

# Large-scale Semantic Parsing via Schema Matching and Lexicon Extension

**Qingqing Cai and Alexander Yates**  
**Temple University**

# What is semantic parsing?



Who directed  
the Titanic?

Semantic Parsing

$\lambda x. \text{directed\_by}(\text{titanic}, x)$

## Combinatory Categorial Grammar (CCG) Semantic Parsing

### CCG Category

Syntax

Semantic: lambda expression

S/(S\NP) :

$\lambda f. \lambda w. \lambda z. f(w, z)$

Who

Lexical Entry

S :

$\lambda x. \text{directed\_by}(\text{titanic}, x)$

S\NP :

$\lambda x. \text{directed\_by}(\text{titanic}, x)$

S\NP/NP :

$\lambda x. \lambda y. \text{directed\_by}(x, y)$

directed

Operation

Application

Composition

Type Shifting

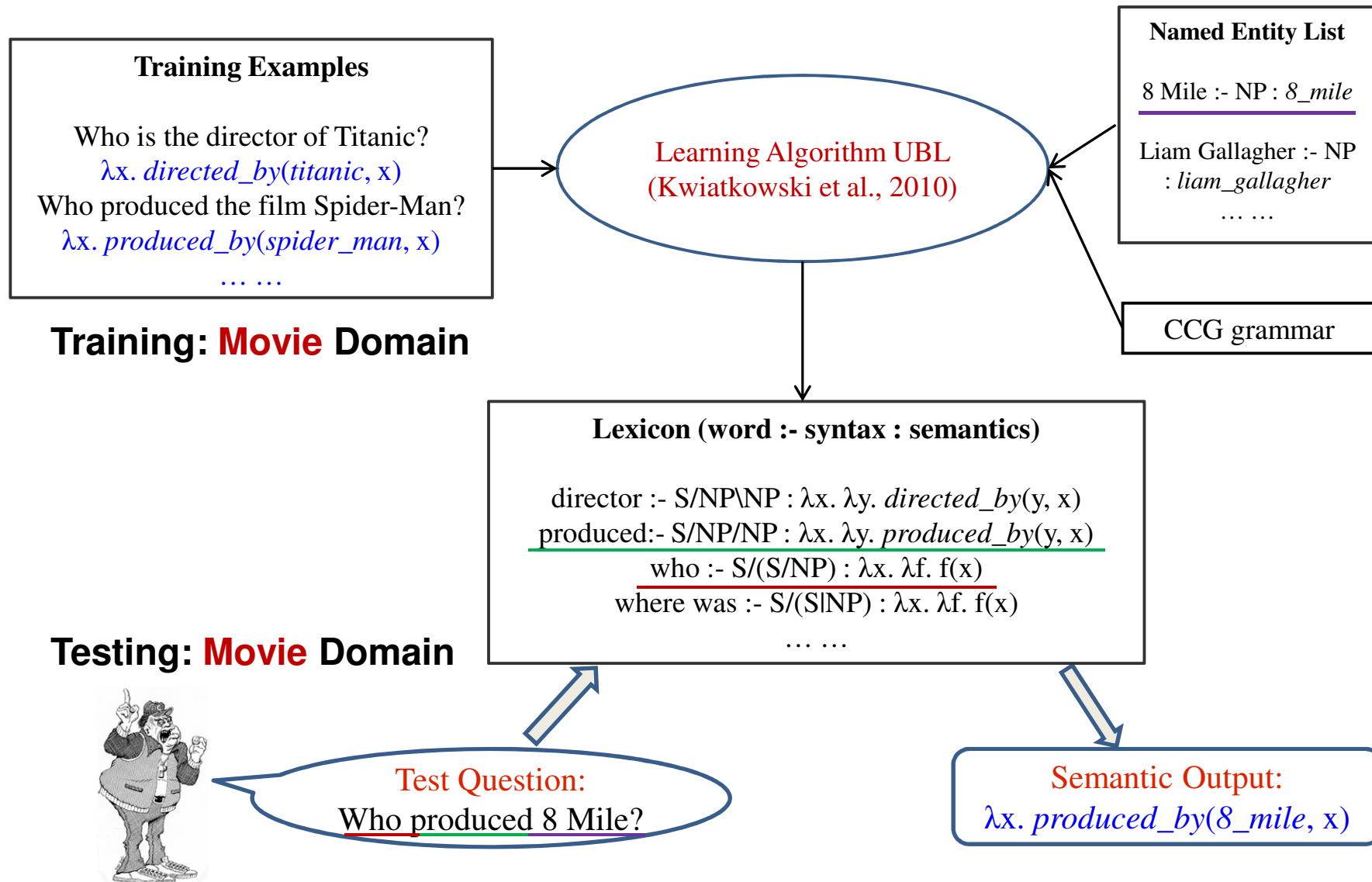
... ..

Forward application

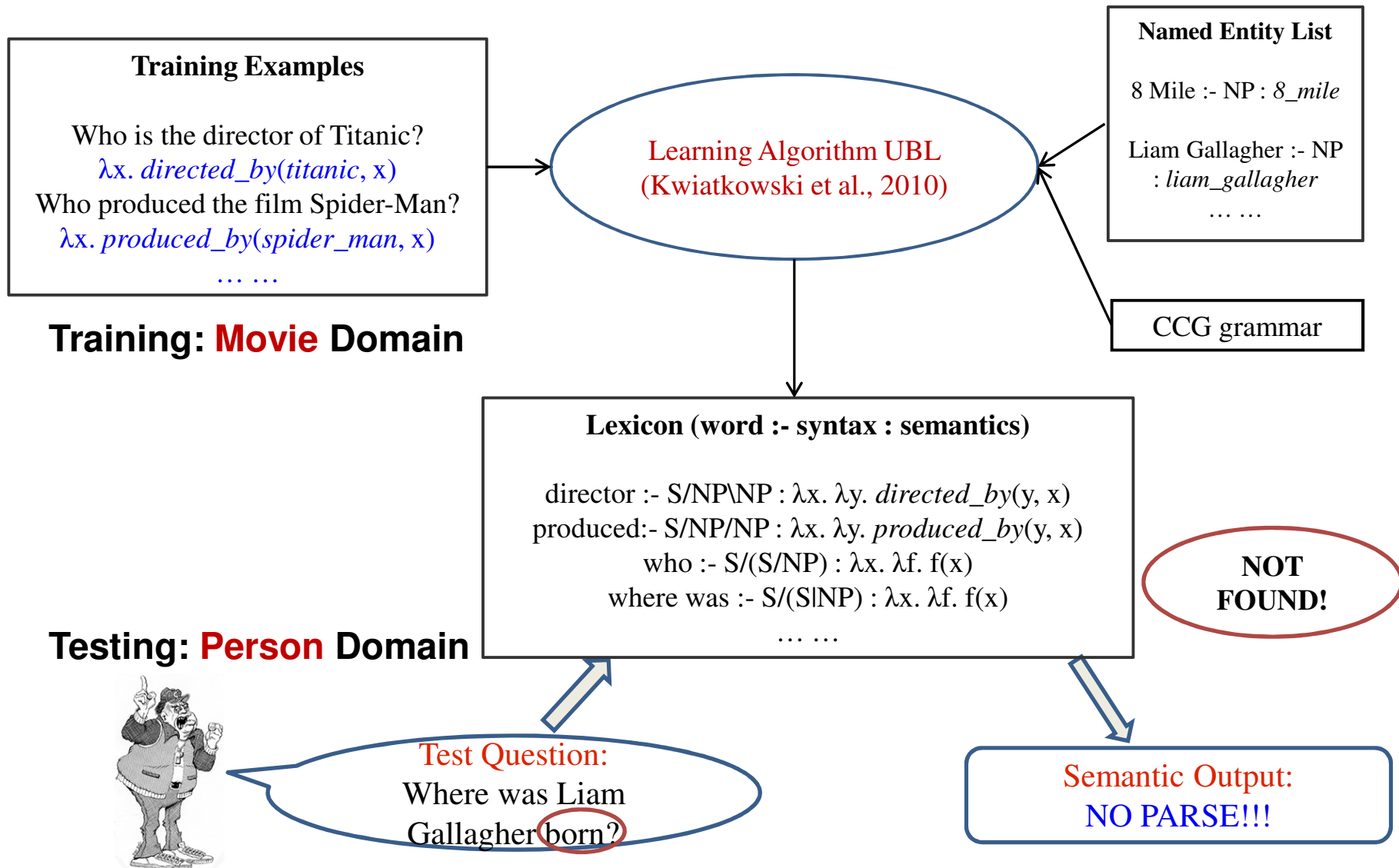
NP : *titanic*

Titanic

# The problem



# The problem



# The Problem

Where was Liam Gallagher born?

born  $\vdash$  S/NP :  $\lambda x. \lambda y. \text{place\_of\_birth}(x,y)$

Not found in training data!!!

- Why can't we correctly parse the questions?
- Why can't we get the correct lexical entry?

## Schema Matching

born  $::$  *place\_of\_birth*

## Lexicon Extension

born  $::$  S/NP :  $\lambda x. \lambda y. \lambda p. p(x,y)$

born  $\vdash$  S/NP :  $\lambda x. \lambda y. \text{place\_of\_birth}(x,y)$

# System

- **Schema Matching - Matcher**
- **Lexicon Extension - Lextender**

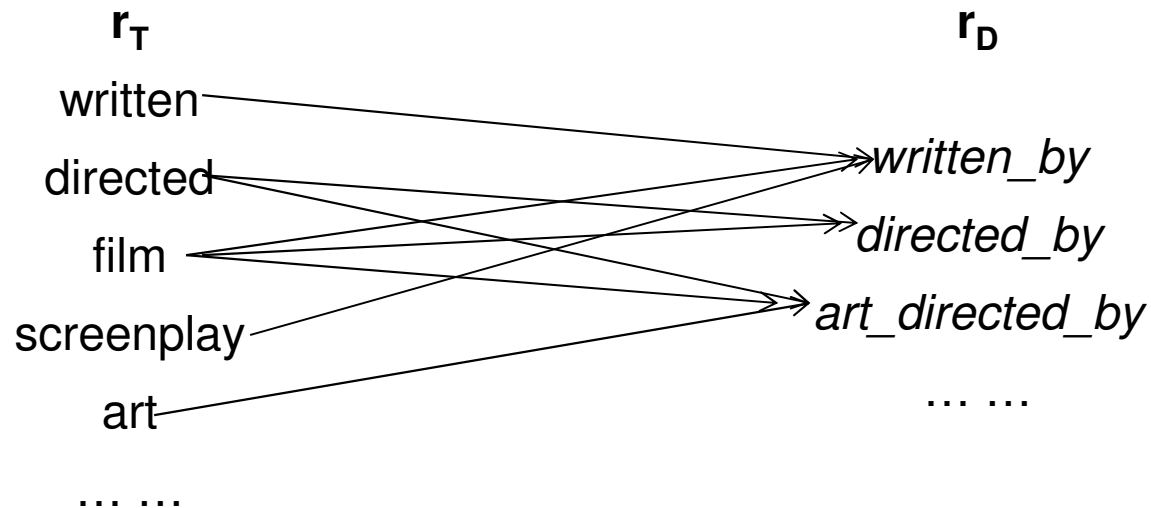
## MATCHER - Task

- **Task**

identify natural language words and phrases ( $r_T$ ) which correspond with relations and entities in a fixed schema ( $r_D$ )

Freebase

- **Example**



# MATCHER - Method

- **Step1: Identify candidate matches**

- **Step2: Matching**

- **Pattern-based matching selection**

- **Instance-based matching selection**

**Based on Logistic Regression**  
Trained on golden-standard matches



# MATCHER – Identify candidate matches

$r_D = \text{“film\_music”}$

## Entities for “film\_music”, from Freebase

| arg1: film           | arg2: contributor |
|----------------------|-------------------|
| Spider-man           | Danny Elfman      |
| Twilight             | Carter Burwell    |
| Sleepless in Seattle | Marc Shaiman      |
| Forrest Gump         | Alan Silvestri    |
| ...                  | ...               |

(1) Generate queries

Eg: “Spider-man” AND “Danny Elfman”

(2) Issue queries to the Web

## Candidates for film\_music

composed:: 0.4  
 directed :: 0.35  
 download :: 0.34  
 film :: 0.32  
 music:: 0.3  
 composers:: 0.21

(1) Remove stop-words and entity-names;

(2) Rank remaining words by frequency;

## Matching snippets for “film\_music”:

- (1) ~~To listen Spiderman Theme Danny Elfman music just click play to download.~~
- (2) ~~Twilight was directed by Catherine Hardwicke and composed by Carter Burwell...~~
- (3) ... ..

# MATCHER - Method

- **Step1: Identify candidate matches**

- **Step2: Matching**

- **Pattern-based matching selection**

- **Instance-based matching selection**

**Based on Logistic Regression**  
Trained on golden-standard matches

## MATCHER – Pattern-based match selection

- **Motivation**

relation words appear in **simple, unambiguous** grammatical **patterns** that connects  $r_T$  with entities from  $r_D$

$r_D = \text{“}film\_music\text{”}$

$r_T = \{ \text{“}composed\text{”}, \text{“}directed\text{”}, \text{“}download\text{”}, \dots \}$

$e$ : an entity from  $r_D$ , with a specific type *contributor*

**$\#(\text{composed by } e) > 100 * \#(\text{directed by } e)$**

# MATCHER – Pattern-based match selection

## • Pattern

| Pattern                     | Condition   | Example                       |
|-----------------------------|---|-------------------------------|
| “r <sub>T</sub> in E”       | r <sub>T</sub> ends with “-ed”, and E has type <i>datetime</i> or <i>location</i> | “founded in 1989”             |
| “r <sub>T</sub> by E”       | r <sub>T</sub> ends with “-ed”  | “invented by Edison”          |
| “r <sub>T</sub> such as E”  | r <sub>T</sub> ends with “-s”   | “directors such as Tarantino” |
| “E is a(n) r <sub>T</sub> ” | all cases   | “Paul Rudd is an acotr”       |

## • Why not use more patterns?

To cover all of Freebase, Machter needs **40 million** queries, or over **1.25 years** if issues 1 query per second

# MATCHER – Pattern-based match selection

## Candidates for *film\_music*

composed:: 0.4  
 directed :: 0.35  
 download :: 0.34  
 film :: 0.32  
 music:: 0.3  
 composers:: 0.21



| Pattern            |
|--------------------|
| “ $r_T$ in E”      |
| “ $r_T$ by E”      |
| “ $r_T$ such as E” |
| “E is a(n) $r_T$ ” |

## (1) Generate pattern queries

Eg: “composed by Danny Elfman”

“composers such as Alan Silvestri”

“Alan Silvestri is a film”

## (2) Issue pattern queries to the Web

| ( $r_T, r_D$ ) \ pattern        | $r_T$ in E | $r_T$ by E | $r_T$ such as E | $r_T$ is a(n) E | total | POS of $r_T$ | ... |
|---------------------------------|------------|------------|-----------------|-----------------|-------|--------------|-----|
| (composed, <i>film_music</i> )  | 0          | 343        | 0               | 0               | 343   | V            | ... |
| (directed, <i>film_music</i> )  | 0          | 30         | 0               | 0               | 20    | V            | ... |
| (film, <i>film_music</i> )      | 0          | 0          | 0               | 431             | 431   | N            | ... |
| (composed, <i>directed_by</i> ) | 0          | 21         | 0               | 0               | 21    | V            | ... |
| ...                             | ...        | ...        | ...             | ...             | ...   | ...          | ... |

$c(p, r_D, r_T)$ : number of instances where  $r_T$  appears in  $p$  that involves entities from  $r_D$

Logistic Regression Model

## Candidates for *film\_music*

composed:: 0.5  
 directed :: 0.4  
 film :: 0.4  
 composers:: 0.37  
 music:: 0.25  
 download :: 0.14

# MATCHER - Method

- **Step1: Identify candidate matches**

- **Step2: Matching**

- **Pattern-based matching selection**

- **Instance-based matching selection**

**Based on Logistic Regression**  
Trained on golden-standard matches

# MATCHER – Instance-based match selection

## • Motivation

*film\_music* : composed, directed, film, ...

Freebase

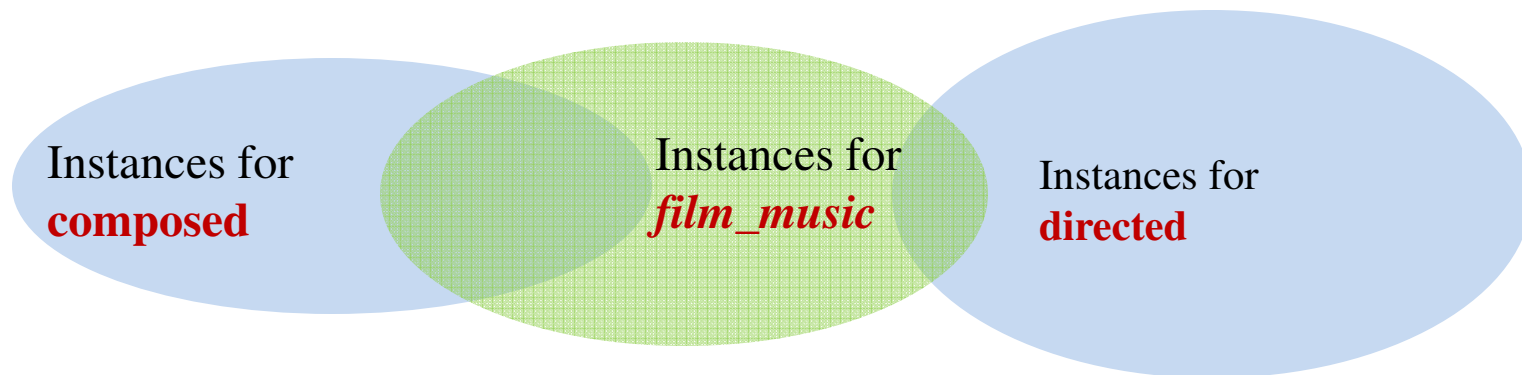
Open IE: Reverb

Open IE: Reverb



| composed            |                    | <i>film_music</i>   |                    | directed            |                    |
|---------------------|--------------------|---------------------|--------------------|---------------------|--------------------|
| <i>film</i>         | <i>contributor</i> | <i>film</i>         | <i>contributor</i> | <i>film</i>         | <i>contributor</i> |
| Million Dollar Tree | Clint Eastwood     | Million Dollar Tree | Clint Eastwood     | Million Dollar Tree | Clint Eastwood     |
| The Crack           | Camilo Sanabria    | The Crack           | Camilo Sanabria    | The Brown Bunny     | Mr. Gallo          |
| ...                 | ...                | George Harrison     | Water              | Titanic             | James Cameron      |
| ...                 | ...                | ...                 | ...                | ...                 | ...                |

## MATCHER – Instance-based match selection



$$PMI(r_D, r_T) = \frac{|I_D(r_D) \cap I_T(r_T)|}{|I_D(r_D)| \cdot |I_T(r_T)|}$$

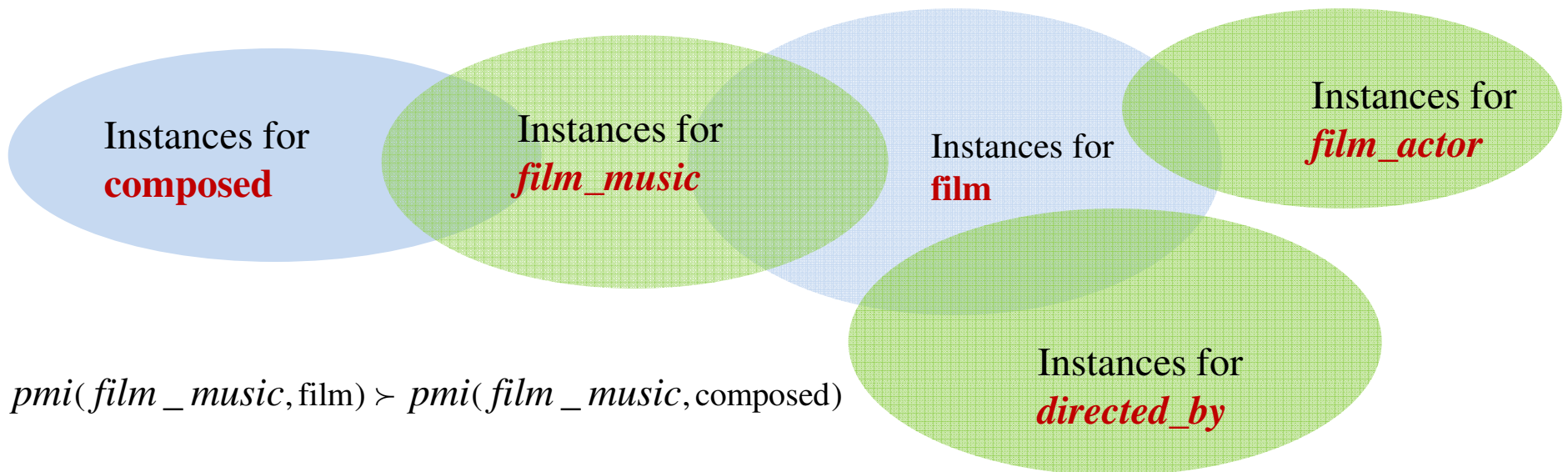
$I_D(r_D)$  : instances extracted from database relations

$I_T(r_T)$  : instances extracted from textual relations

$$pmi(film\_music, composed) \succ pmi(film\_music, directed)$$



# MATCHER – Instance-based match selection



$$pmi(film\_music, film) \succ pmi(film\_music, composed)$$

$$pmi\_polysemy\_penalty(r_D, r_T) = \frac{pmi(r_D, r_T)}{\sum_{r'_D} pmi(r'_D, r_T)}$$

$$pmi\_polysemy\_penalty(film\_music, composed) \succ pmi\_polysemy\_penalty(film\_music, film)$$

## MATCHER – Instance-based match selection

**Candidates for *film\_music***

composed:: 0.5  
directed :: 0.4  
film :: 0.4  
composers:: 0.37  
music:: 0.25  
download :: 0.14

**Matcher Statistics**

$pmi(r_D, r_T)$   
 $pmi\_polysemy\_penalty(r_D, r_T)$

**Logistic Regression  
Model****Candidates for *film\_music***

composed:: 0.6  
composers:: 0.51  
music:: 0.21  
film :: 0.2  
directed :: 0.13  
download :: 0.09

## Our System's Architecture

- **Schema Matching - Matcher**
- **Lexicon Extension - Lextender**

## LEXTENDER – Task

- **Input**

Matching **M** together with **scores** for each pair  $(r_D, r_T)$  in **M**  
Lexical entries learned from UBL

- **Output**

Predict the syntactic category *Syn*, lambda-calculus semantics *Sem*,  
and **weight** for a full lexical entry **L** for each pair in **M**

# LEXTENDER – Method

features: pos of  $r_T$ , suffix of  $r_T$ , #arguments of  $r_D$ ,  
argument types of  $r_D$

|                                 |
|---------------------------------|
| <b><math>(r_D, r_T)</math></b>  |
| <i>(located_in, city)</i>       |
| <i>(award_winner_time, won)</i> |

## Naïve Bayes Classifier

| features  | matching_score | classification_score | ... |
|---|----------------|----------------------|-----|
| <i>(located_in, city)</i> :<br>S\NP/NP : $\lambda p \lambda x \lambda y. p(x, y)$             | 0.7            | 2.1                  | ... |
| <i>(award_winner_time, won)</i><br>(S\NP\NP)/NP : $\lambda p \lambda x \lambda y. p(x, y, z)$ | 0.81           | 3.12                 | ... |

|  |
|--|
| <b><math>(r_D, r_T) : \text{syn} : \text{sem} : \text{weight}</math></b>                                 |
| <i>(located_in, city)</i>  |
| <i>(award_winner_time, won)</i><br>(S\NP\NP)/NP : $\lambda p \lambda x \lambda y. p(x, y, z)$ :<br>0.217 |

|   |
|---|
| <i>(award_winner_time, won)</i><br>(S\NP\NP)/NP : $\lambda p \lambda x \lambda y. p(x, y, z)$ |
|---|

| features                        | pos of $r_T$ | # args | type of arg1    | type of arg2        | ... |
|---------------------------------|--------------|--------|-----------------|---------------------|-----|
| <i>(located_in, city)</i>       | N            | 2      | <i>location</i> | <i>location</i>     | ... |
| <i>(award_winner_time, won)</i> | V            | 3      | <i>award</i>    | <i>Award_winner</i> | ... |

## Data Set

**•Dataset**

917 natural language questions and logical forms  
81 domains and 635 relations

**•Split into three groups for cross-domain validation**

| Group1 |     | Group2       |     | Group3   |     |
|--------|-----|--------------|-----|----------|-----|
| Film   | 49  | business     | 46  | Tv       | 34  |
| people | 29  | location     | 29  | award    | 32  |
| book   | 21  | organization | 24  | medicine | 24  |
| ...    | ... | ...          | ... | ...      | ... |
| Total  | 300 | total        | 310 | total    | 307 |

## Experiment 1 – Matcher Test

- **Setting**

**Frequency:** get candidates

**Pattern:** pattern-based matching selection

**Extractions:** instance-based matching selection

**Matcher:** Pattern + Extractions Selection

- **Evaluation**

M: ( $r_D, r_T$ ) matches predicted by models

G: matches in the gold-standard manual data

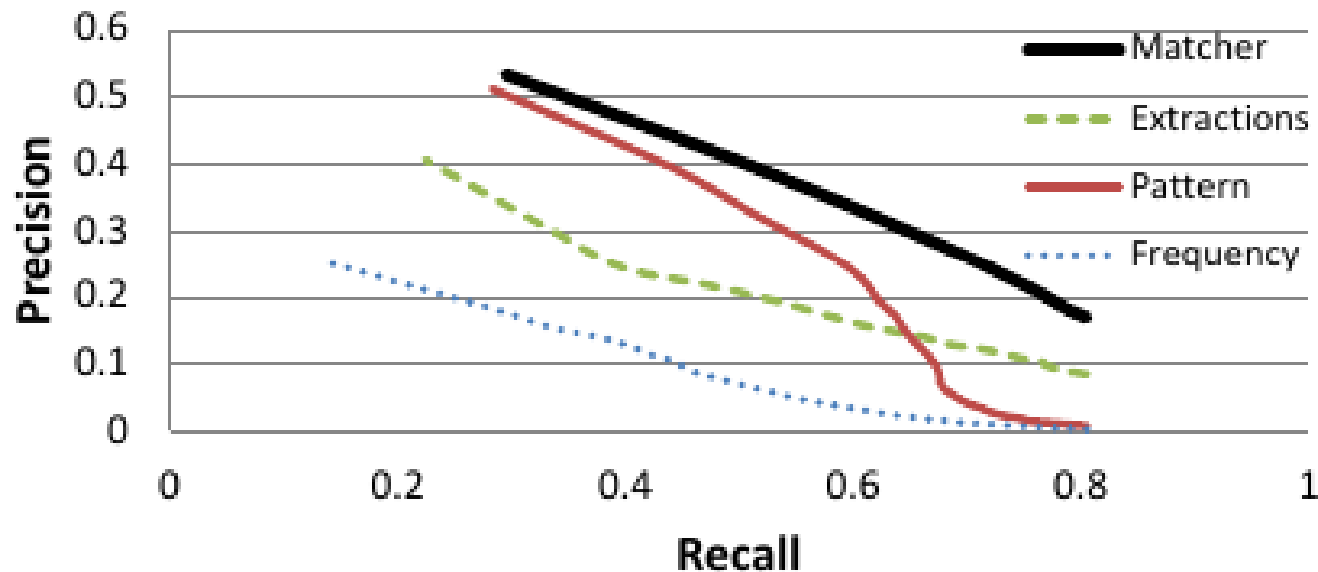
$$\text{precision} : P = \frac{|M \cap G|}{|M|}$$

$$\text{recall} : R = \frac{|M \cap G|}{|G|}$$

## Experiment 1 – Matcher Test

- **Performance**

### Alignment Predictions



- **Lessons**

- (1) Pattern, Extractions, Matcher > Frequency
- (2) Matcher yields the best precision at all levels of recall



## Experiment 2 – Semantic Parsing Test

- **Setting**

**70/30 split and cross-domain validation**

**Baseline:** UBL (supervised semantic parsing)

**Frequency:** get candidates + LEXTENDER + UBL

**Pattern:** pattern-based matching selection + LEXTENDER + UBL

**Extractions:** instance-based matching selection + LEXTENDER + UBL

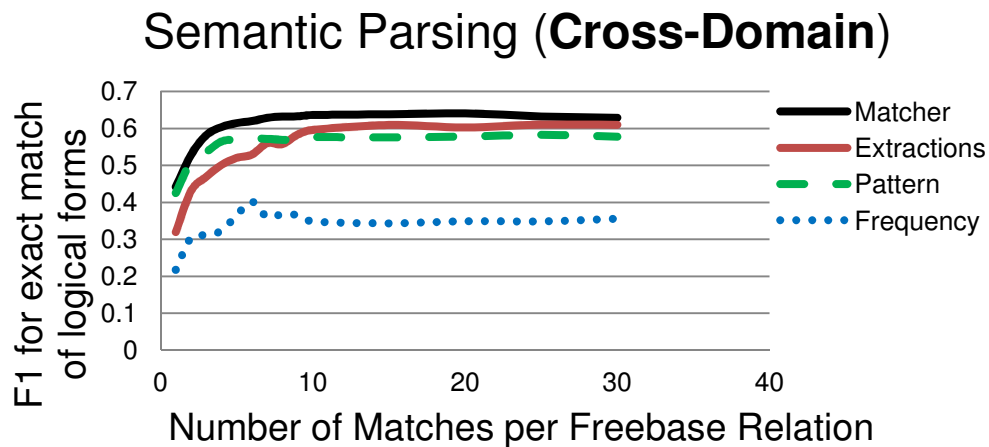
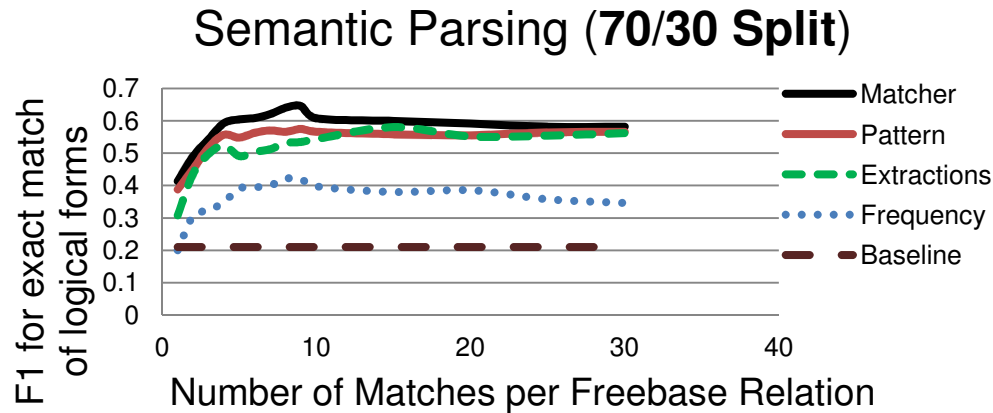
**Matcher:** Pattern + Extractions + LEXTENDER + UBL

- **Evaluation**

Exact match between top-scoring and manually-produced logical forms

# Experiment 2 – Semantic Parsing Test

## •Performance



## Lessons:

- (1) Matcher, Pattern, Extractions, Frequency > Baseline;
- (2) Matcher > Frequency, Pattern, Extractions

# Conclusion

- **Problem**

- (1) How to handle testing questions which are **significantly different** from labeled training questions

- **Solution**

### **MATCHER + LEXTENDER**

- (1) MATCHER: find correspondences between words and database symbols
- (2) LEXTENDER: integrate (word, database symbol) correspondences into a semantic parsing lexicon

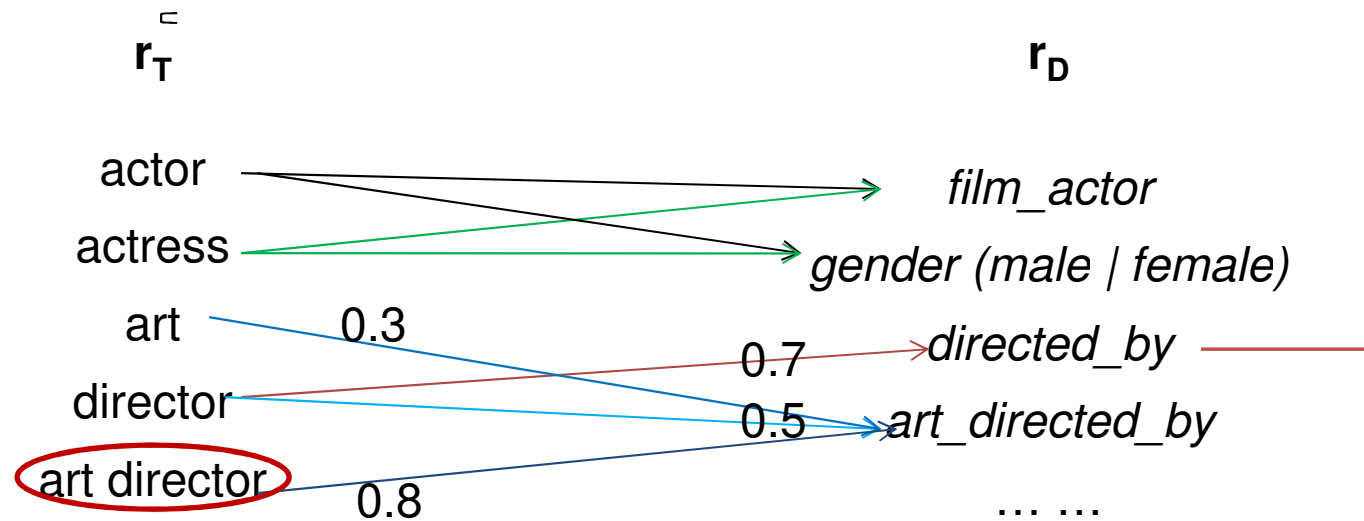
- **Result**

MATCHER + LEXTENDER + UBL: semantic parsers can handle Freebase questions on a large variety of domains with an F1 of **0.63**

## Future Work

- **Future work**

(1) Handle more complex matches between database and textual relations



(2) Handle more complex natural language questions

who is the **art director** of the Titanic?

who is the **director** of the Titanic?

Thanks for your attention!