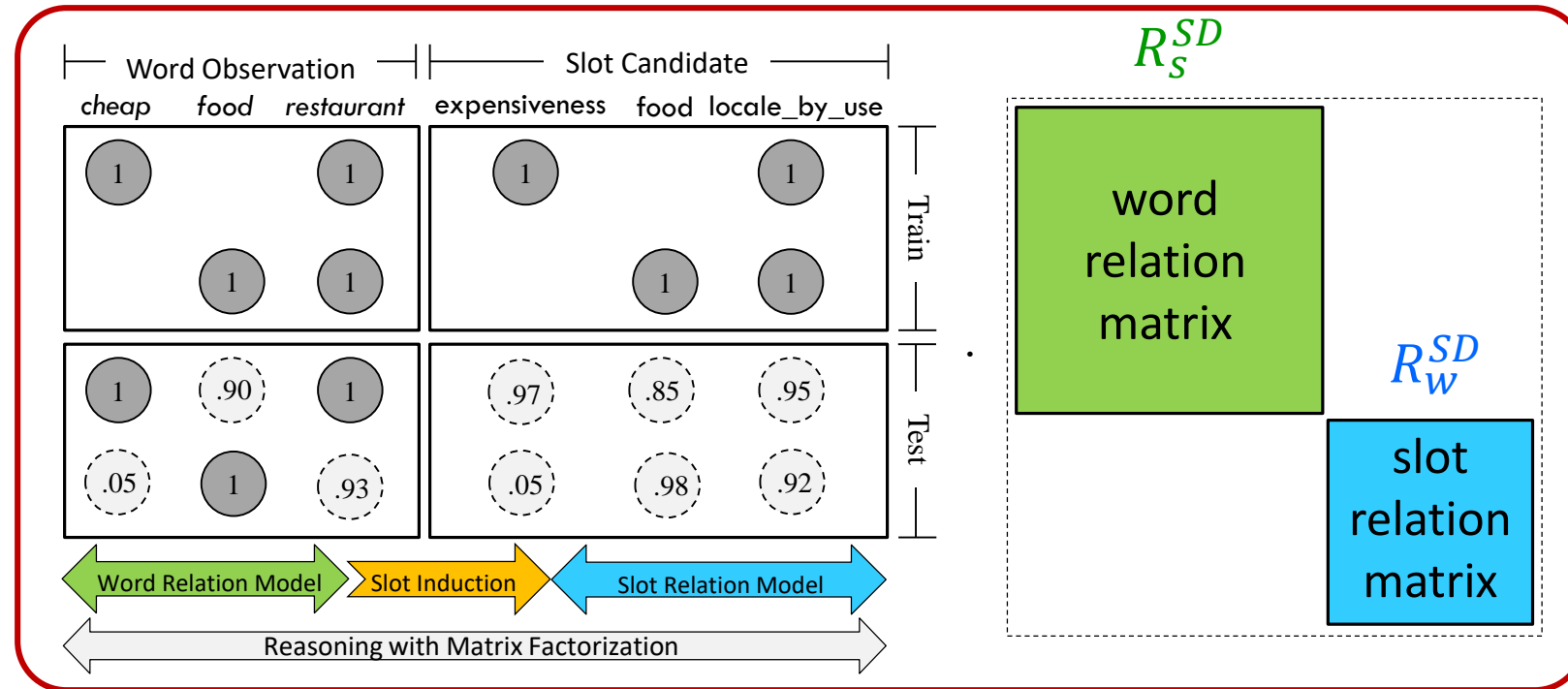
FEATURE MODEL



2nd Issue: hidden semantics cannot be observed but may benefit the understanding performance

2ND ISSUE: HOW TO LEARN THE IMPLICIT SEMANTICS?

MATRIX FACTORIZATION (MF)

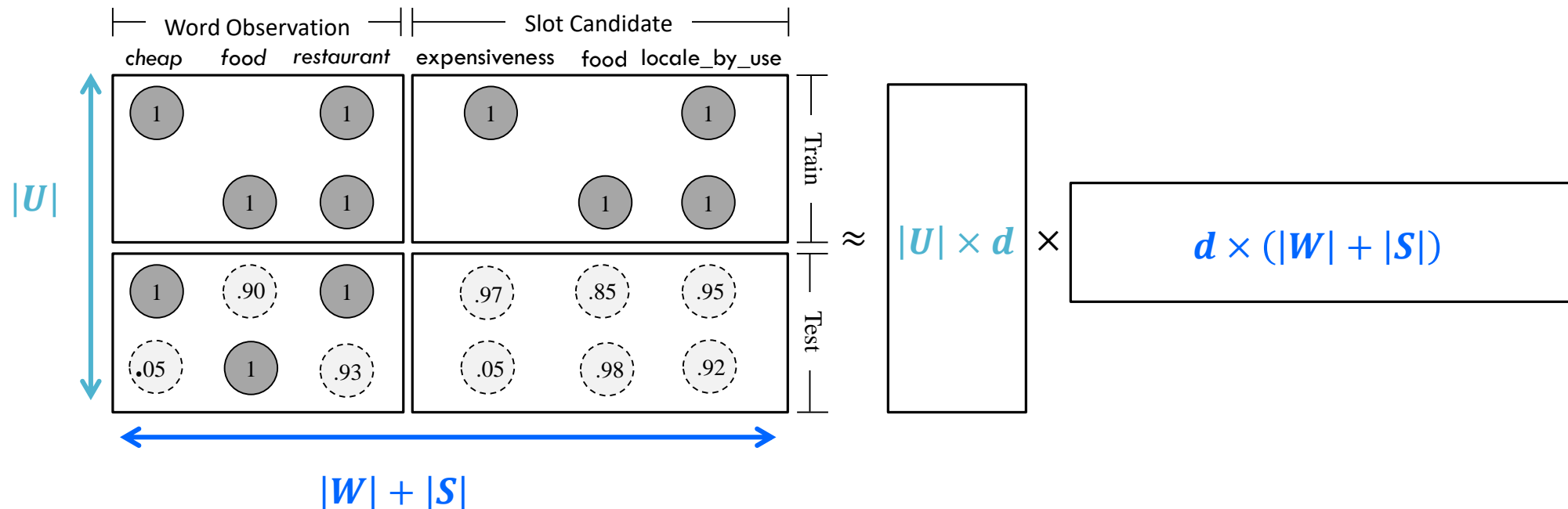


The MF method completes a partially-missing matrix based on the latent semantics by decomposing it into product of two matrices.

MATRIX FACTORIZATION (MF)

The decomposed matrices represent latent semantics for utterances and words/slots respectively

The product of two matrices fills the probability of hidden semantics



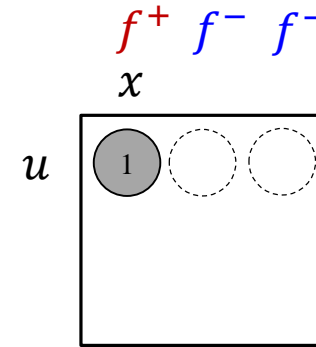
BAYESIAN PERSONALIZED RANKING FOR MF

Model implicit feedback

- not treat unobserved facts as negative samples (true or false)
- give observed facts higher scores than unobserved facts

$$\begin{aligned} f^+ &= \langle u, x^+ \rangle \\ f^- &= \langle u, x^- \rangle \end{aligned} \Rightarrow p(f^+) > p(f^-)$$

$$p(M_{u,x} = 1 \mid \theta_{u,x}) = \sigma(\theta_{u,x}) = \frac{1}{1 + \exp(-\theta_{u,x})}$$



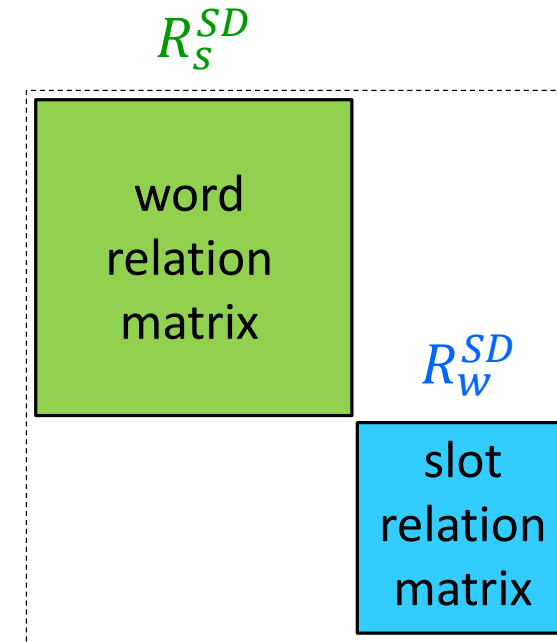
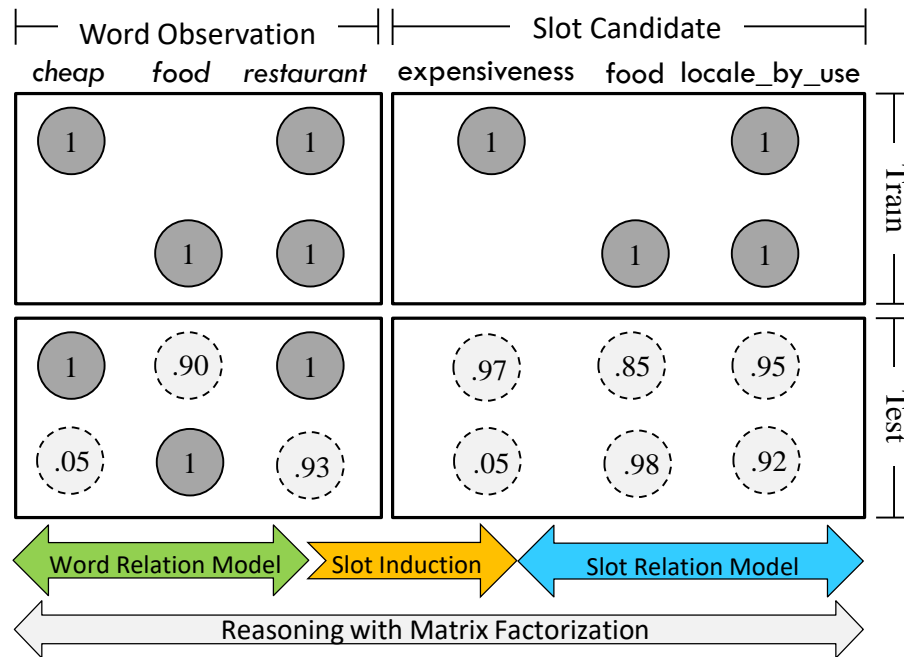
Objective:

$$\sum_{f^+ \in \mathcal{O}} \sum_{f^- \notin \mathcal{O}} \ln \sigma(\theta_{f^+} - \theta_{f^-})$$

The objective is to learn a set of well-ranked semantic slots per utterance.

2ND ISSUE: HOW TO LEARN THE IMPLICIT SEMANTICS?

MATRIX FACTORIZATION (MF)



The MF method completes a partially-missing matrix based on the latent semantics by decomposing it into product of two matrices.

OUTLINE



Introduction



Semantic Decoding [ACL-IJCNLP'15]



- Ontology Induction



- Knowledge Graph Propagation



- Matrix Factorization



- **Experiments**



Future Work

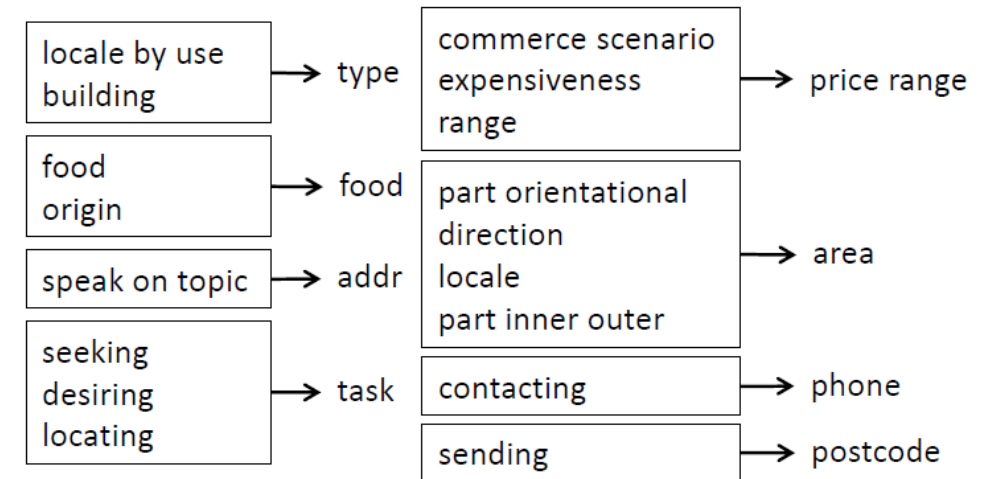


Conclusions

EXPERIMENTAL SETUP

Dataset

- Cambridge University SLU corpus [Henderson, 2012]
 - Restaurant recommendation in an in-car setting in Cambridge
 - WER = 37%
 - vocabulary size = 1868
 - 2,166 dialogues
 - 15,453 utterances
 - dialogue slot: **addr, area, food, name, phone, postcode, price range, task, type**



The mapping table between induced and reference slots



EXPERIMENT 1: QUALITY OF SEMANTICS ESTIMATION

Metric: Mean Average Precision (MAP) of all estimated slot probabilities for each utterance

Approach		ASR		Manual	
		w/o	w/ Explicit	w/o	w/ Explicit
Explicit	Support Vector Machine	32.5		36.6	
	Multinomial Logistic Regression	34.0		38.8	

EXPERIMENT 1: QUALITY OF SEMANTICS ESTIMATION



Metric: Mean Average Precision (MAP) of all estimated slot probabilities for each utterance

Modeling
Implicit
Semantics

Approach			ASR		Manual	
			w/o	w/ Explicit	w/o	w/ Explicit
Explicit	Support Vector Machine		32.5		36.6	
	Multinomial Logistic Regression		34.0		38.8	
Implicit	Baseline	Random				
		Majority				
	MF	Feature Model				
		Feature Model + Knowledge Graph Propagation				

EXPERIMENT 1: QUALITY OF SEMANTICS ESTIMATION



Metric: Mean Average Precision (MAP) of all estimated slot probabilities for each utterance

Approach			ASR		Manual	
			w/o	w/ Explicit	w/o	w/ Explicit
Explicit	Support Vector Machine		32.5		36.6	
	Multinomial Logistic Regression		34.0 		38.8 	
Implicit	Baseline	Random	3.4		2.6	
		Majority	15.4		16.4	
	MF	Feature Model	24.2		22.6	
		Feature Model + Knowledge Graph Propagation	40.5* (+19.1%)		52.1* (+34.3%)	

Modeling Implicit Semantics

EXPERIMENT 1: QUALITY OF SEMANTICS ESTIMATION

Metric: Mean Average Precision (MAP) of all estimated slot probabilities for each utterance

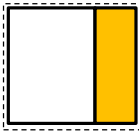
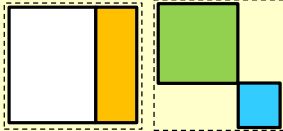
Approach			ASR		Manual	
			w/o	w/ Explicit	w/o	w/ Explicit
Explicit	Support Vector Machine		32.5		36.6	
	Multinomial Logistic Regression		34.0 		38.8 	
Implicit	Baseline	Random	3.4	22.5	2.6	25.1
		Majority	15.4	32.9	16.4	38.4
	MF	Feature Model	24.2	37.6*	22.6	45.3*
		Feature Model + Knowledge Graph Propagation	40.5* (+19.1%)	43.5* (+27.9%)	52.1* (+34.3%)	53.4* (+37.6%)

Modeling Implicit Semantics

The MF approach effectively models hidden semantics to improve SLU.

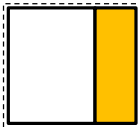
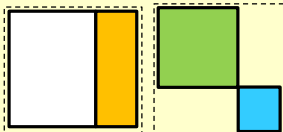
Adding a knowledge graph propagation model further improves the results.

EXPERIMENT 2: EFFECTIVENESS OF RELATIONS

Approach			ASR	Manual
Feature Model 			37.6	45.3
Feature + Knowledge Graph Propagation 	Semantic	$\begin{bmatrix} R_w^S & 0 \\ 0 & R_s^S \end{bmatrix}$	41.4*	51.6*
	Dependency	$\begin{bmatrix} R_w^D & 0 \\ 0 & R_s^D \end{bmatrix}$	41.6*	49.0*
	Word	$\begin{bmatrix} R_w^{SD} & 0 \\ 0 & 0 \end{bmatrix}$	39.2*	45.2
	Slot	$\begin{bmatrix} 0 & 0 \\ 0 & R_s^{SD} \end{bmatrix}$	42.1*	49.9*
	Both	$\begin{bmatrix} R_w^{SD} & 0 \\ 0 & R_s^{SD} \end{bmatrix}$		

All types of relations are useful to infer hidden semantics.

EXPERIMENT 2: EFFECTIVENESS OF RELATIONS

Approach			ASR	Manual
Feature Model 			37.6	45.3
Feature + Knowledge Graph Propagation 	Semantic	$\begin{bmatrix} R_w^S & 0 \\ 0 & R_s^S \end{bmatrix}$	41.4*	51.6*
	Dependency	$\begin{bmatrix} R_w^D & 0 \\ 0 & R_s^D \end{bmatrix}$	41.6*	49.0*
	Word	$\begin{bmatrix} R_w^{SD} & 0 \\ 0 & 0 \end{bmatrix}$	39.2*	45.2
	Slot	$\begin{bmatrix} 0 & 0 \\ 0 & R_s^{SD} \end{bmatrix}$	42.1*	49.9*
	Both	$\begin{bmatrix} R_w^{SD} & 0 \\ 0 & R_s^{SD} \end{bmatrix}$	43.5* (+15.7%)	53.4* (+17.9%)

All types of relations are useful to infer hidden semantics.

Combining different relations further improves the performance.

OUTLINE



Introduction



Semantic Decoding [ACL-IJCNLP'15]



- Ontology Induction



- Knowledge Graph Propagation



- Bayesian Personalized Ranking for Matrix Factorization



- Experiments



Future Work

Conclusions

LOW- AND HIGH-LEVEL UNDERSTANDING

Semantic concepts for individual utterances do not consider high-level semantics (user intents)

The follow-up behaviors are observable and usually correspond to user intents



"can i have a cheap restaurant"



price=*"cheap"*
target=*"restaurant"*

behavior=*navigation*

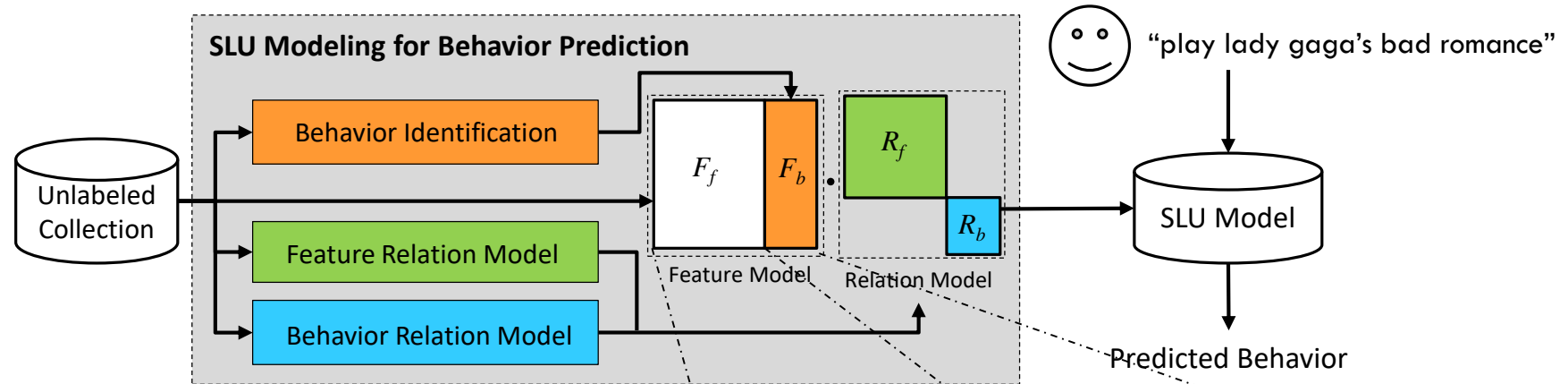
"i plan to dine in din tai fung tonight"



restaurant=*"din tai fung"*
time=*"tonight"*

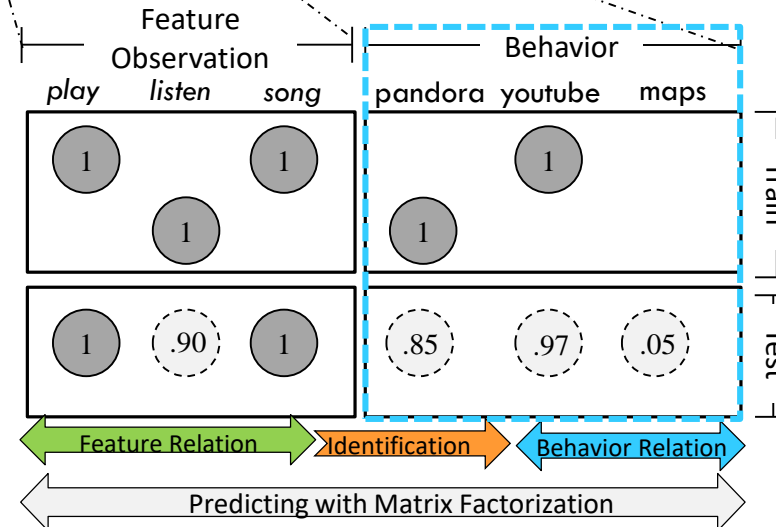
behavior=*reservation*

BEHAVIOR PREDICTION



Utterance 1
play lady gaga's song bad romance
 Utterance 2
i'd like to listen to lady gaga's bad romance

⋮



OUTLINE



Introduction



Semantic Decoding [ACL-IJCNLP'15]



- Ontology Induction



- Knowledge Graph Propagation



- Bayesian Personalized Ranking for Matrix Factorization



- Experiments



Future Work



Conclusions

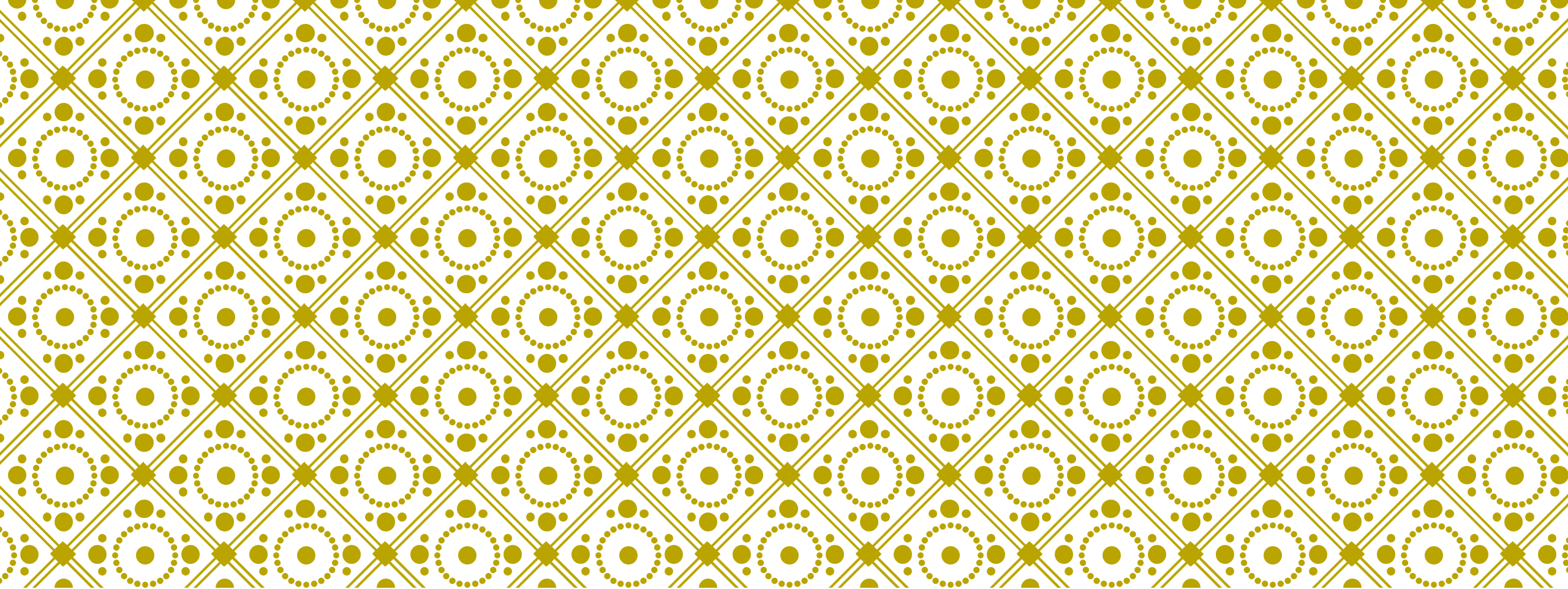
CONCLUSIONS

The ontology induction and knowledge graph construction enable systems to automatically acquire open domain knowledge.

The MF technique for SLU modeling provides a principle model that is able to unify the automatically acquired knowledge, and then allows systems to consider implicit semantics for better understanding.

- Better semantic representations for individual utterances
- Better follow-up behavior prediction

The work shows the feasibility and the potential of improving *generalization, maintenance, efficiency, and scalability* of SDSs.



Q & A

Thanks for your attentions!!

CAMBRIDGE UNIVERSITY SLU CORPUS

hi i'd like a restaurant in the cheap price range in the centre part of town

type=restaurant, pricerange=cheap, area=centre

um i'd like chinese food please

food=chinese

how much is the main cost

pricerange

okay and uh what's the address

addr

great uh and if i wanted to uh go to an italian restaurant instead

food=italian, type=restaurant

italian please

food=italian

what's the address

addr

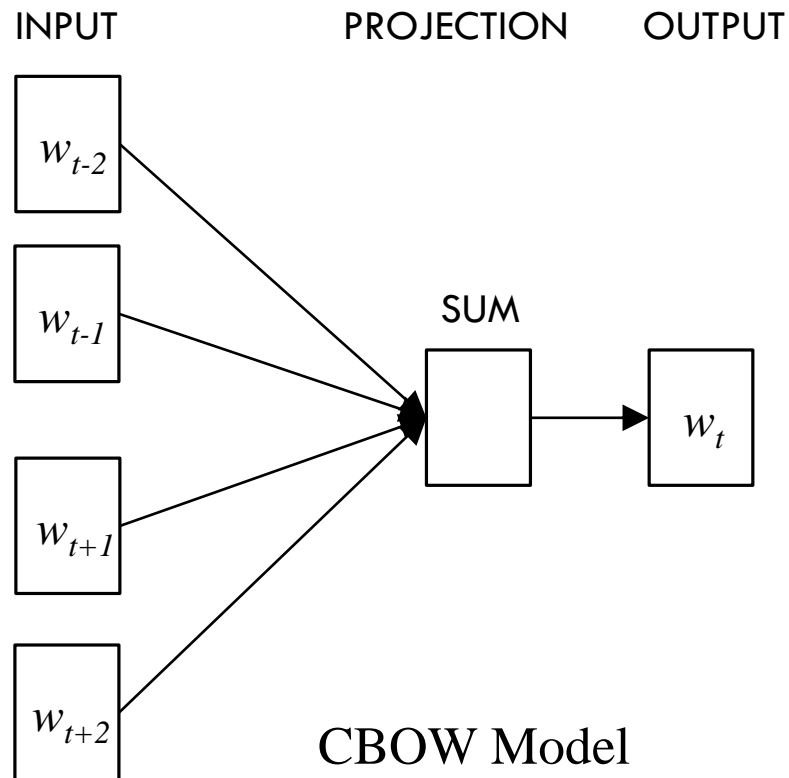
i would like a cheap chinese restaurant

pricerange=cheap, food=chinese, type=restaurant

something in the riverside

area=centre

WORD EMBEDDINGS



Training Process

- Each word w is associated with a vector
- The contexts within the window size c are considered as the training data D
- Objective function:

$$\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq i \leq c, i \neq 0} \log p(w_t | w_{t+i})$$

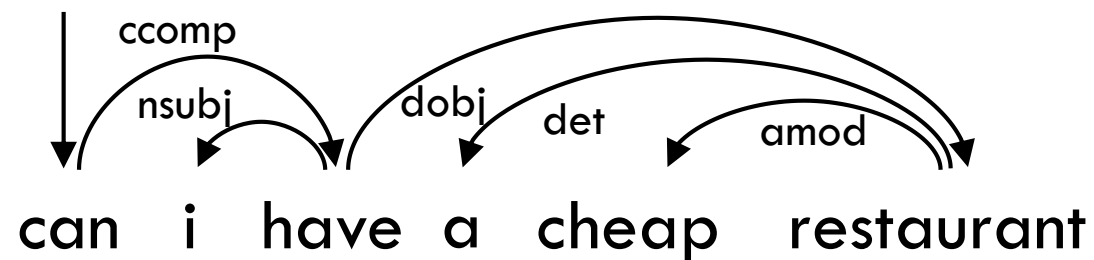
Mikolov et al., "Efficient Estimation of Word Representations in Vector Space," in *Proc. of ICLR*, 2013.

Mikolov et al., "Distributed Representations of Words and Phrases and their Compositionality," in *Proc. of NIPS*, 2013.

Mikolov et al., "Linguistic Regularities in Continuous Space Word Representations," in *Proc. of NAACL-HLT*, 2013.

DEPENDENCY-BASED EMBEDDINGS

Word & Context Extraction



Word	Contexts
can	have/ccomp
i	have/nsubj ⁻¹
have	can/ccomp ⁻¹ , i/nsubj, restaurant/dobj
a	restaurant/det ⁻¹
cheap	restaurant/amod ⁻¹
restaurant	have/dobj ⁻¹ , a/det, cheap/amod

DEPENDENCY-BASED EMBEDDINGS

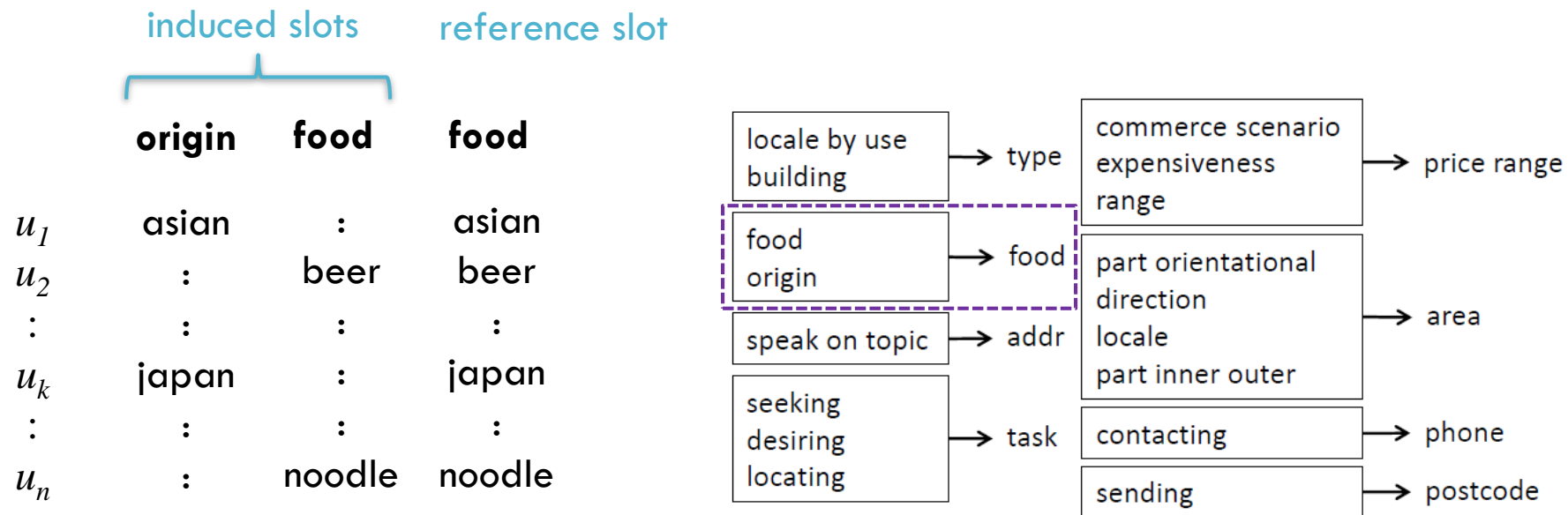
Training Process

- Each word w is associated with a vector v_w and each context c is represented as a vector v_c
- Learn vector representations for both words and contexts such that the dot product $v_w \cdot v_c$ associated with **good** word-context pairs belonging to the training data D is maximized
- Objective function:

$$\arg \max_{v_w, v_c} \sum_{(w, c) \in D} \log \frac{1}{1 + \exp(-v_c \cdot v_w)}$$

SLOT MAPPING TABLE

Create the mapping if slot fillers of the induced slot are included by the reference slot



SEMAFOR PERFORMANCE

The SEMAFOR evaluation

Table 5

Frame identification results on both the SemEval 2007 data set and the FrameNet 1.5 release. Precision, recall, and F_1 were evaluated under exact and partial frame matching; see Section 3.3. **Bold** indicates best results on the SemEval 2007 data, which are also statistically significant with respect to the baseline ($p < 0.05$).

FRAME IDENTIFICATION (§5.2)		exact matching			partial matching		
		P	R	F_1	P	R	F_1
SemEval 2007 Data	gold targets	60.21	60.21	60.21	74.21	74.21	74.21
	automatic targets (§4)	69.75	54.91	61.44	77.51	61.03	68.29
	J&N'07 targets	65.34	49.91	56.59	74.30	56.74	64.34
	Baseline: J&N'07	66.22	50.57	57.34	73.86	56.41	63.97
FrameNet 1.5 Release	gold targets	82.97	82.97	82.97	90.51	90.51	90.51
	– unsupported features	80.30	80.30	80.30	88.91	88.91	88.91
	& – latent variable	75.54	75.54	75.54	85.92	85.92	85.92