# Question Answering Over Knowledge Graph

## Lei Zou

# Knowledge Graph

Google launches **Knowledge Graph** project at 2012.

# Knowledge Graph

Essentially, KG is a sematic network, which models **the entities (including properties) and the relation between each other.**

# Resource Description Framework (RDF)

- RDF is an **de facto standard** for Knowledge Graph (KG).
- RDF is a **language** for the conceptual modeling of information about web resources
- A **building block** of semantic web
- Make the information on the web and the interrelationships among them "**Machine Understandable**"

# RDF & SPARQL

## RDF Datasets

| Subject | Predicate | Object |
|---------|-----------|--------|
| Resident_Evil:_Retribution | type | film |
| Resident_Evil:_Retribution | budget | "6.5E7" |
| Resident_Evil:_Retribution | director | Paul_W._S._Anderson |
| Paul_W._S._Anderson | type | director |
| Resident_Evil | director | Paul_W._S._Anderson |
| Paul_Anderson_(actor) | type | actor |
| The_Revenant | strarring | Philadelphia |
| Priestley Medal | awards | Paul S. Anderson |
| Maclovia_(1948_film) | distributor | Filmex |

"What is the budget of the film directed by Paul Anderson ?."

## SPARQL

SELECT ?y WHERE
{
?x director Paul_W._S._Anderson .
?x type film .
?x budget ?y.
}

# Interdisciplinary Research

**Database**

RDF Database

Data Integration 、 Knowledge Fusion

**Natural Language Processing**

Information Extraction

Semantic Parsing



**Machine Learning**

Knowledge Representation (Graph Embedding)

**Knowledge Engineering**

KB construction

Rule-based Reasoning

# KG-based Question/Answering

- SPARQL syntax are too complex for ordinary users
- RDF KG is "schema-less" data, not like schema-first relational database.

# KG-based Question/Answering

- An **Easy-to-Use** Interface to Access Knowledge Graph
- It is interesting to both **academia** and **industry**.
- **Interdisciplinary research** between database and NLP (natural language processing) communities.

# KG-based Question/Answering



Oren Etzioni, AAAI Fellow

"(Researchers) They must invest much more in bold strategies that can achieve **natural-language searching and answering**"
---Oren Etzioni, Search needs a shake up, **NATURE**, Vol 476, p25-26, 2011.

# Facebook **Graph Search**

## "My friends who live in Canada"

" Facebook Graph Search"
------announced by Mark
Zuckerberg on January 16,
2013

# Facebook Graph Search

"Photos of my friends who live in Canada"

# EVI---(originally, True Knowledge)

where is the capital of united states ?

You asked: where is the capital of united states

Washington, D.C..

→ website   → wikipedia

Washington, D.C.

| | Venture Capital |
|---|---|
| 2007-09 | 1.2 Million USD |
| 2008-07 | 4 Million USD |
| 2012-01 | Acquired by Amazon |

William Tunstall-Pedoe: *True Knowledge: Open-Domain Question Answering using Structured Knowledge and Inference*. AI Magazine 31(3): 80-92 (2010)

# KG-based Question/Answering

- Information Retrieval-based
  - Generate candidate answers
  - Ranking

- Semantic Parsing-based
  - Translate NLQ to logical forms
  - Executing

# Knowledge-based QA (KB-QA)

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



(Cite: Nan Duan, MSRA)

# KG-based Question/Answering

- Information Retrieval-based



"What is the budget of the film directed by Paul Anderson?"

film ← type ← Resident_Evil → director → Paul.W.S.Anderson → type → director

Resident_Evil → budget → "6.5E7"

# KG-based Question/Answering

- Information Retrieval-based



"What is the budget of the film directed by Paul Anderson?"

film

director

type

type

Resident_Evil → director → Paul.W.S.Anderson

Step. 1
Mentioned entity

budget

"6.5E7"

Step. 2
Candidate Answer Selection
(within 2-hops)

Step. 3
Ranking Answers

"6.5E7"

# Question Answering with Subgraph Embeddings [Bordes et al. EMNLP 2014]



Figure 1: Illustration of the subgraph embedding model scoring a candidate answer: (i) locate entity in the question; (ii) compute path from entity to answer; (iii) represent answer as path plus all connected entities to the answer (the subgraph); (iv) embed both the question and the answer subgraph separately using the learnt embedding vectors, and score the match via their dot product.

# Question Answering with Subgraph Embeddings
# [Bordes et al. EMNLP 2014]

Let $W$ be a matrix $\Re^{k \times N}$

k:  the dimension of the embedding space

N:  $N = N_W + N_S$

$N_W$  is the number of words

$N_S$  is the number of entities and relation
types

Embedding a question q

$$f(q) = W\phi(q)$$

$\phi(q)$  is a sparse vector indicating
the presence of words
(usually 0 or 1).

# Question Answering with Subgraph Embeddings [Bordes et al. EMNLP 2014]

Embedding a candidate answer a

$$g(a) = W\varphi(a)$$

$\varphi(a)$ is a sparse vector representation of the answer $a$

- **Single Entity**

The answer is represented as a single entity:

$\varphi(a)$ is a 1-of-Ns coded vector with 1 corresponding the answer.

- **Path Representation**

The answer is represented as a path from the entity mentioned in the question to the answer entity $a$.

$\varphi(a)$ is a 3-of-Ns (or 4-of-Ns) coded vector, expressing the start and the end entities of the path and the relation types (but not entities) in-between.

film → director

Resident_Evil → director → Paul.W.S.Anderson

type ↑ film

type ↑ director

budget → "6.5E7"

Candidate Answer

2-hop paths

# Question Answering with Subgraph Embeddings [Bordes et al. EMNLP 2014]

- Subgraph Representation

The answer is represented both the path and 1-hop neighbors around the answer a.

Embedding a candidate answer a

$$g(a) = W\varphi(a)$$

$\varphi(a)$ is a sparse vector
representation of the answer $a$

film

director

type

type

Candidate
Answer

Resident_Evil

director

Paul.W.S.Anderson

budget

"6.5E7"

1-hop
neighbors

# Question Answering with Subgraph Embeddings [Bordes et al. EMNLP 2014]

Scoring Function

candidate answer

$$S(q,a) = f(q)^T g(a)$$

question sentence

The loss function

$$\sum_{i=1}^{|D|} \sum_{a' \in A'(a_i)} \max\{0, m - S(q_i, a_i) + S(q_i, a')\}$$

$A'(a_i)$ is a set of incorrect canidates to question $q$.

# Question Answering over Freebase with Multi-Column Convolutional Neural Networks [Dong et al., ACL 2015]



Figure 1: Overview for the question-answer pair *(when did Avatar release in UK, 2009-12-17)*. Left: network architecture for question understanding. Right: embedding candidate answers.

# Question Answering over Freebase with Multi-Column Convolutional Neural Networks [Dong et al., ACL 2015]

## Scoring Function

question sentence

candidate answer

$$S(q, a) =$$

$$f_1(q)^T g_1(a) + f_2(q)^T g_2(a) + f_3(q)^T g_3(a)$$

answer path   answer context   answer type

# Question Answering over Freebase with Multi-Column Convolutional Neural Networks [Dong et al., ACL 2015]

Let the question $q = w_1 w_2 ... w_n$

The look layer transform every word into a vector

$$w_j = W_v u(w_j)$$

$$W_v \in \Re^{d_v \times |V|},$$

$d_v$ is the word embedding dimention and

|V| is the vocabulary size

# Question Answering over Freebase with Multi-Column Convolutional Neural Networks [Dong et al., ACL 2015]

MCCNNs for
Question Understanding

Let the question $q = w_1 w_2 ... w_n$

The convolutional layer computes representation of the words in sliding windows.

$$x_j = h(W[w_{j-s}^T ... w_j^T ... w_{j+s}^T] + b)$$

The max-pooling layer

$$f(q) = \max_{j=1,...,n} \{x_j\}$$

# Question Answering over Freebase with Multi-Column Convolutional Neural Networks [Dong et al., ACL 2015]

Embedding Candidate Answers

Answer Path

$$g_1(a) = \frac{1}{\|u_p(a)\|_1} W_p u_p(a)$$

$u_p(a)$ is a length-|R| binary vector, indicating the presence or absence of every relation in the answer path.

$W_p \in \mathfrak{R}^{d_q \times |R|}$ is the parameter matrix

# Question Answering over Freebase with Multi-Column Convolutional Neural Networks [Dong et al., ACL 2015]

Embedding Candidate Answers

Answer Context

The 1-hop entities and relations connected to the answer path are regarded as the *answer context*.

$$g_2(a) = \frac{1}{\|u_c(a)\|_1} W_c u_c(a)$$

$u_c(a)$ is a length-|C| binary vector, indicating the presence or absence of every entity or relation in the context.

$W_c \in \Re^{d_q \times |C|}$ is the parameter matrix

# Question Answering over Freebase with Multi-Column Convolutional Neural Networks [Dong et al., ACL 2015]

## Answer Type

Type information is an important clue to score candidate answers.

$$g_3(a) = \frac{1}{\left\| u_t(a) \right\|_1} W_t u_t(a)$$

$u_t(a)$ is a length-|T| binary vector, indicating the presence or absence of answer type.

$W_t \in \mathfrak{R}^{d_t \times |T|}$ is the parameter matrix

# Question Answering over Freebase with Multi-Column Convolutional Neural Networks [Dong et al., ACL 2015]

## Model Training

For every correct answer a of the question q, we randomly sample k wrong a' from the set of candidate answers, and use them as the negative instances to estimate parameters.
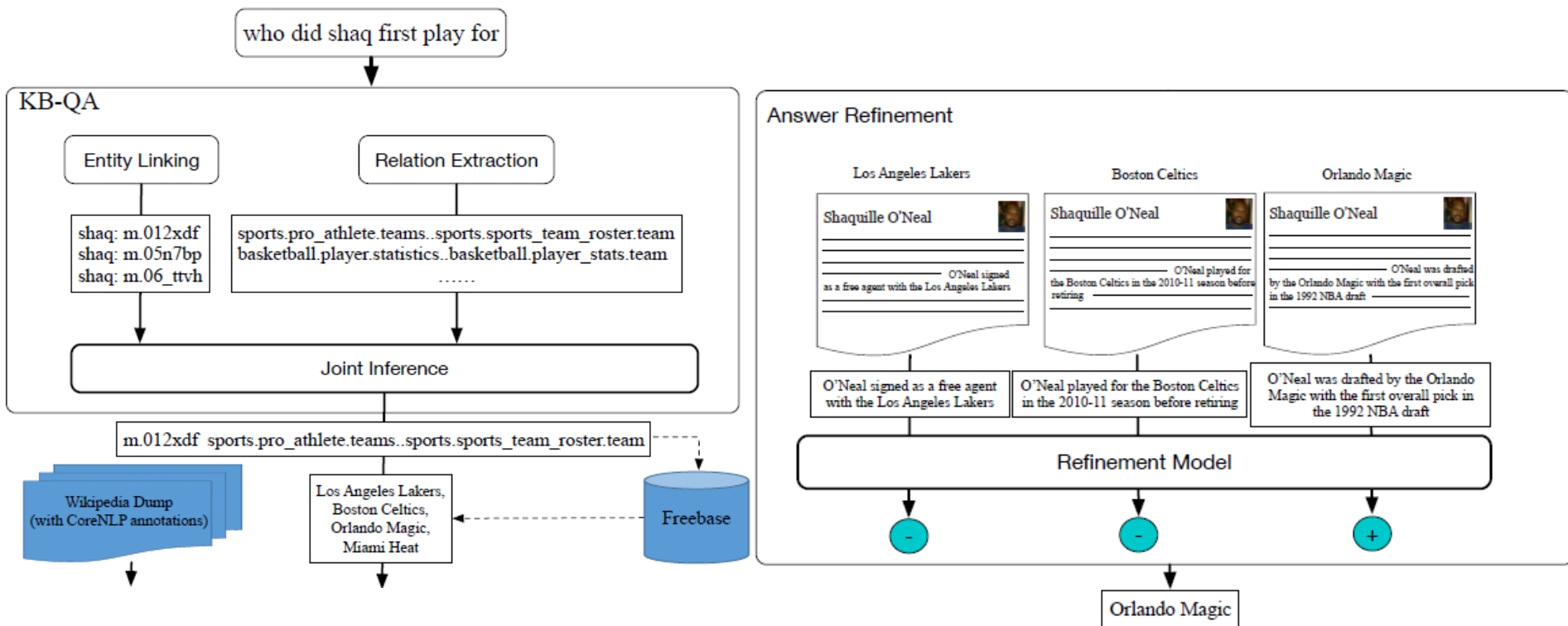
$$l(q,a,a') = (m - S(q,a) + S(q,a'))_+$$

$$\min \sum_q \frac{1}{|A_q|} \sum_{a \in A_q} \sum_{a' \in R_q} l(q,a,a')$$

$$R_q \subseteq C_q \setminus A_q$$
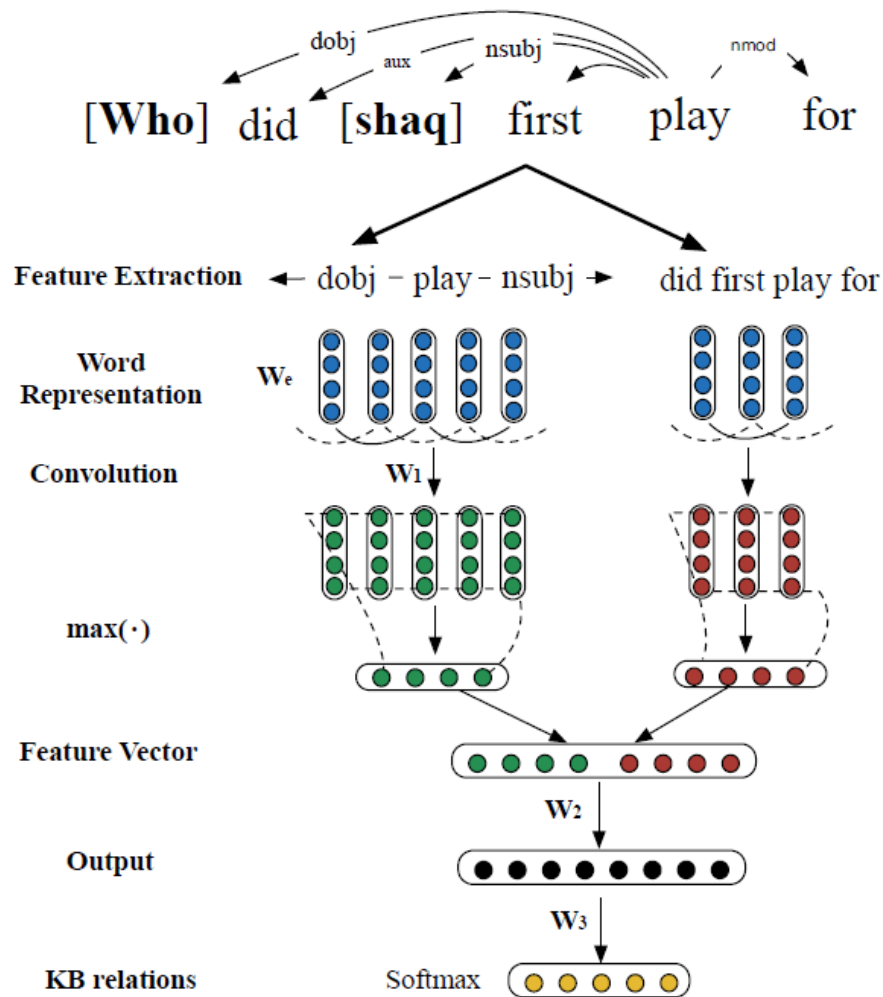
$A_q$ is the correct answer set to question $q$.

$C_q$ is the set of canidate answer set to question $q$.

# Question Answering on Freebase via Relation Extraction and Textual Evidence [Xu et al., ACL 2016]

胡森

# Question Answering on Freebase via Relation Extraction and Textual Evidence[Xu et al., ACL 2016]

Relation Extraction

# Question Answering on Freebase via Relation Extraction and Textual Evidence[Xu et al., ACL 2016]

Question Decomposition

"who plays ken barlow in coronation street? "

decompose
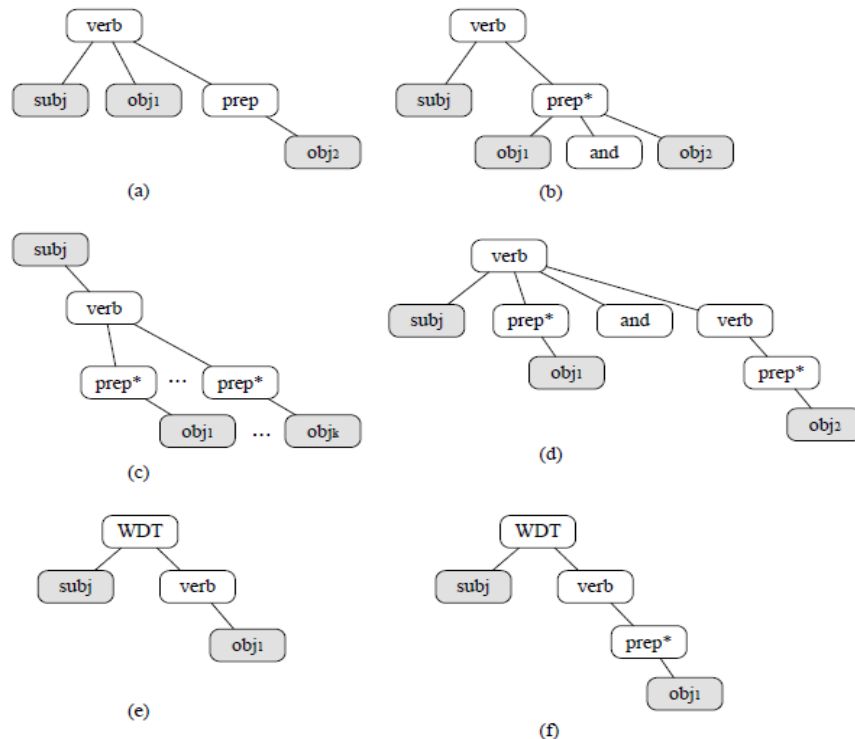
"who plays ken barlow"
+
"who plays in coronation street"



Figure 3: Syntax-based patterns for question decomposition.

# KG-based Question/Answering

- Information Retrieval-based
  - Generate candidate answers
  - Ranking

- Semantic Parsing-based
  - Translate NLQ to logical forms
  - Executing

# Semantic Parsing

[Zettlemoyer et al., UAI 05]

Transforming natural language (NL) sentences into computer executable complete meaning representations (MRs) for domain-specic applications.

E.g., "Which states borders New Mexico ?"

Lambda-**calculus**   [Alonzo Church, 1940 ]

$$\lambda x.state(x) \wedge borders(x, new\_mexico)$$

"**Simply typed Lambda-calculus** can express varies database query languages such as **relational algebra**, fixpoint logic and the complex object algebra." [Hillebrand et al., 1996]

# Semantic Parsing

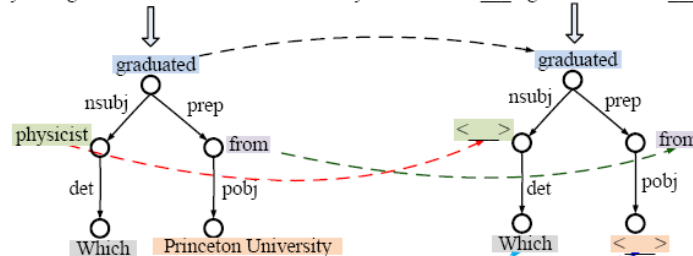- **Manually constructed rules**
  [Pedoe, AI magazine 2010]

- **Grammar-based, e.g.,**
  Combinatory Categorial Grammar
  [Zettlemoyer and Collins, UAI 2005]

- **Supervised Learning**
  [Berant and Liang, ACL 2014]



Template-based Approach [cite: Weiguo Zheng, Lei Zou, et al., SIGMOD 15]

# Semantic Parsing via Staged Query Graph Generation: Question Answering with Knowledge Base [Yih et al., ACL 2015]
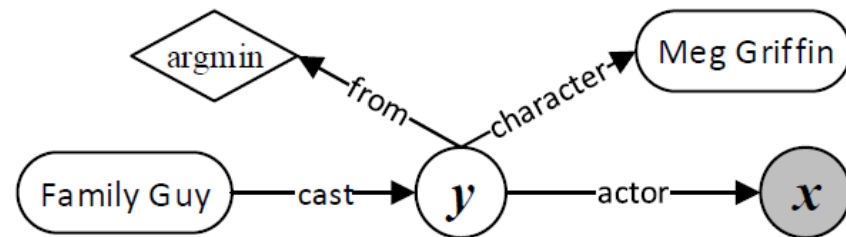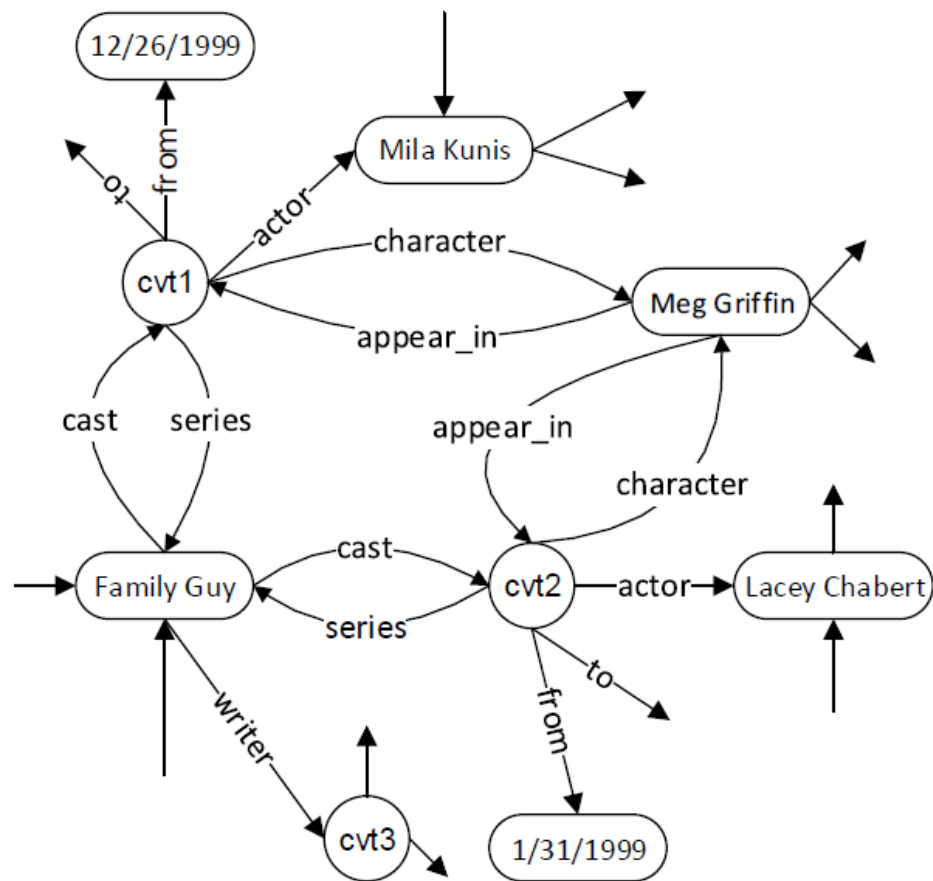


Figure 2: Query graph that represents the question "Who first voiced Meg on Family Guy?"

# Semantic Parsing via Staged Query Graph Generation: Question Answering with Knowledge Base [Yih et al., ACL 2015]
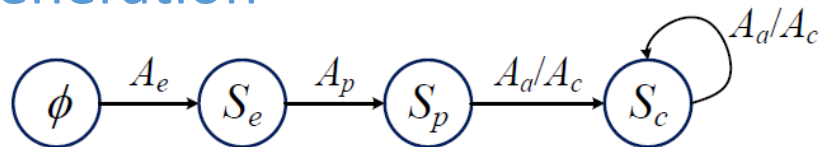
## Query Graph Generation



Figure 3: The legitimate actions to *grow* a query graph. See text for detail.
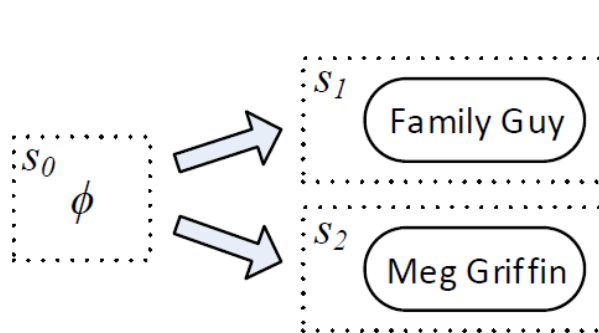


Figure 4: Two possible topic entity linking actions applied to an empty graph, for question "Who first voiced [Meg] on [Family Guy]?"
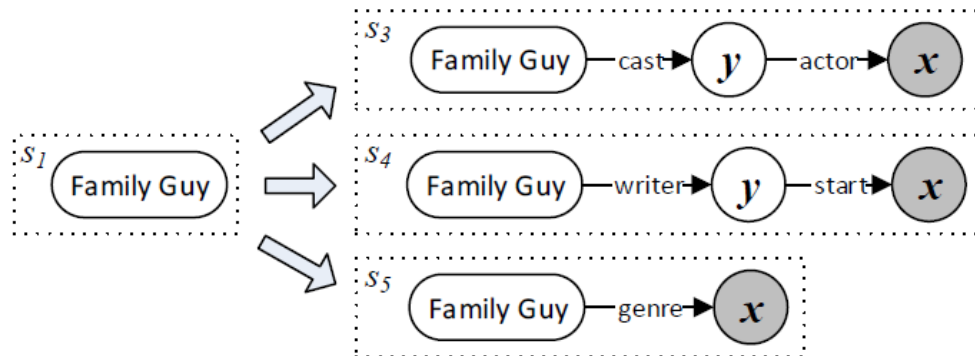


Figure 5: Candidate core inferential chains start from the entity FamilyGuy.

# Semantic Parsing via Staged Query Graph Generation: Question Answering with Knowledge Base [Yih et al., ACL 2015]
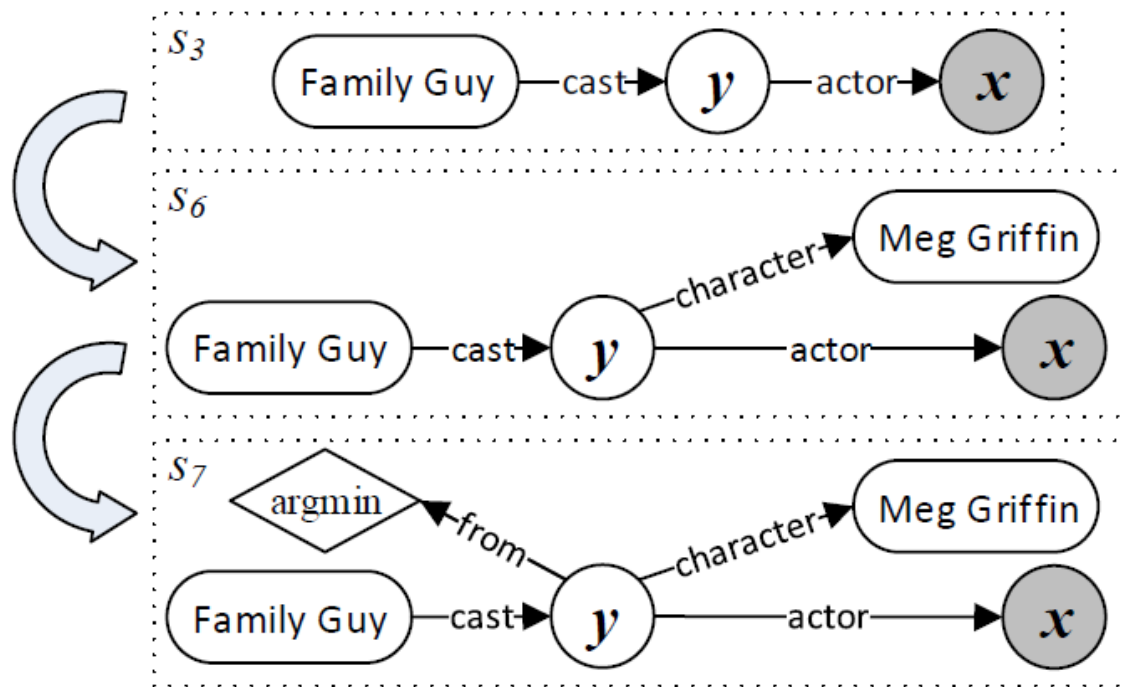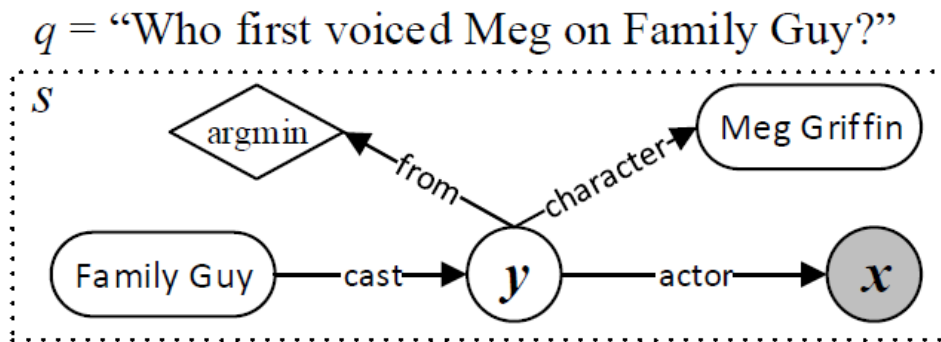
## Query Graph Generation



Figure 7: Extending an inferential chain with constraints and aggregation functions.

# Semantic Parsing via Staged Query Graph Generation: Question Answering with Knowledge Base [Yih et al., ACL 2015]

## Reward Function



$q$ = "Who first voiced Meg on Family Guy?"

(1) EntityLinkingScore(FamilyGuy, "Family Guy") = 0.9
(2) PatChain("who first voiced meg on <e>", cast-actor) = 0.7
(3) QuesEP($q$, "family guy cast-actor") = 0.6
(4) ClueWeb("who first voiced meg on <e>", cast-actor) = 0.2
(5) ConstraintEntityWord("Meg Griffin", $q$) = 0.5
(6) ConstraintEntityInQ("Meg Griffin", $q$) = 1
(7) AggregationKeyword(argmin, $q$) = 1
(8) NumNodes(s) = 5
(9) NumAns(s) = 1

Figure 8: Active features of a query graph $s$. (1) is the entity linking score of the topic entity. (2)-(4) are different model scores of the core chain. (5) indicates 50% of the words in "Meg Griffin" appear in the question $q$. (6) is 1 when the mention "Meg" in $q$ is correctly linked to MegGriffin by the entity linking component. (8) is the number of nodes in $s$. The knowledge base returns only 1 entity when issuing this query, so (9) is 1.

# Semantic Parsing via Staged Query Graph Generation: Question Answering with Knowledge Base [Yih et al., ACL 2015]

**Identifying Core Inferential Chain**
**(Relation Extraction)**

two neural networks
1) question
2) inferential chain

Compute Similarity
(e.g. cosine)

Semantic layer: $y$

Semantic projection matrix: $W_s$
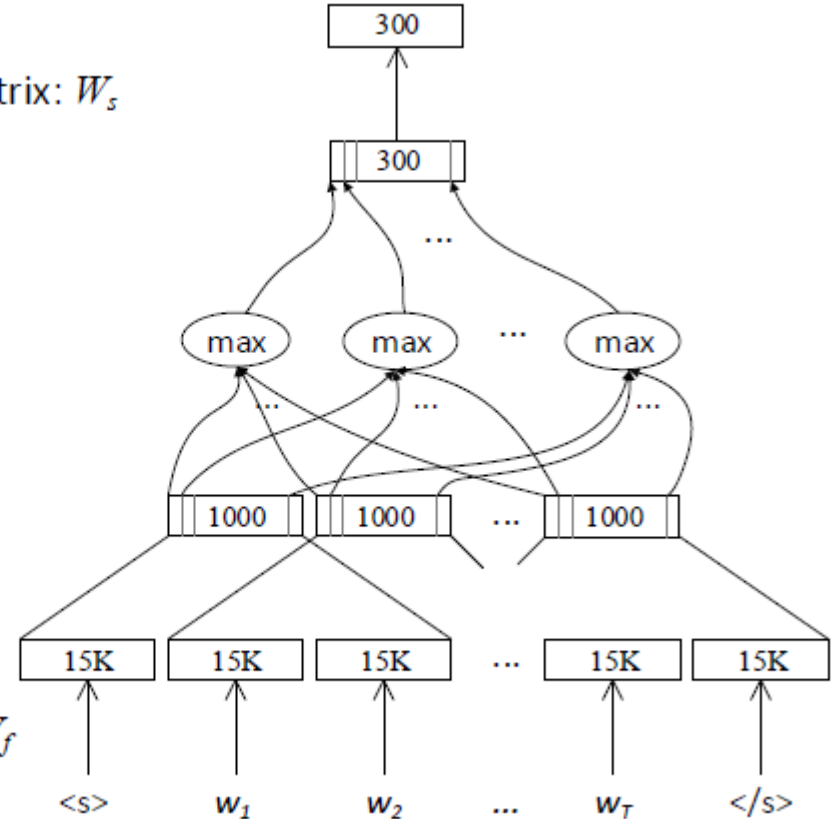
Max pooling layer: $v$

Max pooling operation

Convolutional layer: $h_t$

Convolution matrix: $W_c$

Word hashing layer: $f_t$

Word hashing matrix: $W_f$

Word sequence: $x_t$

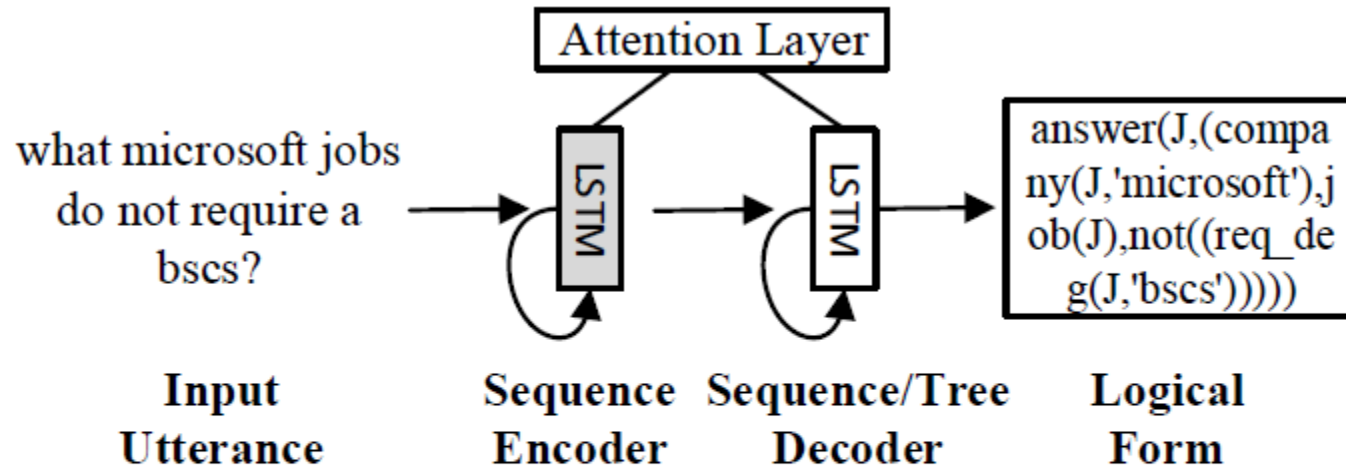# Language to Logical Form with Neural Attention
## [Dong et al., ACL 2016]



Figure 1: Input utterances and their logical forms are encoded and decoded with neural networks. An attention layer is used to learn soft alignments.

# Language to Logical Form with Neural Attention
# [Dong et al., ACL 2016]

*dallas to san francisco leaving after 4 in the afternoon please*
(lambda $0 e (and ($>$(departure_time $0) 1600:ti) (from $0 dallas:ci) (to $0 san_francisco:ci)))

**Algorithm 1** Decoding for SEQ2TREE

**Input:** $q$: Natural language utterance
**Output:** $\hat{a}$: Decoding result

1: ▷ *Push the encoding result to a queue*
2: $Q.init(\{hid : \mathsf{SeqEnc}(q)\})$
3: ▷ *Decode until no more nonterminals*
4: **while** $(c \leftarrow Q.pop()) \neq \varnothing$ **do**
5:     ▷ *Call sequence decoder*
6:     $c.child \leftarrow \mathsf{SeqDec}(c.hid)$
7:     ▷ *Push new nonterminals to queue*
8:     **for** $n \leftarrow$ nonterminal in $c.child$ **do**
9:         $Q.push(\{hid : \mathsf{HidVec}(n)\})$
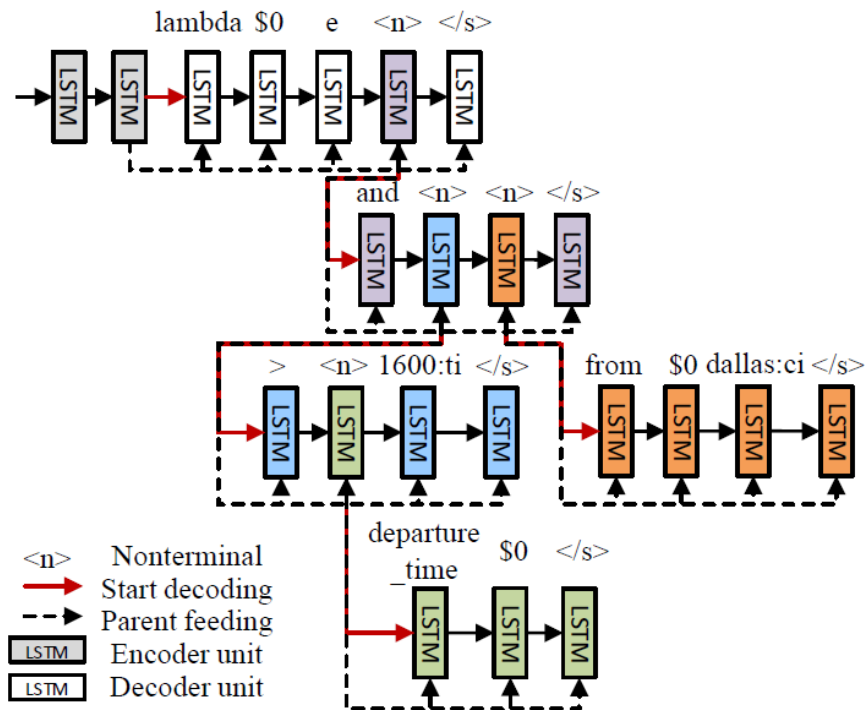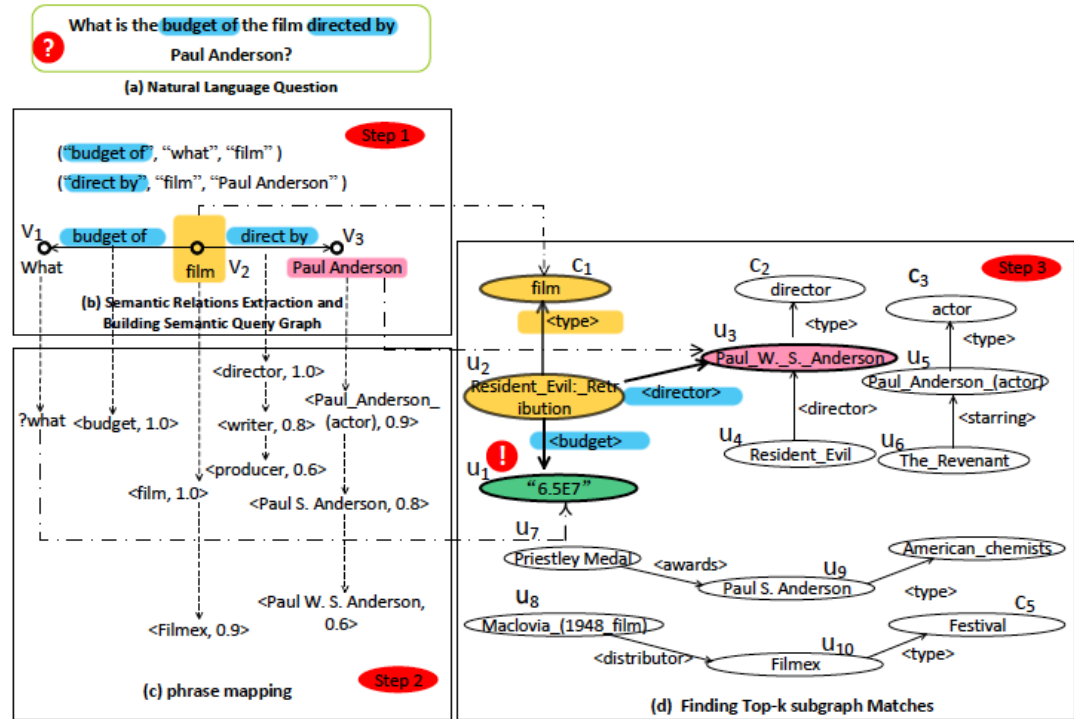10: $\hat{a} \leftarrow$ convert decoding tree to output sequence



Figure 3: Sequence-to-tree (SEQ2TREE) model with a hierarchical tree decoder.

# Our Approach- Data Driven & Relation-first framework gAnswer [Zou et al, SIGMOD 14]
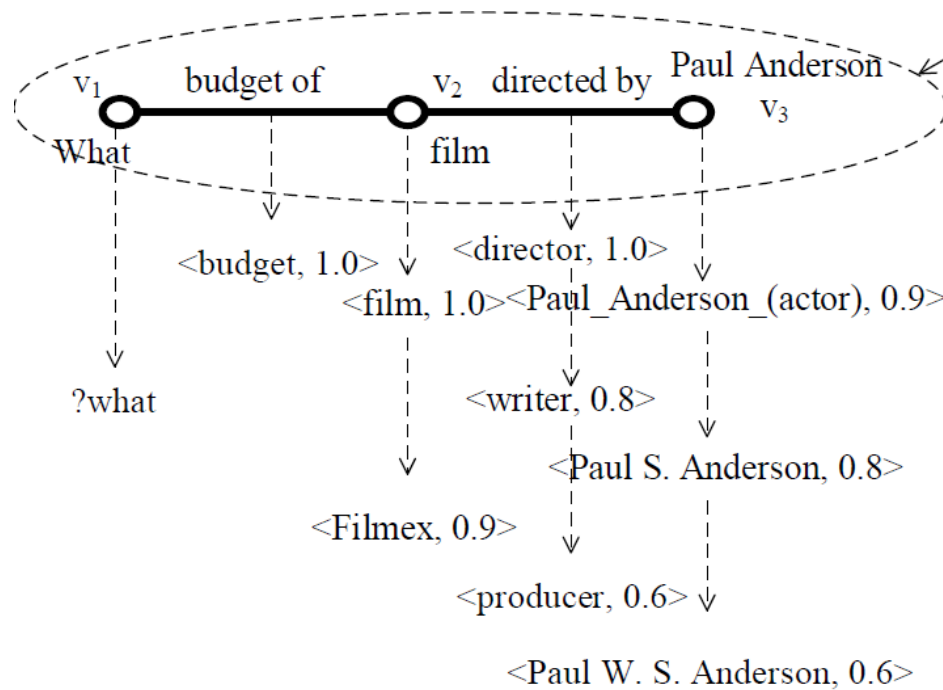
- Using graph matching-based method

- Graph Matching-based Disambiguation

- Combing Disambiguation and Query together

# Our Approach- Data Driven & Relation-first framework gAnswer [Zou et al, SIGMOD 14]

Semantic Query Graph

# Our Approach- Data Driven & Relation-first framework gAnswer [Zou et al, SIGMOD 14]
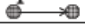
Besides KG, we require two dictionaries.

- **Entity Mention Dictionary**

  It helps the entity linking task
  [Spitkovsky et al., LERC 12; Chisholm et al, TACL 15].

- **Relation Mention Dictionary**

  Mapping the natural language relation
  phrases to predicate in RDF dataset.
  [Nakashole et al., EMNLP-CoNLL 2012]

| Relation Phrases | Predicates or Predicate Paths | Confidence Probability |
|---|---|---|
| "be married to" | \<spouse\> | 1.0 |
| "play in" | \<starring\> | 0.9 |
| "play in" | \<director\> | 0.5 |
| "uncle of" | \<hasChild\> \<hasChild\> \<hasChild\> | 0.8 |
| … … | … … | … … |

# Our Approach- Data Driven & Relation-first framework gAnswer [Zou et al, SIGMOD 14]

- ## Question Understanding
  - Relation extraction



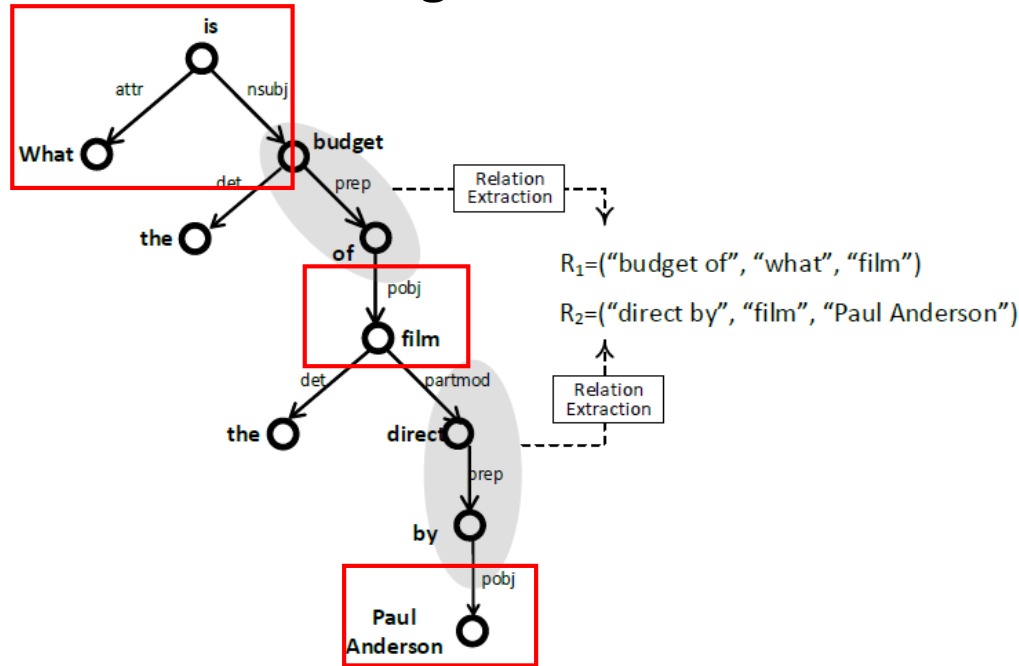| Relation Phrases | Predicates or Predicate Paths | Confidence Probability |
|---|---|---|
| "directed by" | <director> | 1.0 |
| "starred by" | <starring> | 0.9 |
| "budget of" | <budget> | 0.8 |
| "uncle of" | <hasChild> <hasChild> <hasChild> | 0.8 |
| … … | … … | … … |

Relation Paraphrase Dictionary

# Our Approach- Data Driven & Relation-first framework gAnswer [Zou et al, SIGMOD 14]

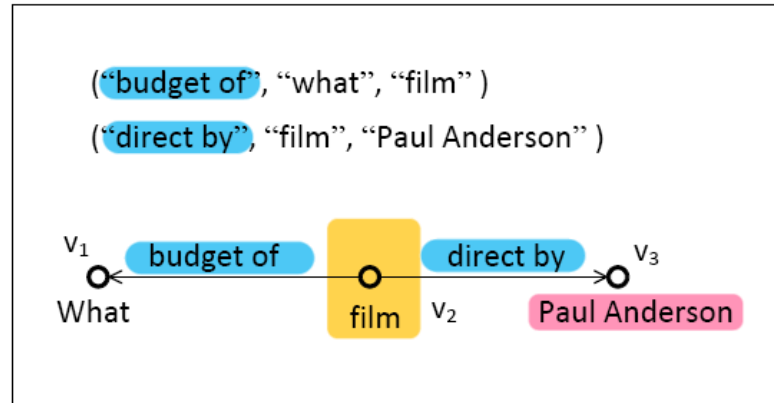- Question Understanding
  - Find associated arguments



$R_1$=("budget of", "what", "film")

$R_2$=("direct by", "film", "Paul Anderson")

# Our Approach- Data Driven & Relation-first framework gAnswer [Zou et al, SIGMOD 14]

- ## Question Understanding
  - Query Graph Assembly



**Semantic Relations Extraction and Building Semantic Query Graph**

# Our Approach- Data Driven & Relation-first framework gAnswer [Zou et al, SIGMOD 14]

- Query Execution

**Algorithm 3** Generating Top-k SPARQL Queries

**Require: Input**: A semantic query graph $Q^S$ and a RDF $G$. **Output**: Top-k SPARQL Queries, i.e., the top-k matches from $Q^S$ to $G$.

1: Sorting all candidates in a non-ascending order
2: Set the threshold $\theta = -\infty$
3: $n = |E(Q^S)| + |V(Q^S)|$
4: Initialize $n$ bit vector $\Gamma$ with zero
5: Initialize maximum heap $H$ with one element $(\Gamma, score(\Gamma))$
6: **while** $(\Gamma, s) \leftarrow H.pop()$ **do**
7: $\quad QG = \text{BuildQueryGraph}(Q^S, \Gamma)$
8: $\quad \text{SubgraphMatching}(G, QG)$ // Any subgraph isomorphism algorithm such as VF2
9: $\quad$ Update the threshold $\theta$ to be the top-k match sore so far.
10: $\quad$ **for** Each candidate list $L_i$ **do**
11: $\quad\quad \Gamma = \Gamma + (1 \leftarrow i)$
12: $\quad\quad H.push(\Gamma, score(\Gamma))$
13: $\quad$ **if** already find k matches **then**
14: $\quad\quad$ Break
15: Output the top-k matches

# Our Approach- Data Driven & Relation-first framework gAnswer [Zou et al, SIGMOD 14]

- Limitations
  - Still highly relied on parser and heuristic rules
  - Can not handle implicit relations

What is the budget of the film directed by Paul Anderson
and starred by a **Chinese girl**

<?girl, dbo:country, dbr:China>

# Our Approach- Data Driven & Node-first framework gAnswer+ [Hu and Zou et al, TKDE 17]
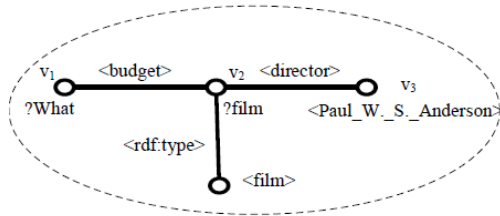
- Data Driven!
  - The structure of query graph can be modified in execution stage.
  - First recognize nodes.

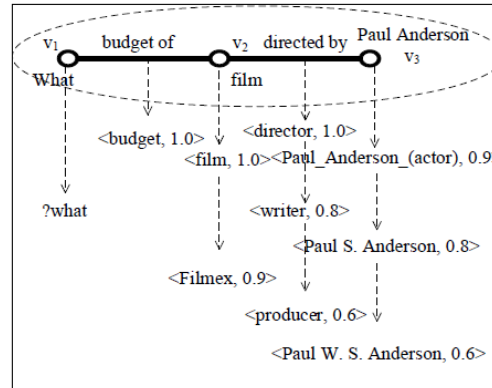# Our Approach- Data Driven & Node-first framework gAnswer+ [Hu and Zou et al, TKDE 17]
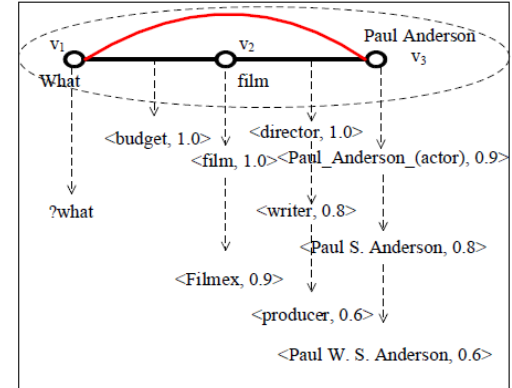
Hyper Query Graph

- Extend SQG by *allowing false edges*.



query graph

semantic query graph

hyper query graph

# Our Approach- Data Driven & Node-first framework gAnswer+ [Hu and Zou et al, TKDE 17]

- Question Understanding
  - Node recognizing

entity extraction + conflict resolution
  - entity, type, literal, wildcard
  - constant, variable
  - modified, hidden information

What is the budget of the film directed by Paul Anderson and starred by a Chinese girl?

| variable | variable | variable | constant | constant | variable |
|----------|----------|----------|----------|----------|----------|
|          |          | type     | entity   | entity   |          |

# Our Approach- Data Driven & Node-first framework gAnswer+ [Hu and Zou et al, TKDE 17]

- Question Understanding
    - Build structure of HQG

connect which two nodes?

**Definition 10.** (**Assumption 1**) Given a question sentence $N$ with appropriate query graph $G$, if $T$ is a correct dependency tree of $N$, the following condition should be satisfied: There is no such three nodes $\{n_1, n_2, n_3\}$ where $n_1$ connect $n_2$ in $G$ and $n_3 \in ShortestPath(n_1, n_2)$ in $T$.

# Our Approach- Data Driven & Node-first framework gAnswer+ [Hu and Zou et al, TKDE 17]

# Our Approach- Data Driven & Node-first framework gAnswer+ [Hu and Zou et al, TKDE 17]

- Question Understanding
  - Finding relations

Explicit relation

# Our Approach- Data Driven & Node-first framework gAnswer+ [Hu and Zou et al, TKDE 17]

- ## Question Understanding
  - Finding relations

## Implicit relation

- Locating the two nodes in KG and finding the frequent predicate between them.

<starring>

starred

by

girl

Chinese

<China>

# Our Approach- Data Driven & Node-first framework gAnswer+ [Hu and Zou et al, TKDE 17]

- ## Query Executing
  - A top-down algorithm

- ## Naïve method
  - (1) Enumerate spanning subgraph of HQG,
  - (2) Call algorithm SQG executing algorithm
  - (3) Sort and select top-k matches

- ## Advanced method
  - (1) Add <drop, 0> to the candidate list of unsteady edges
  - (2) Call algorithm 3

# Our Approach- Data Driven & Node-first framework gAnswer+[Hu and Zou et al, TKDE 17]

- ## Query Executing
  - A top-down algorithm

  ## Drawbacks
  - Query graphs with higher scores may have no matches

| s | p | o |
|---|---|---|
| $e_1$ | $p_1$ | var |
| ... | ... | |
| $e_n$ | $p_m$ | |

# Our Approach- Data Driven & Node-first framework gAnswer+ [Hu and Zou et al, TKDE 17]

- Query Executing
  - A bottom-up algorithm

  Intuition
  - Growing structures step by step
  - Keep correct structures when growing
  - Find matches of multi-label query graph (SQG)
  - Drop useless candidates as early as possible

# Our Approach- Data Driven & Node-first framework gAnswer+ [Hu and Zou et al, TKDE 17]

- ## Query Executing
  - A bottom-up algorithm

```
1: Initialize result set MS, query graph QG, queue que
2: QG ← start node st
3: que.push(st)
4: while x = que.pop() do
5:     /*Try to expand current query graph*/
6:     for each v_i x ∈ E(Q^H) ∧ v_i x ∉ QG do
7:         TQG = QG ← v_i x
8:         if GraphExplore(G, TQG) == TRUE then
9:             QG = TQG
10:        else
11:            QG = Backtrack(QG, v_i x)
12:        if v_i x ∈ QG then
13:            que ← v_i
14: Sort the graph explore results of QG and select top-k matches
```

# Our Approach- Data Driven & Node-first framework gAnswer+ [Hu and Zou et al, TKDE 17]

- Query Executing
  - A bottom-up algorithm

  Optimization
  - Call GraphExplore() only when adding unsteady edges
  - Design cost model to estimate the best explore order

# Experiments

QALD is a series of evaluation campaigns on question answering over linked data.

TABLE 7
Evaluating QALD-6 Testing Questions (Total Question Number=100)

| | Processed | Right | Recall | Precision | F-1 |
|---|---|---|---|---|---|
| **NFF** | 100 | 68 | **0.70** | **0.89** | **0.78** |
| RFF | 100 | 40 | 0.43 | 0.77 | 0.55 |
| CANaLI | 100 | 83 | 0.89 | 0.89 | 0.89 |
| UTQA | 100 | 63 | 0.69 | 0.82 | 0.75 |
| KWGAnswer | 100 | 52 | 0.59 | 0.85 | 0.70 |
| SemGraphQA | 100 | 20 | 0.25 | 0.70 | 0.37 |
| UIQA1 | 44 | 21 | 0.63 | 0.54 | 0.25 |
| UIQA2 | 36 | 14 | 0.53 | 0.43 | 0.17 |
| DEANNA | 100 | 20 | 0.21 | 0.74 | 0.33 |
| Aqqu | 100 | 36 | 0.37 | 0.39 | 0.38 |

QALD-6 Competition Results

# Experiments

WebQuestions is widely used in Question Answering literatures and does not contain golden SPARQL queries.

TABLE 8
Evaluating WebQuestions Testing Questions

|  | Average F1 |
|---|---|
| **NFF** | 49.6% |
| RFF | 31.2% |
| Sempre | 35.7% |
| ParaSempre | 39.9% |
| Aqqu | 49.4% |
| STAGG | 52.5% |
| Yavuz et al. (2016) | **52.6%** |

WebQuestions Results

# Online Demo

URL: http://ganswer.gstore-pku.com/

# Is it Possible ?

Semantic Parsing（NLP）+Query Evaluation（DB）



Where is the nearest post office ?

$$\arg\min(\lambda x.POST(x) \wedge dis(HERE, x))$$

**SPARQL**

SELECT ?x WHERE {
?x rdf:type Post.
?x :longitude ?o.
?x :latitude ?a. }
ORDERY BY Dist(HERE, [?o, ?a])
LIMIT 1

# 与深圳狗尾草公司合作.

刘德华的女儿是？

柬埔寨首都在哪儿？

SPARQL

gStre

An open-source Graph
RDF database

公子小白

6600 万 Triples
Zhishi.me

参考文献：

- Antoine Bordes, Nicolas Usunier, Alberto García-Durán, Jason Weston, Oksana Yakhnenko: Translating Embeddings for Modeling Multi-relational Data. NIPS 2013: 2787-2795

- Luke S. Zettlemoyer, Michael Collins: Learning to Map Sentences to Logical Form: Structured Classification with Probabilistic Categorial Grammars. UAI 2005: 658-666

- Pablo N. Mendes, Max Jakob, Christian Bizer: DBpedia: A Multilingual Cross-domain Knowledge Base. LREC 2012: 1813-1817

- Fabian M. Suchanek, Gjergji Kasneci and Gerhard Weikum, Yago - A Core of Semantic Knowledge,  16th international World Wide Web conference (WWW 2007)

- Kurt D. Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, Jamie Taylor: Freebase: a collaboratively created graph database for structuring human knowledge. SIGMOD Conference 2008: 1247-1250

- Peter Buneman, Gao Cong, Wenfei Fan, Anastasios Kementsietsidis: Using Partial Evaluation in Distributed Query Evaluation. VLDB 2006: 211-222

- Yuk Wah Wong, Raymond J. Mooney: Learning for Semantic Parsing with Statistical Machine Translation. HLT-NAACL 2006

- C. Unger, L. Bühmann, J. Lehmann, A.-C. N. Ngomo, D. Gerber, and P. Cimiano. Template-based question answering over RDF data. In WWW, pages 639–648, 2012

- Lei Zou, Ruizhe Huang, Haixun Wang, Jeffrey Xu Yu, Wenqiang He and Dongyan Zhao, Natural Language Question Answering over RDF ---- A Graph Data Driven Approach , SIGMOD (2014)

- Lei Zou, Jinghui Mo, Lei Chen,M. Tamer Özsu, Dongyan Zhao, gStore: Answering SPARQL Queries Via Subgraph Matching, in Proceedings of 37th International Conference on Very Large Databases (VLDB), 2011.

- Peng Peng, Lei Zou, Tamer Ozsu, Lei Chen, Dongyan Zhao, Processing SPARQL queries over distributed RDF graphs, accepted by VLDB Journal

- Church, A. "A Formulation of the Simple Theory of Types". Journal of Symbolic Logic 5: 1940. doi:10.2307/2266170

参考文献：

- C. Unger, L. Bühmann, J. Lehmann, A.-C. N. Ngomo, D. Gerber, and P. Cimiano. Template-based question answering over RDF data. In WWW, pages 639–648, 2012

- Lei Zou, Ruizhe Huang, Haixun Wang, Jeffrey Xu Yu, Wenqiang He and Dongyan Zhao, Natural Language Question Answering over RDF ---- A Graph Data Driven Approach , SIGMOD (2014)

- Lei Zou, Jinghui Mo, Lei Chen,M. Tamer Özsu, Dongyan Zhao, gStore: Answering SPARQL Queries Via Subgraph Matching, in Proceedings of 37th International Conference on Very Large Databases (VLDB), 2011.

- Antoine Bordes, Sumit Chopra, Jason Weston: Question Answering with Subgraph Embeddings. EMNLP 2014: 615-620

- William Tunstall-Pedoe: True Knowledge: Open-Domain Question Answering using Structured Knowledge and Inference. AI Magazine 31(3): 80-92 (2010)

- Luke S. Zettlemoyer, Michael Collins: Learning Context-Dependent Mappings from Sentences to Logical Form. ACL/IJCNLP 2009: 976-984

- Jonathan Berant and Percy Liang. 2014. Semantic parsing via paraphrasing. In Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (ACL'14), Baltimore, USA

- Kun Xu, Siva Reddy, Yansong Feng, Songfang Huang, Dongyan Zhao: Question Answering on Freebase via Relation Extraction and Textual Evidence. ACL (1) 2016

- Sen Hu, Lei Zou, Jeffrey Xu Yu, Haixun Wang, Dongyan Zhao, Answering Natural Language Questions by Subgraph Matching over Knowledge Graphs, accepted by IEEE Transactions on Knowledge and Data Engineering, 2017

- Wen-tau Yih, Ming-Wei Chang, Xiaodong He, Jianfeng Gao: Semantic Parsing via Staged Query Graph Generation: Question Answering with Knowledge Base. ACL (1) 2015: 1321-1331

- Li Dong, Mirella Lapata: Language to Logical Form with Neural Attention. ACL (1) 2016

# Thanks！

zoulei@pku.edu.cn