

Question Answering Systems


Neural KG-QA systems

Rishiraj Saha Roy

Max Planck Institute for Informatics, Germany



Quick word(s) on reviewing

- 1 ■ Do not pick disadvantages from future work
- 2 ■ Judge time and scope of paper *course* *chronological X*
 - No complex, no feedback, no paraphrases are not necessarily valid points
- 3 ■ Take a stand: do not pose same thing as advantage or disadvantage
- 4 ■ *A1) X*
A2
A3
X (D1 D2 D3) Write mutually exclusive points based on understanding of paper *"semantic"* 
- 5 ■ Read papers beforehand, attend class, clarify questions during class
- 6 ■ Assignment all 4's does not directly mean a grade of 4.0; don't worry 😊

Question of the day

How do we design KG-QA systems with
neural learning?

You'll find this covered in

- ① ■ Semantic Parsing via Staged Query Graph Generation: Question Answering with Knowledge Base
 - Yih et al.
 - ACL 2015
 - <https://www.aclweb.org/anthology/P15-1128.pdf>
- ② ■ Knowledge Graph Embedding Based Question Answering
 - Huang et al.
 - WSDM 2019
 - <http://research.baidu.com/Public/uploads/5c1c9a58317b3.pdf>

msR

Baidu Research

Research paper 1

Semantic Parsing via Staged Query Graph Generation: Question Answering with Knowledge Base

STAGS

**Semantic parsing via staged query graph
generation: Question answering with knowledge base**

[PDF] microsoft.com

SW Yih, MW Chang, X He, J Gao - 2015 - microsoft.com

We propose a novel **semantic parsing** framework for question answering using a knowledge base. We define a **query graph** that resembles subgraphs of the knowledge base and can be directly mapped to a logical form. **Semantic parsing** is reduced to **query graph** ...

☆ 77 Cited by 353 Related articles All 8 versions >>

Best performance on WebQuestions*

Method	Prec.	Rec.	F ₁
(Berant et al., 2013)	48.0	41.3	35.7
(Bordes et al., 2014b)	-	-	29.7
(Yao and Van Durme, 2014)	-	-	33.0
(Berant and Liang, 2014)	40.5	46.6	39.9
(Bao et al., 2014)	-	-	37.5
(Bordes et al., 2014a)	-	-	39.2
(Yang et al., 2014)	-	-	41.3
(Wang et al., 2014)	-	-	45.3
Our approach – STAGG	52.8	60.7	52.5

In Paper
Leaderboard

w/o addl.
evidence

~ Jul-Aug 2015

→ only KG

Satish Jain

NAACL

* Until <https://www.aclweb.org/anthology/N16-2016.pdf>

The STAGG System: What's new?

- First notable system to use graph representations of KBs for QA *KG*
- Previously restricted to RDF triples and SPARQL querying *(logical form)*
- Equivalence of subgraph search and logical form
- Introduces *neural* learning

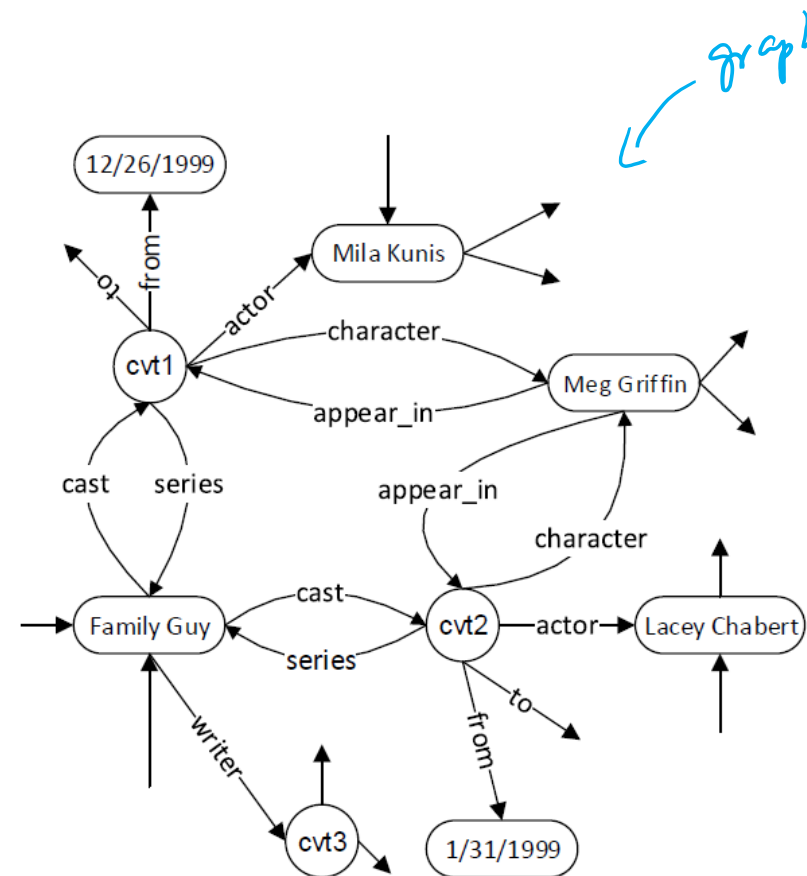


Figure 1: Freebase subgraph of Family Guy

The STAGG System: What's new?

- Defines notion of context subgraph
- No need to search over whole KG
- Subgraph search posed as staged procedure of growing query graph

- Find topic entity
- Find predicate
- Consider additional constraints

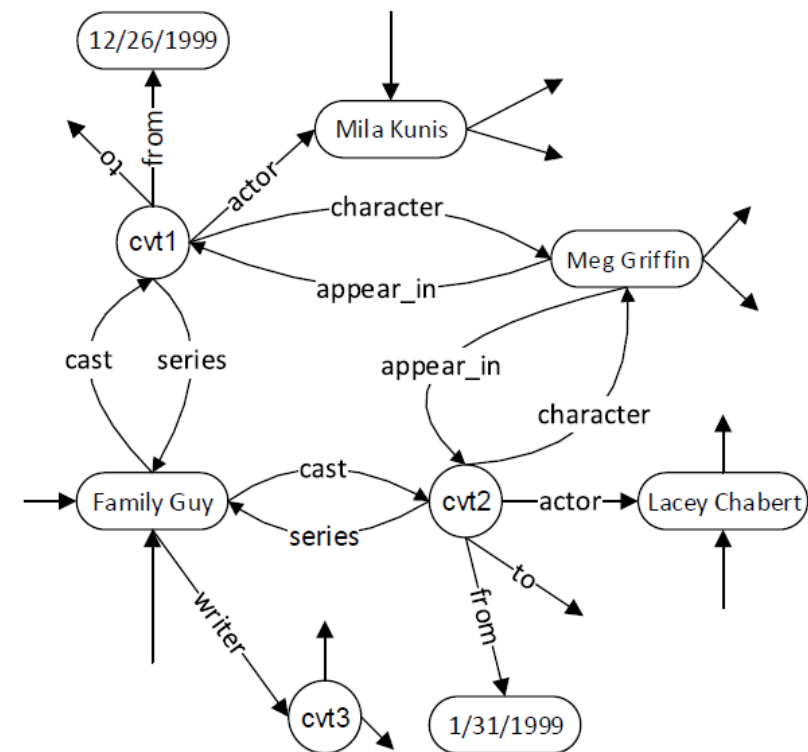
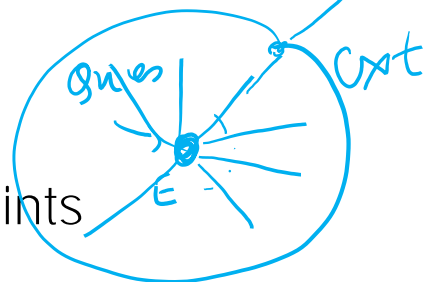


Figure 1: Freebase subgraph of Family Guy

Why focus on a context graph?

eg →

Who first voiced Meg on Family Guy?

re don't need to!!

① Focus on promising area of KG space

② Resolving other entities and predicates also becomes easier

③ Two sub-problems: Entity linking, BA

2 Relation matching

④ What motivates thinking on these lines?

→ go neural!
→ classification/ranking/sim (prediction)
→ not trivial for BA

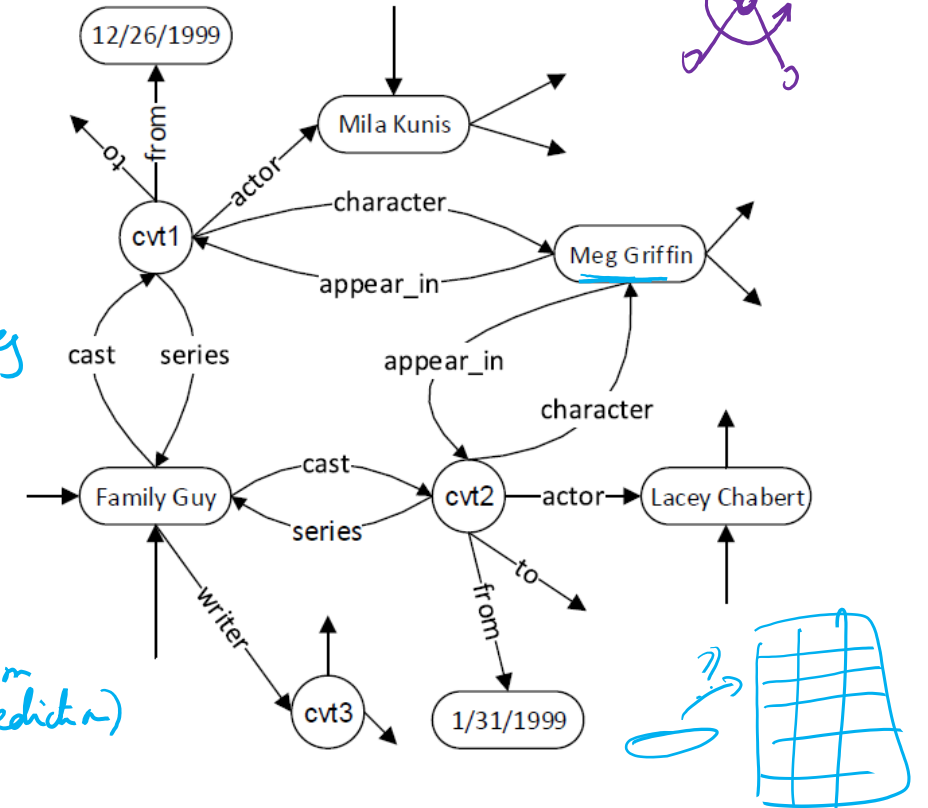


Figure 1: Freebase subgraph of Family Guy

Graph model

- 1 ■ Entities, classes, literals, CVTs as nodes
- 2 ■ Predicates as edge labels
- Directed edges between nodes
- No need to differentiate between entities and literals (and CVTs)
- STAGG works over Freebase
 - 46M topics, 2.6B facts *(xt)* *efficiency*

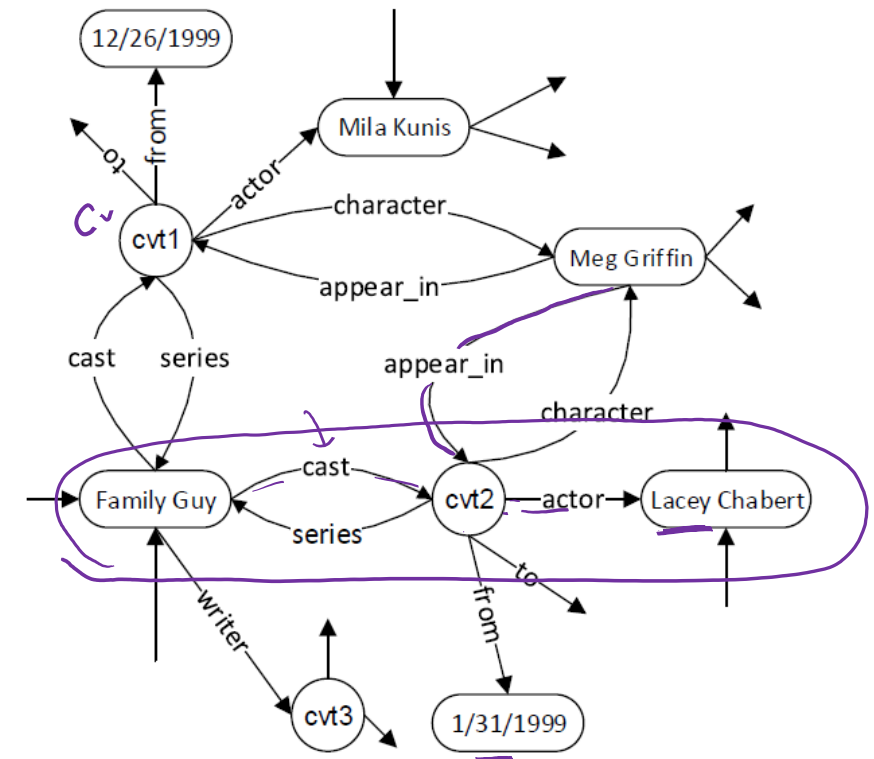


Figure 1: Freebase subgraph of Family Guy

cxt
+

Query subgraph

- 1 ■ Grounded entity ^{linked}
 - 2 ■ Existential variable
 - 3 ■ Lambda variable
 - 4 ■ Aggregation function ^{+ constraints}
- Answer maps to lambda variables (Berant et al. 2013)

∴ graph → flexible search

can be executed over KGs

graphical query

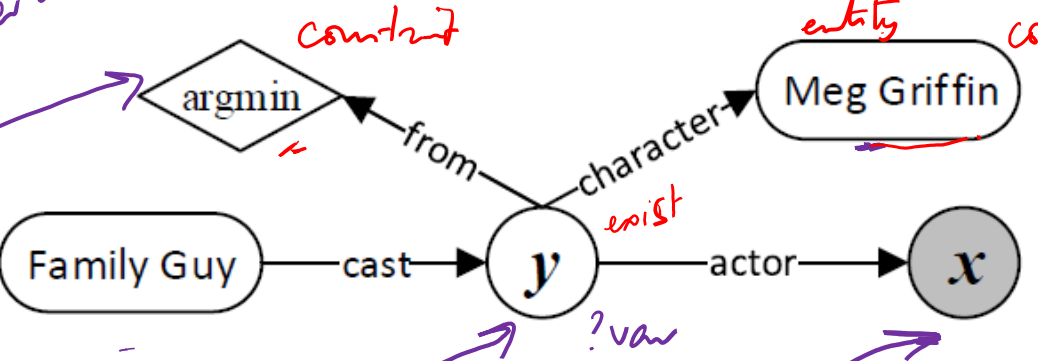


Figure 2: Query graph that represents the question “Who first voiced Meg on Family Guy?”

count
λ-calculus → λ-DCS
simplify

graphical query ≡ logical form ≡ SPARQL query
 $\lambda x. \exists y. \text{cast}(\text{FamilyGuy}, y) \wedge \text{actor}(y, x) \wedge \text{char}(y, \text{MegGr})$

SELECT ?x WHERE
{ FG cast ?y. ?y actor ?x.
?y char MG. ?y from ?z } ORDER BY ?z LIMIT 1

“first” → ??

STAG 6

Staged query graph generation

- Build a tree graph
- Root = topic entity
- One lambda variable = answer
- One directed path from root to answer: Core inferential chain
- Possibly multiple ^{derivation} existential variables in between

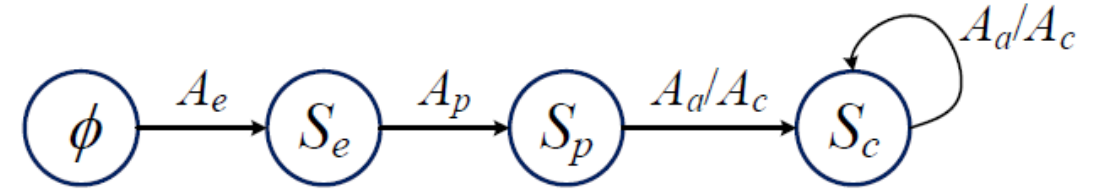


Figure 3: The legitimate actions to *grow* a query graph. See text for detail.

Staged query graph generation

- nodes
 All are variables on core chain except root (grounded) inf...
- Additional constraints: Entity or aggregation nodes can be attached to each variable meg
- Grow graph with actions
 - Pick entity, pick predicate (core chain), pick condition

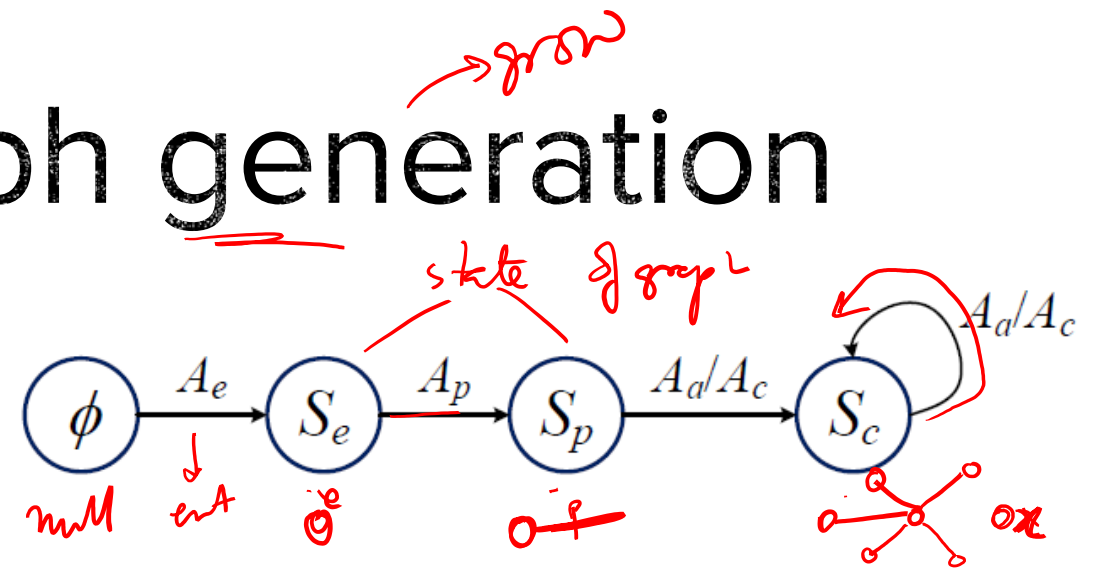
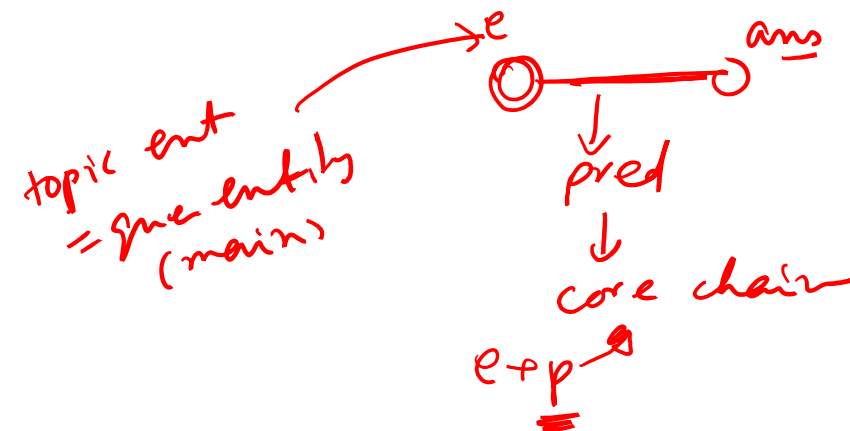


Figure 3: The legitimate actions to *grow* a query graph. See text for detail.



Why this growth order?

- Entity, predicate, conditions
- Efficiency considerations
- Implicit in graph-based QA: To the extent of being “natural”
- What happens if you choose predicates or classes first?

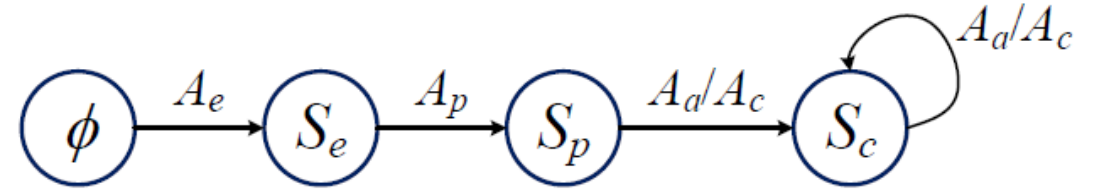
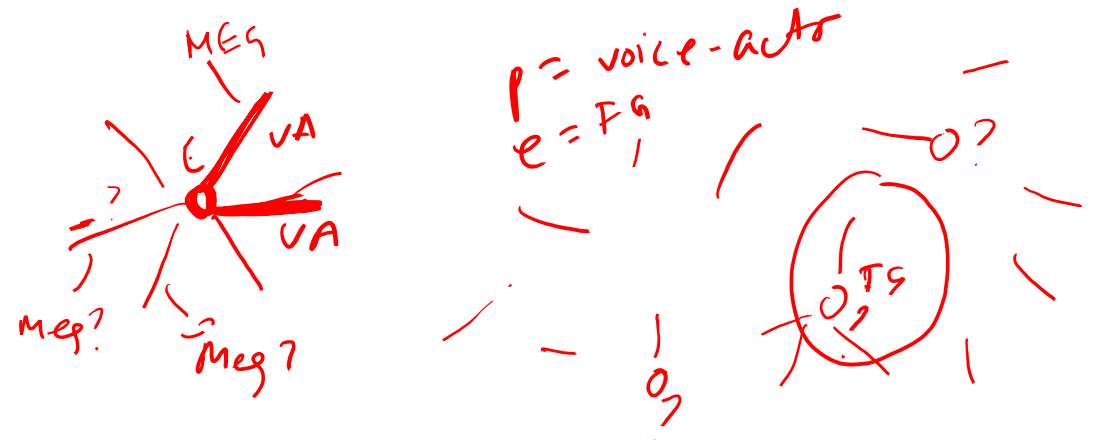


Figure 3: The legitimate actions to *grow* a query graph. See text for detail.



Linking topic entity

NERD

- S-MART: NED system for short and noisy texts (Tweets)
- Standard lexicon-based scoring with Wikipedia metadata
- To account for NED mistakes, retains up to top-10 NED!

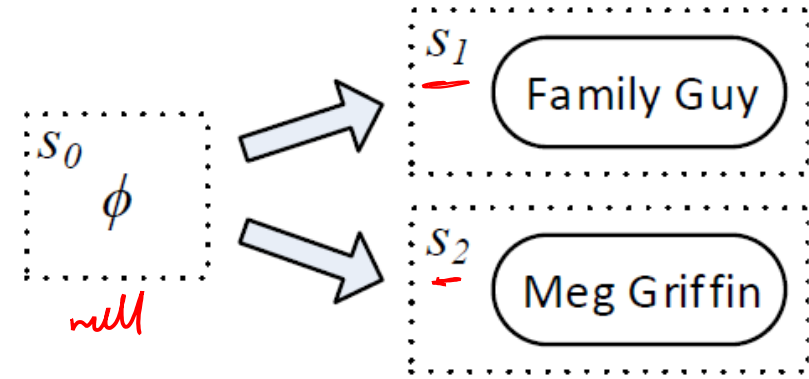


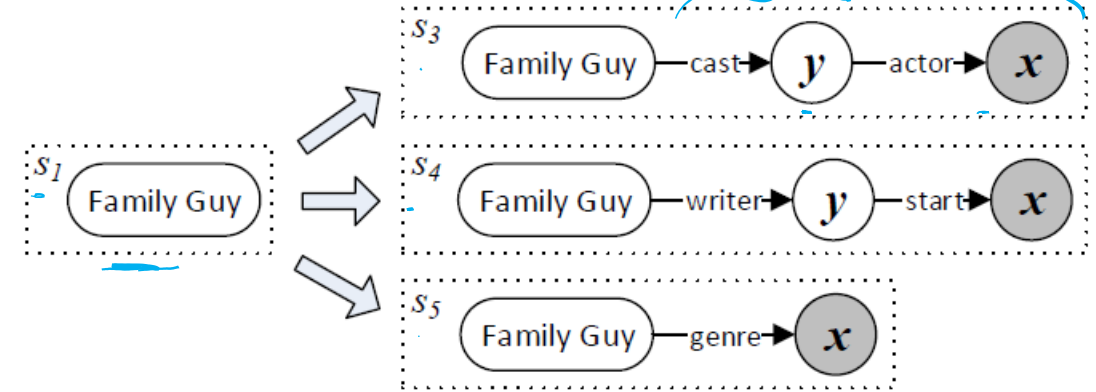
Figure 4: Two possible topic entity linking actions applied to an empty graph, for question “Who first voiced [Meg] on [Family Guy]?”

Identifying core inferential chain

→ predicate candidates

paths / ind. chains

- cx
- Only need to explore legitimate predicate sequences (paths) that start from linked entity!



→ qualified

- All paths of length 2 if CVT in between, length 1 otherwise

dist=1?
In CVT MG

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

Identifying core inferential chain

- Problem boils down to matching question with path similarity
- Similarity using predictors based on neural networks

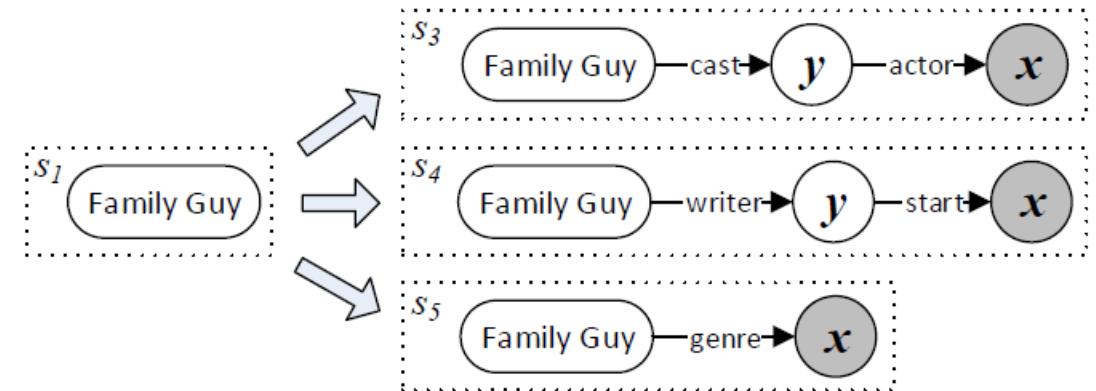


Figure 5: Candidate core inferential chains start from the entity FamilyGuy.

Deep convolutional neural networks

- 7 layer
 - Map question to pattern by replacing entity with <e> *wildcard tag*
 - Two neural networks
 - Question pattern
 - Inferential chain (path)
 - Siamese network architecture
 - Inspired by mismatch in question and KG vocabulary, lots of paraphrases
- 2 networks same pattern*

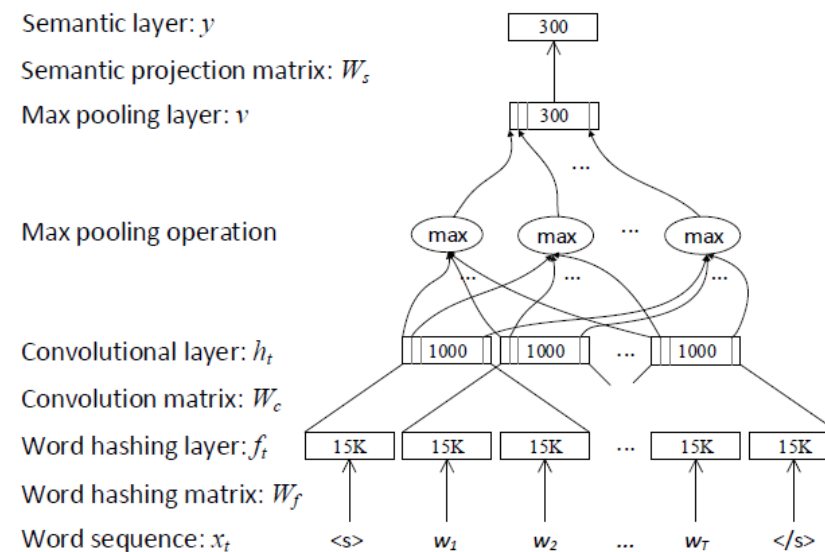


Figure 6: The architecture of the convolutional neural networks (CNN) used in this work. The CNN model maps a variable-length word sequence (e.g., a pattern or predicate sequence) to a low-dimensional vector in a latent semantic space. See text for the description of each layer.

Deep convolutional neural networks

Master's

- Word hashing technique (Huang et al. 2013)

- Break word into character trigrams

- Who \rightarrow #-w-h, w-h-o, h-o-#

- # Why?

- Reduces high-dimensionality of input space (think! #words more or #trigrams more?)

- Subword semantics without tricky + slow techniques like stemming, lemmatization, ...

- Robust to typos

brings
bringing

gentles
now: word space
now: character trigram

bringing
bringly $\rightarrow x!$

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

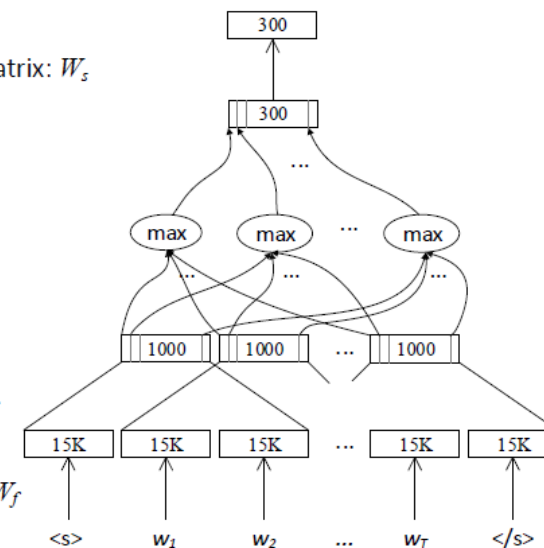


Figure 6: The architecture of the convolutional neural networks (CNN) used in this work. The CNN model maps a variable-length word sequence (e.g., a pattern or predicate sequence) to a low-dimensional vector in a latent semantic space. See text for the description of each layer.

Deep convolutional neural networks

- Uses convolution layer to project the letter-trigram vectors of words within a context window of 3 words
- Creates a local contextual feature vector $h_t \rightarrow h_t$
- Local to global context vector using max pooling operation
- Final fully connected feedforward layer to create final “semantic” vector

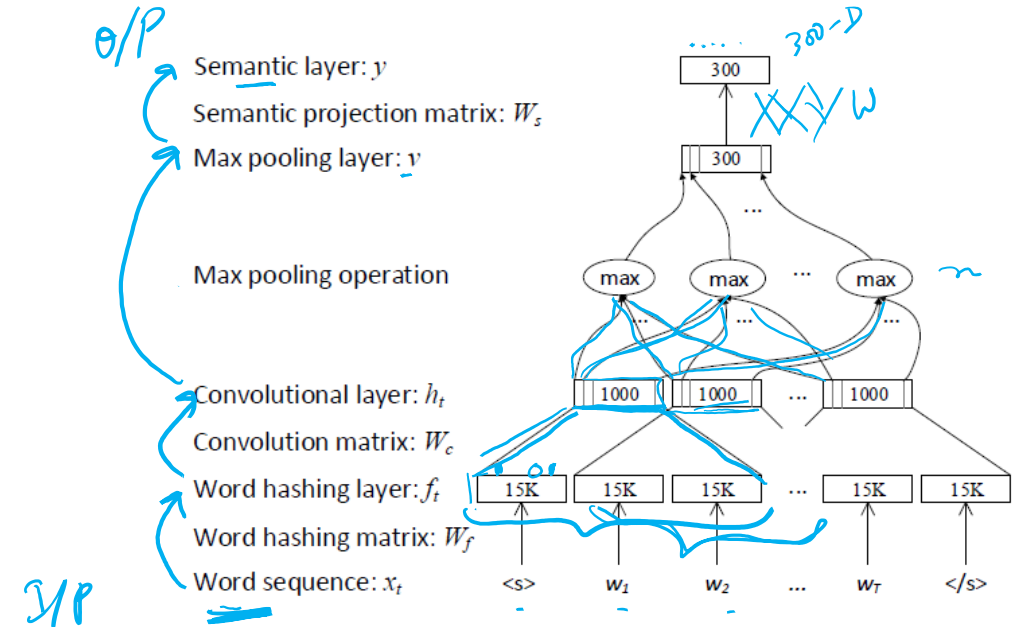
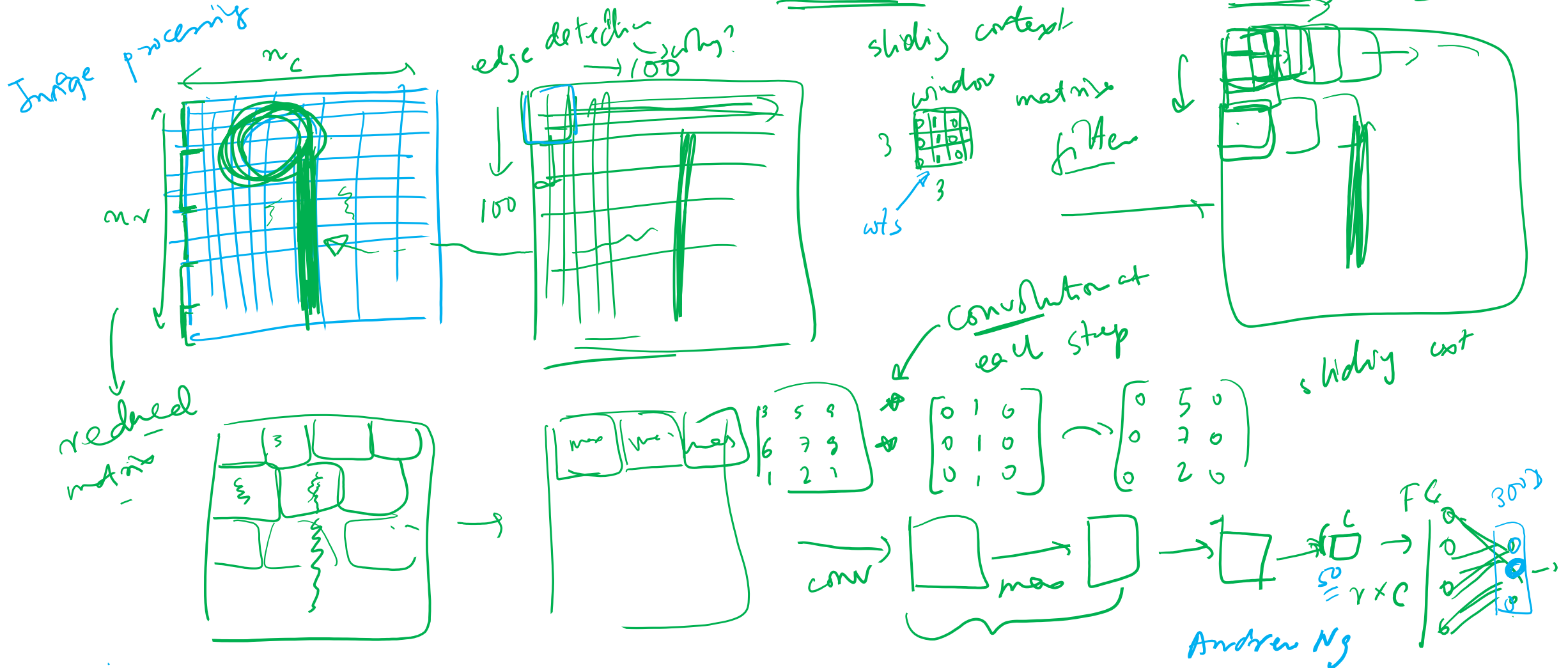


Figure 6: The architecture of the convolutional neural networks (CNN) used in this work. The CNN model maps a variable-length word sequence (e.g., a pattern or predicate sequence) to a low-dimensional vector in a latent semantic space. See text for the description of each layer.

Convolution and max pooling



Deep convolutional neural networks

- Training model needs positive (and negative) pairs: (question, core-chain)
- Sample +ve pair
(who first voiced meg on <e>, cast-actor)
- Obtained from SPARQL queries (semantic parses = correct subgraph patterns)
- If not available, ~~create Q-a pairs by using~~ Q-A paths in KG: Distant supervision!

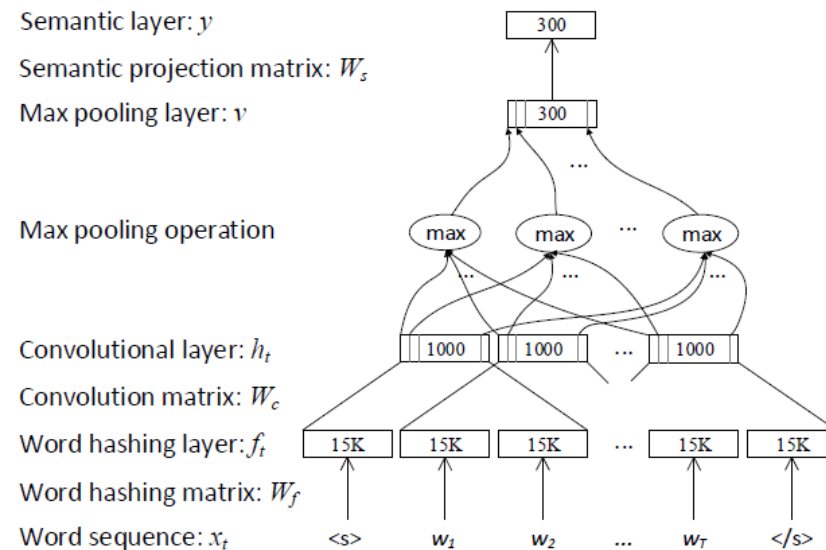


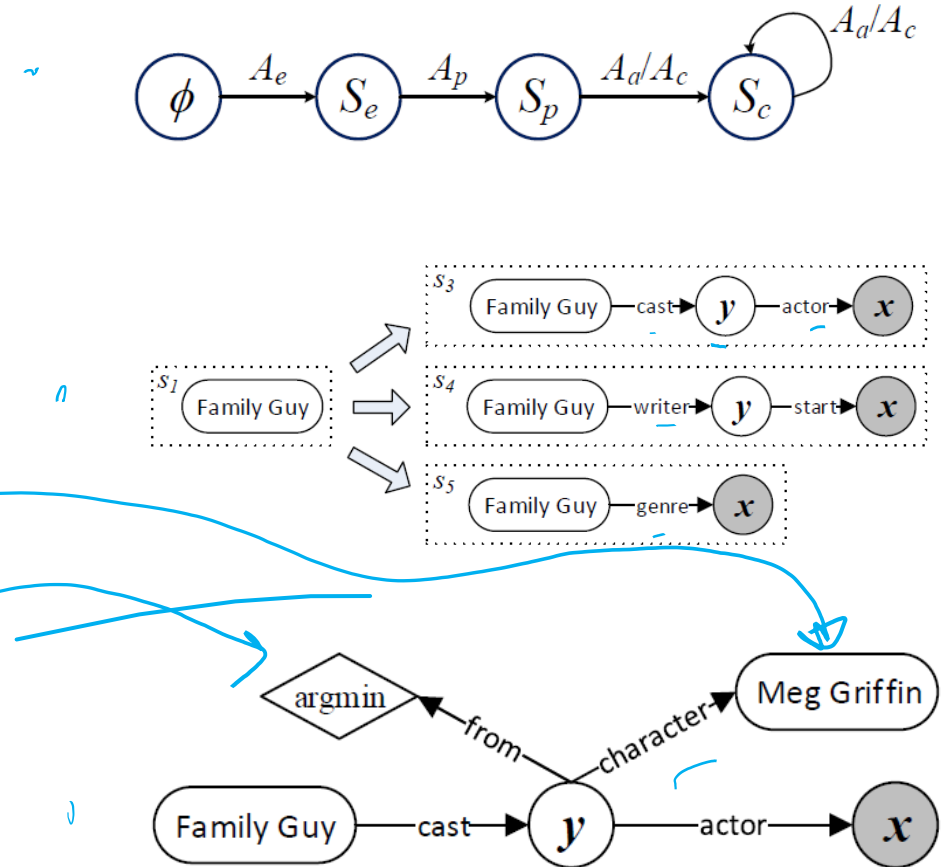
Figure 6: The architecture of the convolutional neural networks (CNN) used in this work. The CNN model maps a variable-length word sequence (e.g., a pattern or predicate sequence) to a low-dimensional vector in a latent semantic space. See text for the description of each layer.

Augmenting constraints and aggregations

- Graph with only core chain can be executed to get candidate answers: Still too many!

- Use conditions as filters

- 1 ■ Entity constraint *Meg*
- 2 ■ Keywords like first, last, ... (rule-based)



Learning reward function

Query ranking features

Topic Entity E

Core Inferential Chain

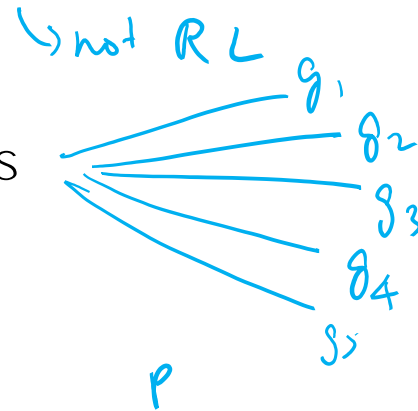
PatChain (2 CNNs) \downarrow Siamese

QuesEP (2 more CNNs) \downarrow Siamese

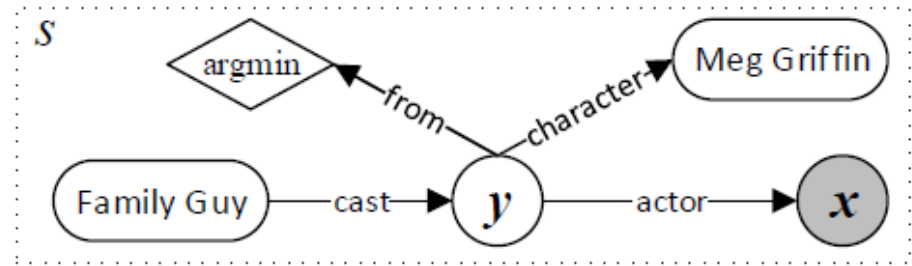
ClueWeb (1 more CNN) \downarrow FACL

Constraints and aggregations: Rules

Overall: #Answers fetched and #nodes in graph



$q = \text{"Who first voiced Meg on Family Guy?"}$

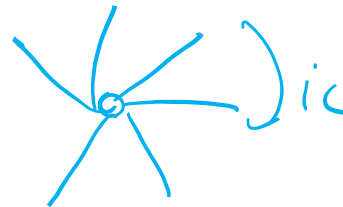


- (1) EntityLinkingScore(FamilyGuy, "Family Guy") = 0.9
- (2) PatChain("who first voiced meg on <e>", [cast-actor]) = 0.7
- (3) QuesEP(q , "family guy cast-actor") = 0.6
- (4) ClueWeb("who first voiced meg on <e>", cast-actor) = 0.2
- (5) ConstraintEntityWord("Meg Griffin", q) = 0.5
- (6) ConstraintEntityInQ("Meg 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.

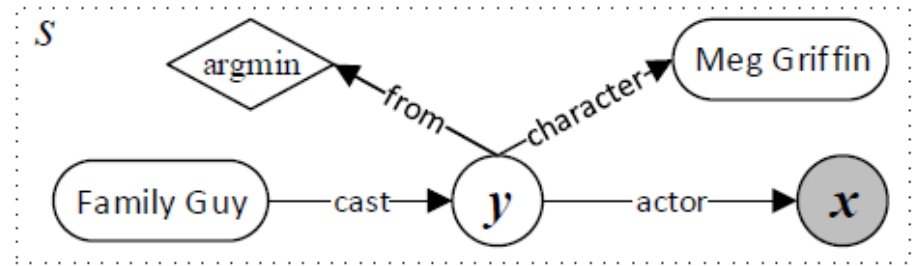
Learning reward function

- STAGG treats prediction as ranking task and ^{not} a binary classification problem where only the correct query graphs are labeled as positive



- Why? Closer to training data
- Parse-graph, F1-score ^(ret ~ an, gold an wq)
- ^{UTR} One-layer neural network model based on Lambda-rank → Get best graph, execute and bingo!

$q = \text{"Who first voiced Meg on Family Guy?"}$



- (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.

Cool ideas in STAGG

- ✓ ■ Graph-based search
- ✓ ■ Context subgraph
- ✓ ■ Execution order (s, p, c)
- ✓ ■ Top-k NED
- ✓ ■ Character trigrams
- ✓ ■ (Preliminary) ideas for handling complexity: core^{simple} chain and conditions ✓
- ✓ ■ Distant supervision for training ^{Q-Q from Q-A}
- ✓ ■ CNN models for similarity (representation)
- ✓ ■ Pose prediction as neural LTR

Research paper 2

Knowledge Graph Embedding Based Question Answering

[Knowledge graph embedding based question answering](#)

[\[PDF\] acm.org](#)

[X Huang, J Zhang, D Li, P Li - ... Conference on Web Search and Data ..., 2019 - dl.acm.org](#)

Question answering over knowledge graph (QA-KG) aims to use facts in the knowledge graph (KG) to answer natural language questions. It helps end users more efficiently and more easily access the substantial and valuable knowledge in the KG, without knowing its data structures. QA-KG is a nontrivial problem since capturing the semantic meaning of natural language is difficult for a machine. Meanwhile, many knowledge graph embedding methods have been proposed. The key idea is to represent each predicate/entity as a low ...

☆  Cited by [39](#) [Related articles](#) [All 6 versions](#)

The KEQA System

- Leverages knowledge graph embeddings!

- S Head P entity, O predicate, tail entity

- Uses the TransE Model

$QA \sim KG \sim KG \sim QA? \sim KB \sim QA?$

for QA

KG C, link prediction

$e_h \xrightarrow{p} e_t$
sem. capital -> behavior

Trans H, R, TAE

Translating embeddings for modeling multi-relational data

[PDF] nips.cc

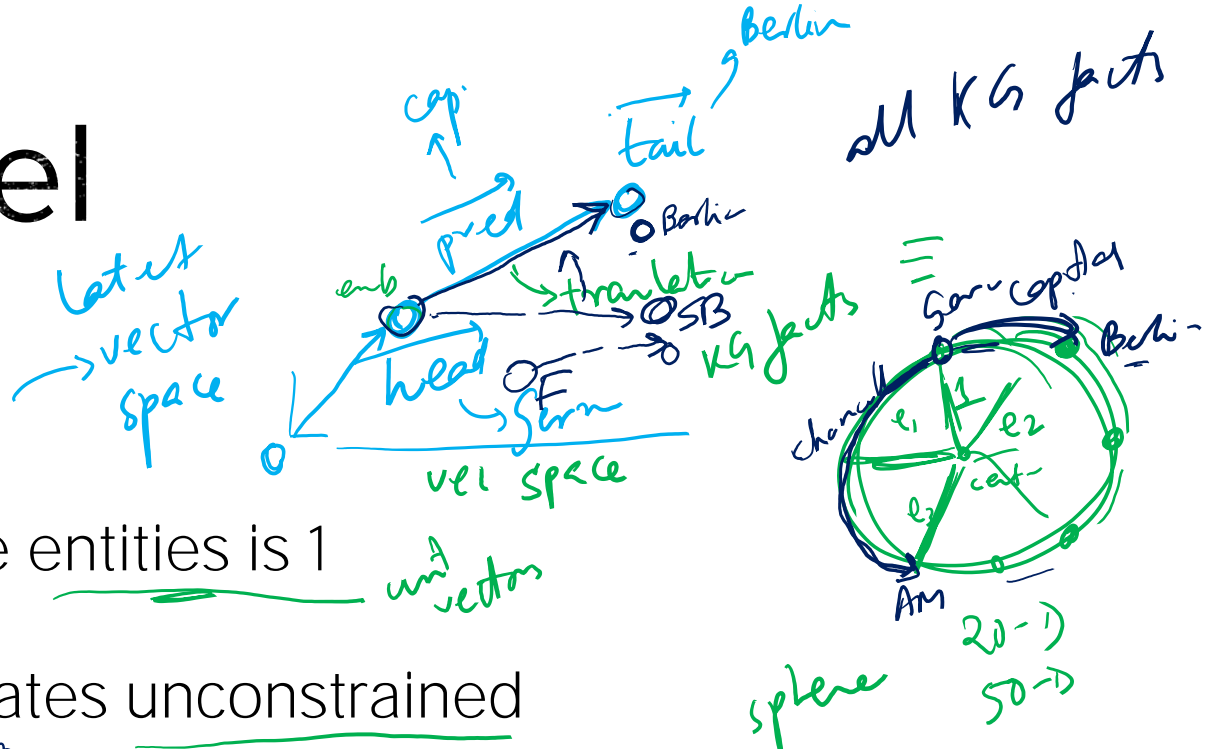
A Bordes, N Usunier, A Garcia-Duran... - Advances in neural ..., 2013 - papers.nips.cc

We consider the problem of embedding entities and relationships of multi-relational data in low-dimensional vector spaces. Our objective is to propose a canonical model which is easy to train, contains a reduced number of parameters and can scale up to very large databases. Hence, we propose, TransE, a method which models relationships by interpreting them as translations operating on the low-dimensional embeddings of the entities. Despite its simplicity, this assumption proves to be powerful since extensive experiments show that ...

☆ 77 Cited by 2086 Related articles All 19 versions >>

The TransE model

1. Head entity, predicate, tail entity
2. L2-norm of the embeddings of the entities is 1
- L2-norm of embeddings of predicates unconstrained



trve pair → KG facts
-ve pair → corrupted

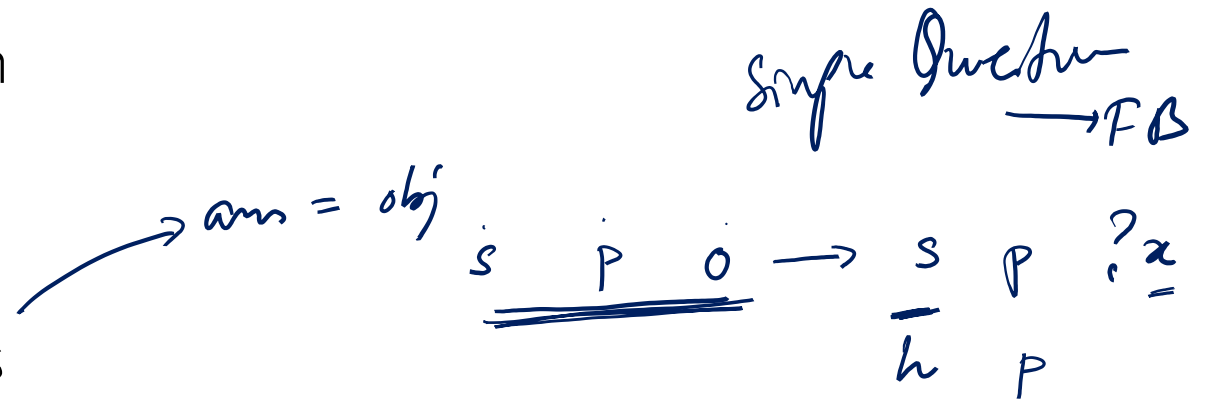
$$\mathcal{L} = \sum_{(h, \ell, t) \in S} \sum_{(h', \ell, t') \in S'_{(h, \ell, t)}} [\gamma + \underbrace{d(\underline{h} + \underline{\ell}, \underline{t})}_S - \underbrace{d(\underline{h}' + \underline{\ell}, \underline{t'})}_{F, C, B}]_+$$

{ Ger cap Berl
→ Frank cap Berl
Sem cap SB }

$$S'_{(h, \ell, t)} = \{(h', \ell, t) | h' \in E\} \cup \{(h, \ell, t') | t' \in E\}$$

The KEQA System

- Leverages knowledge graph embeddings!
- Uses the TransE Model (or TransE-like ...)
- Similar motivations
 - Name ambiguity
 - Predicate vocabulary mismatch
- From Baidu Research
- Simple questions, no qualifiers



KEQA: Notations

- Given a question, we want to predict corr. fact's head entity and predicate

$$\underline{\mathbf{e}_t} \approx f(\mathbf{e}_h, \mathbf{p}_\ell)$$

Handwritten annotations: An arrow points from the word "tail" to \mathbf{e}_t . Another arrow points from a "+" sign to the function f . A third arrow points from the expression $\mathbf{e}_h + \mathbf{p}_\ell$ to the function f .

Table 1: The important symbols and their definitions.

Notations	Definitions
\mathcal{G}	a knowledge graph
(h, ℓ, t)	a fact, i.e., (head entity, <u>predicate</u> , tail entity)
\tilde{Q}	a set of simple questions with ground truth facts
M	total number of predicates in \mathcal{G}
N	total number of entities in \mathcal{G}
d	dimension of the embedding representations
$\mathbf{P} \in \mathbb{R}^{M \times d}$	embedding representations of all predicates in \mathcal{G}
$\mathbf{E} \in \mathbb{R}^{N \times d}$	embedding representations of all entities in \mathcal{G}
$f(\cdot)$	relation function, given (h, ℓ, t) , $\Rightarrow \mathbf{e}_t \approx f(\mathbf{e}_h, \mathbf{p}_\ell)$
$\hat{\mathbf{p}}_\ell \in \mathbb{R}^{1 \times d}$	predicted predicate representation
$\hat{\mathbf{e}}_h \in \mathbb{R}^{1 \times d}$	predicted head entity representation
HED	Head Entity Detection model
$\text{HED}_{\text{entity}}$	head entity name tokens returned by the HED
HED_{non}	non entity name tokens returned by the HED

KEQA: Outline

- Based on Q and their corresponding predicates' embeddings, KEQA trains predicate learning model
- Takes a question as the input and returns a predicate vector that lies in KG embedding space

Table 1: The important symbols and their definitions.

Notations	Definitions
\mathcal{G}	a knowledge graph
(h, ℓ, t)	a fact, i.e., (head entity, predicate, tail entity)
Q	a set of simple questions with ground truth facts
M	total number of predicates in \mathcal{G}
N	total number of entities in \mathcal{G}
d	dimension of the embedding representations
$P \in \mathbb{R}^{M \times d}$	embedding representations of all predicates in \mathcal{G}
$E \in \mathbb{R}^{N \times d}$	embedding representations of all entities in \mathcal{G}
$f(\cdot)$	relation function, given $(h, \ell, t), \Rightarrow \mathbf{e}_t \approx f(\mathbf{e}_h, \mathbf{p}_\ell)$
$\hat{\mathbf{p}}_\ell \in \mathbb{R}^{1 \times d}$	predicted predicate representation
$\hat{\mathbf{e}}_h \in \mathbb{R}^{1 \times d}$	predicted head entity representation
HED	Head Entity Detection model
$\text{HED}_{\text{entity}}$	head entity name tokens returned by the HED
HED_{non}	non entity name tokens returned by the HED

KEQA: Outline

Table 1: The important symbols and their definitions.

Notations	Definitions
\mathcal{G}	a knowledge graph
(h, ℓ, t)	a fact, i.e., (head entity, predicate, tail entity)
Q	a set of simple questions with ground truth facts
M	total number of predicates in \mathcal{G}
N	total number of entities in \mathcal{G}
d	dimension of the embedding representations
$P \in \mathbb{R}^{M \times d}$	embedding representations of all predicates in \mathcal{G}
$E \in \mathbb{R}^{N \times d}$	embedding representations of all entities in \mathcal{G}
$f(\cdot)$	relation function, given $(h, \ell, t), \Rightarrow e_t \approx f(e_h, p_\ell)$
$\hat{p}_\ell \in \mathbb{R}^{1 \times d}$	predicted predicate representation
$\hat{e}_h \in \mathbb{R}^{1 \times d}$	predicted head entity representation
HED	Head Entity Detection model
HED _{entity}	head entity name tokens returned by the HED
HED _{non}	non entity name tokens returned by the HED

- Similarly for head entity: head entity learning
 - Head entity detection $\sim NED$
 - Get tail entity $\text{predicted representation}$
 - Get closest fact in KG using distance function
- $$\hat{e}_t = f(\hat{e}_h, \hat{p}_\ell)$$
- $$(\hat{e}_h, \hat{p}_\ell, \hat{e}_t) \Rightarrow KG$$

KEQA: Outline

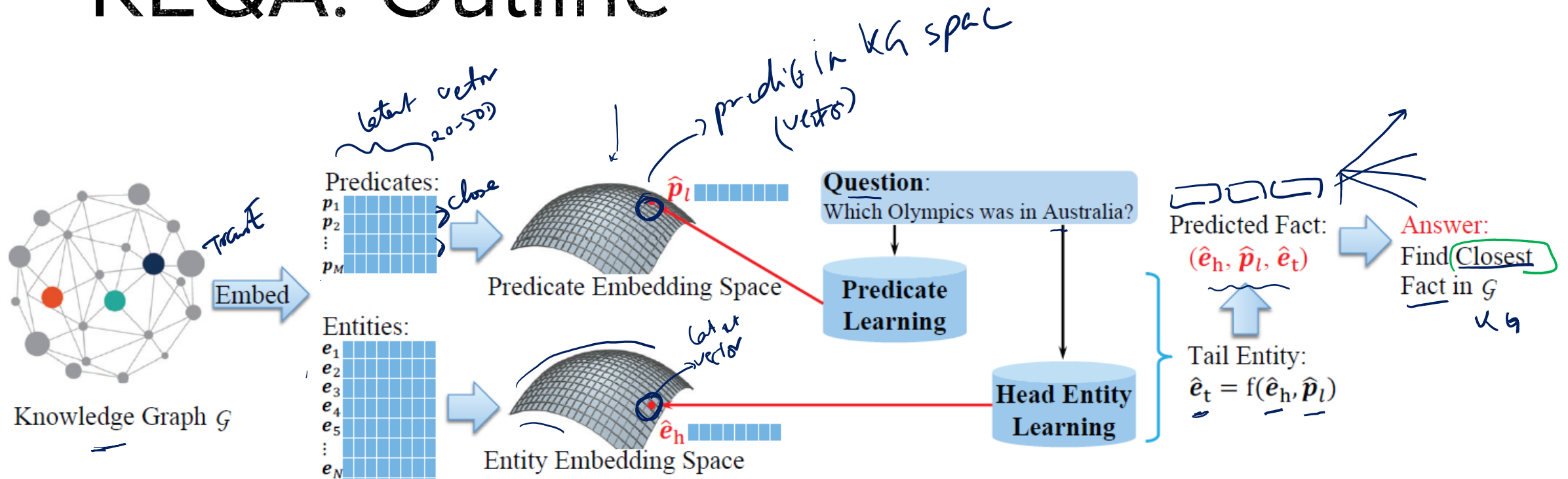


Figure 1: Instead of inferring the head entity and predicate directly, KEQA targets at jointly recovering the question's head entity, predicate, and tail entity representations $(\hat{e}_h, \hat{p}, \hat{e}_t)$ in the knowledge graph embedding spaces.



Predicate and head entity learning models

- Bi-LSTM for word order → in \mathcal{G}
- Useful model for sequence models for NLP
- Attention for word importance
imp wts
- Learning representations
generalization to unseen predicates at test time!!

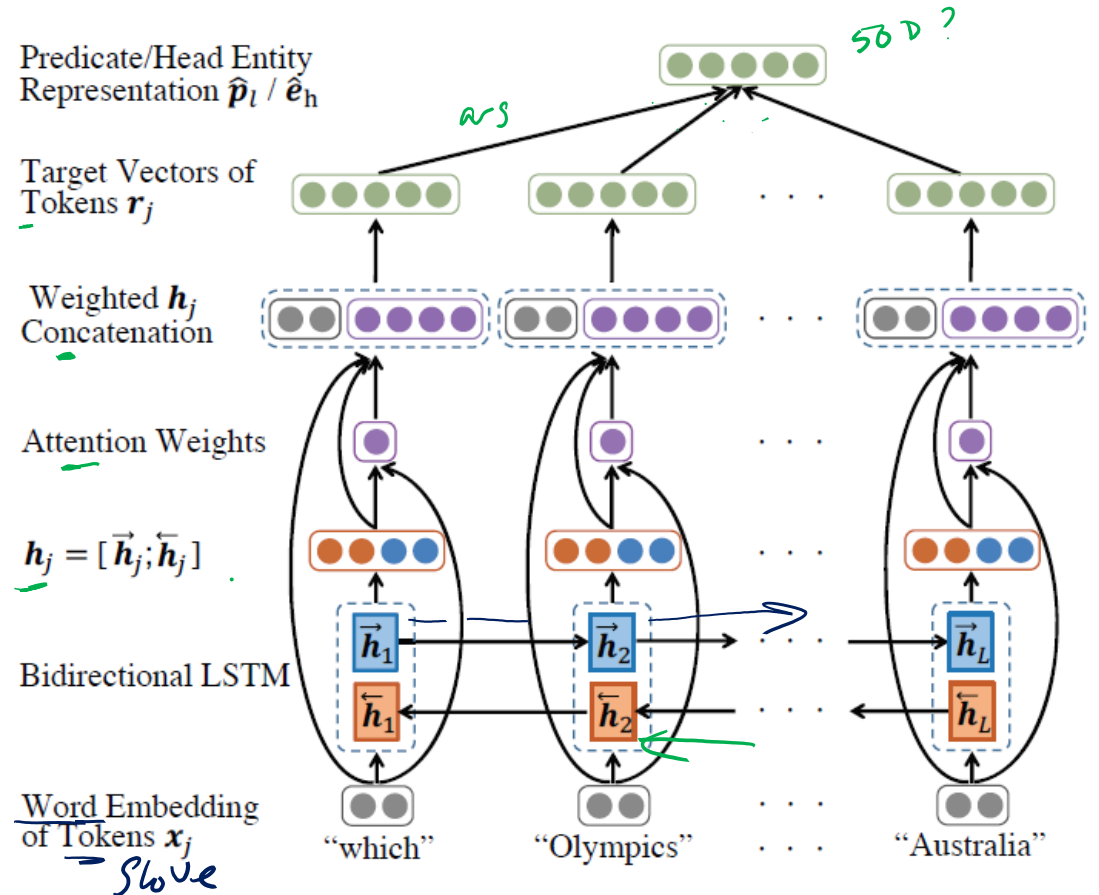


Figure 2: The architecture of the proposed predicate and head entity learning models.

Head entity detection model

- Getting the right entity (embedding) is crucial
- Need to reduce entity embedding search space
- Only some words matter
- No attention component
- 2 values predicted per token

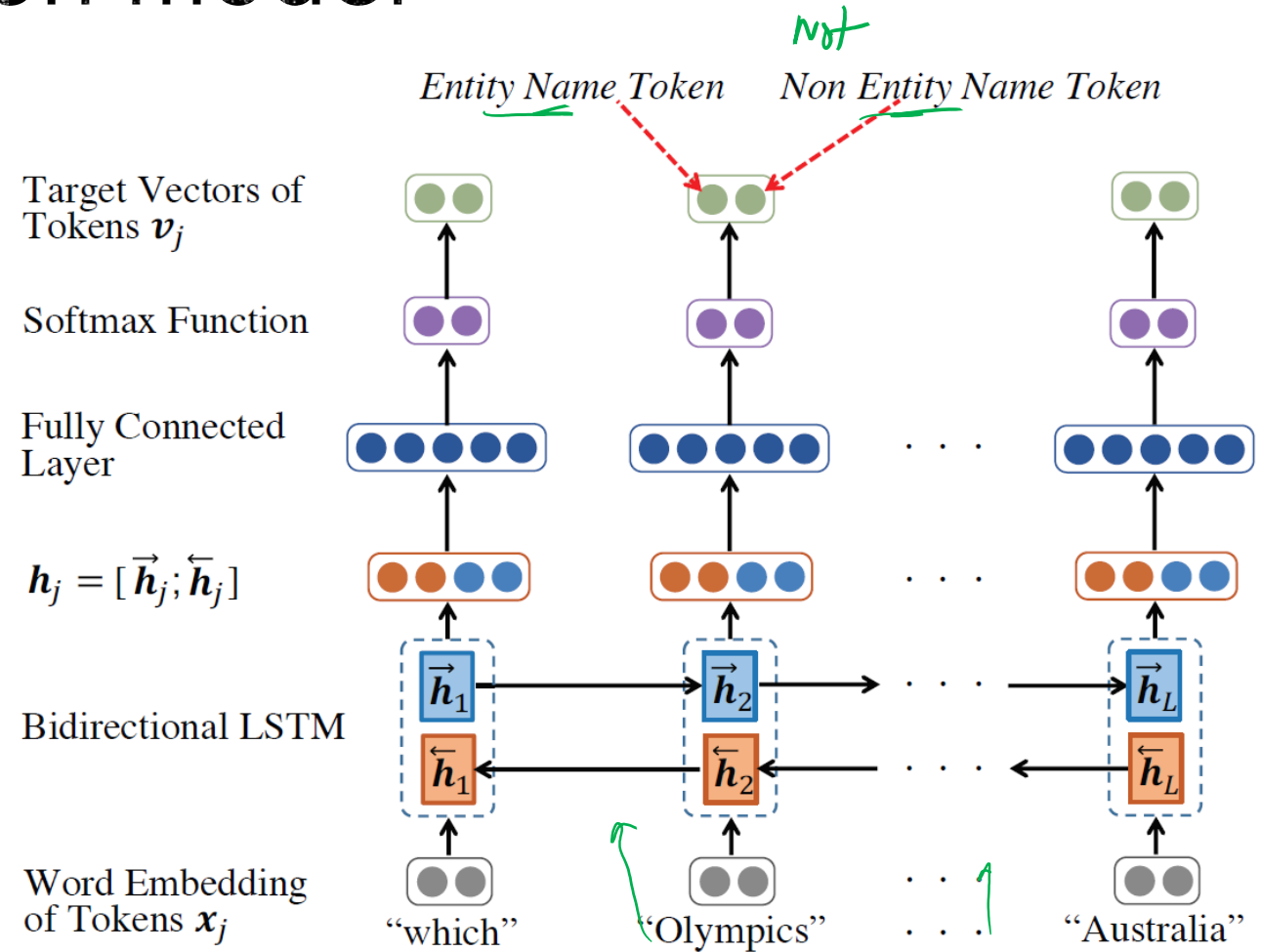


Figure 3: Structure of Head Entity Detection (HED) model.

Joint search on embedding spaces

- We now have predicted entity and predicate vectors
- Also, the head entities
- Now all we need is to define a distance metric for search!

$$\begin{aligned}
 & \text{arg} \underset{(h, \ell, t) \in C}{\text{minimize}} \quad \underbrace{\|\mathbf{p}_\ell - \hat{\mathbf{p}}_\ell\|_2}_{\text{embedding sim}} + \underbrace{\beta_1 \|\mathbf{e}_h - \hat{\mathbf{e}}_h\|_2}_{\text{entity char}} + \underbrace{\beta_2 \|f(\mathbf{e}_h, \mathbf{p}_\ell) - \hat{\mathbf{e}}_t\|_2}_{\text{tail}} \\
 & \quad - \underbrace{\beta_3 \text{sim}[n(h), \text{HED}_{\text{entity}}] + \beta_4 \text{sim}[n(\ell), \text{HED}_{\text{non}}]}_{\text{string similarity}}.
 \end{aligned}$$

Conclusions

- Neural methods have shown great promise in KG-QA
- Think: Guidelines for where to apply neural learning in QA pipeline
- Graph representations of KBs very flexible and efficient for search
- Need large training data
- Lot depends on effective sampling of positive and negative pairs to train loss function

*Thank
you*