

# Entity Representation and Retrieval from Knowledge Graphs

Alexander Kotov

Textual Data Analytics Lab, Department of Computer Science, Wayne State University



RuSSIR 2016

# OVERVIEW

Entities and Entity Retrieval

Knowledge Graphs

Entity Representation

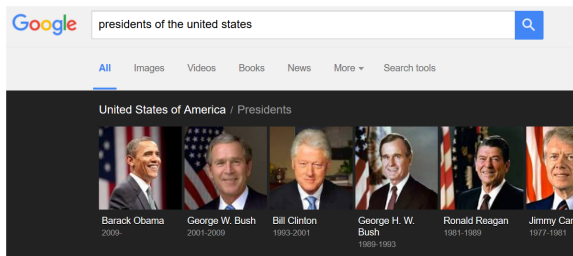
Entity Retrieval

Conclusion

# ENTITIES AND ENTITY RETRIEVAL

- ▶ **Entities:** material objects or concepts that exist in the real world or fiction (e.g. people, books, conferences, colors etc.)
- ▶ Entities (named entities) are typically designated by proper nouns or proper noun phrases (e.g. Barack Obama)
- ▶ **Entity retrieval:** answering arbitrary information needs related to particular aspects of objects (entities), expressed in unconstrained natural language and resolved using a collection of entities [Pound, Mika et al., WWW'10]

# AD-HOC ENTITY RETRIEVAL



- ▶ **Query:** keyword query corresponding to an entity name, description of property (properties) of the target entity or a set of entities
- ▶ “Telegraphic” queries – neither well-formed, nor grammatically correct sentences or questions
- ▶ **Results:** rank list of entities (entity representations) instead of or in addition to documents.

# AD-HOC ENTITY RETRIEVAL



Alan Moore graphic novels adapted to film



About 186,000 results (0.37 seconds)

## Alan Moore - Wikipedia, the free encyclopedia

[https://en.wikipedia.org/wiki/Alan\\_Moore](https://en.wikipedia.org/wiki/Alan_Moore) - Wikipedia

For other people named **Alan Moore**, see **Alan Moore** (disambiguation). ... Frequently described as the best **graphic novel** writer in history, he has been called "one of the most important British writers of the ... 4 **Film adaptations**; 5 **Personal life**.

## From Hell - Wikipedia, the free encyclopedia

[https://en.wikipedia.org/wiki/From\\_Hell](https://en.wikipedia.org/wiki/From_Hell) - Wikipedia

From Hell is a **graphic novel** by writer **Alan Moore** and artist Eddie Campbell, ... The comic was loosely **adapted** into a **film** of the same title, released in 2001.

## Watchmen - Wikipedia, the free encyclopedia

<https://en.wikipedia.org/wiki/Watchmen> - Wikipedia

For the 2009 **film adaptation**, see **Watchmen (film)**. ... The series was created by writer **Alan Moore**, artist Dave Gibbons, and colorist John Higgins. .... Moore reasoned that **MLJ Comics'** **Mighty Crusaders** might be available for such a project, ...

## Alan Moore - IMDb

[www.imdb.com/name/nm0600872/](http://www.imdb.com/name/nm0600872/) - Internet Movie Database

**Alan Moore** was born on November 18, 1953 in Northampton, England. ... Includes clips from **GLORY DAZE**, **COMMUNITY** and horror film **URBAN** ... 2009 **Tales of the Black Freighter** (Video short) (**graphic novel** "Watchmen" - uncredited).

[Biography](#) - [Awards](#) - [Photo Gallery](#) - [Publicity](#)

## IMDb: Graphic Novel Adaptations - a list by deano11

[www.imdb.com/list/ls005377599/](http://www.imdb.com/list/ls005377599/) - Internet Movie Database

# OVERVIEW

Entities and Entity Retrieval

Knowledge Graphs

Entity Representation

Entity Retrieval

Conclusion

# SUBJECT-PREDICATE-OBJECT (RDF) TRIPLES

- ▶ One way to represent knowledge in machine readable way
- ▶ *Subjects* correspond to entities designated by an identifier (URI `http://dbpedia.org/page/Barack_Obama` in case of DBpedia)
- ▶ Entities are connected with other entities, literals or scalars by relations or *predicates* (e.g. *hasGenre*, *knownFor*, *marriedTo*, *isPCmemberOf* etc.)
- ▶ Each triple represents a simple fact (e.g. `<http://dbpedia.org/page/Barack_Obama, marriedTo, http://dbpedia.org/page/Michelle_Obama>`)
- ▶ Many SPO triples → *knowledge graph*



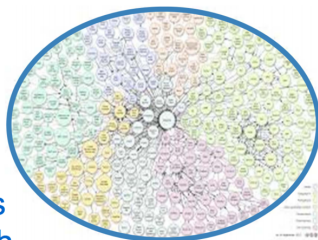
# EXISTING KNOWLEDGE GRAPHS

 Freebase

  
yago  
select knowledge

  
DBpedia

Facebook's  
Entity Graph



Microsoft's  
Satori



OpenIE  
(Reverb, OLLIE)

Google's  
Knowledge Graph



# LINKED OPEN DATA

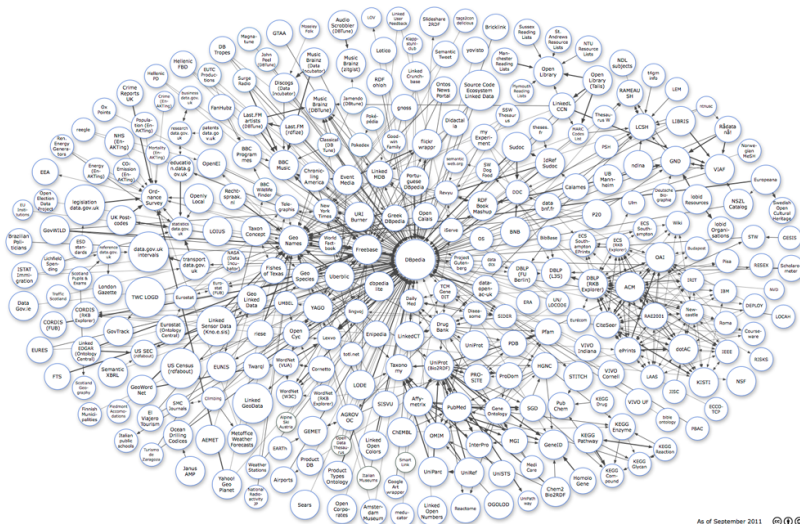
- ▶ Individual knowledge repositories can be published in machine readable form (RDF)
- ▶ The repositories can be connected to each other → Linked Open Data (LOD) cloud



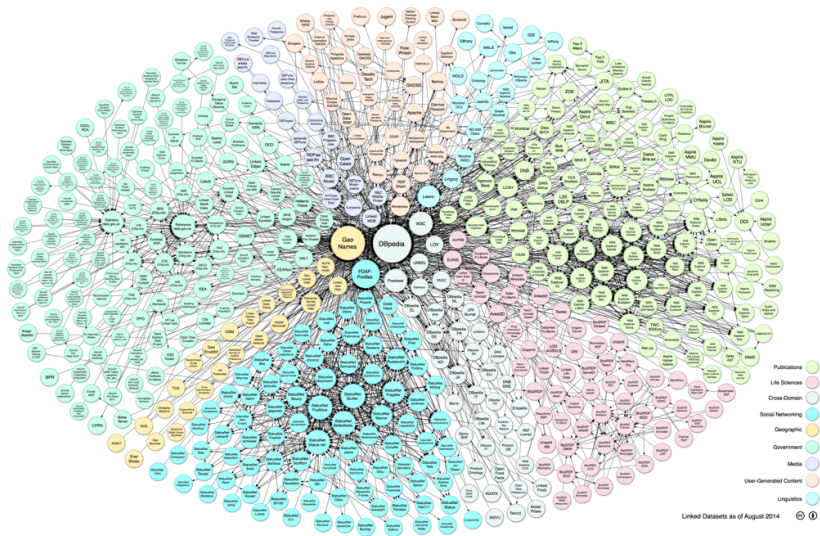




# LOD Cloud (CIRCA 2011)



# LOD CLOUD (CURRENT STATE)



# DBPEDIA ENTITY PAGE

dbpedia.org/page/Barack\_Obama

Search

DBpedia

Browse using - Formats -

Faceted Browser Sparql Endpoint

## About: Barack Obama

An Entity of Type : [agent](#), from Named Graph : <http://dbpedia.org>, within Data Space : [dbpedia.org](http://dbpedia.org)

Barack Hussein Obama II (US /bəˈrɑːk huːˈseɪn əˈbɑːmə/; born August 4, 1961) is the 44th and current President of the United States, and the first African American to hold the office. Born in Honolulu, Hawaii, Obama is a graduate of Columbia University and Harvard Law School, where he served as president of the Harvard Law Review. He was a community organizer in Chicago before earning his law degree.

Property	Value
<a href="#">dbpedia:abstract</a>	<ul style="list-style-type: none"> <li>Barack Hussein Obama II (US <span><span>/<span><span>b</span><span>ə</span><span>ˈ</span><span>r</span><span>ɑː</span><span>k</span><span> </span><span>h</span><span>uː</span><span>ˈ</span><span>s</span><span>eɪ</span><span>n</span><span> </span><span>ə</span><span>ˈ</span><span>b</span><span>ɑː</span><span>m</span><span>ə</span></span>/</span></span>; born August 4, 1961) is the 44th and current President of the United States, and the first African American to hold the office. Born in Honolulu, Hawaii, Obama is a graduate of Columbia University and Harvard Law School, where he served as president of the Harvard Law Review. He was a community organizer in Chicago before earning his law degree. He worked as a civil rights attorney and taught constitutional law at University of Chicago Law School between 1992 and 2004. He served three terms representing the 13th District in the Illinois Senate from 1997 to 2004, running unsuccessfully for the United States House of Representatives in 2000.In 2004, Obama received national attention during his campaign to represent Illinois in the United States Senate with his victory in the March Democratic Party primary, his keynote address at the Democratic National Convention in July, and his election to the Senate in November. He began his presidential campaign in 2007 and, after a close primary campaign against Hillary Rodham Clinton in 2008, he won sufficient delegates in the Democratic Party primaries to receive the presidential nomination. He then defeated Republican nominee John McCain in the general election, and was inaugurated as president on January 20, 2009. Nine months after his inauguration, Obama was named the 2009 Nobel Peace Prize laureate.During his first two</li> </ul>





# ENTITY RETRIEVAL FROM KNOWLEDGE GRAPH(S) (ERKG)

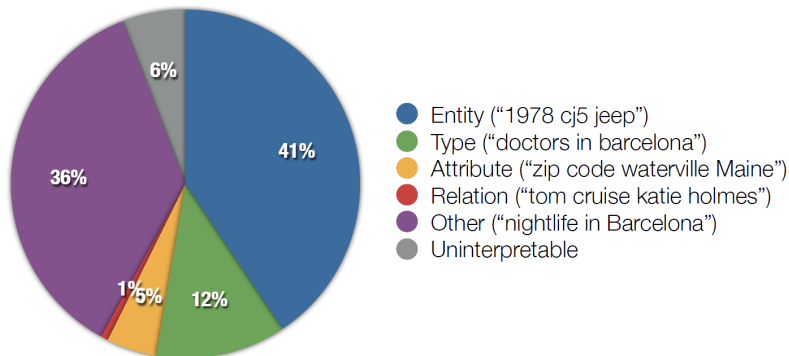
- ▶ Knowledge graphs are perfectly suited for addressing the information needs that aim at finding specific objects (entities) rather than documents
- ▶ Entity retrieval is a unique and interesting IR problem, since there is *no notion of a document*
- ▶ *Ad-hoc* Entity Retrieval assumes keyword queries (structured queries are studied more in the DB community)

# TYPICAL ERKG TASKS

- ▶ **Entity Search:** simple queries aimed at finding a particular entity or an entity which is an attribute of another entity
  - ▶ *“Ben Franklin”*
  - ▶ *“England football player highest paid”*
  - ▶ *“Einstein Relativity theory”*
- ▶ **List Search:** descriptive queries with several relevant entities
  - ▶ *“US presidents since 1960”*
  - ▶ *“animals lay eggs mammals”*
  - ▶ *“Formula 1 drivers that won the Monaco Grand Prix”*
- ▶ **Question Answering:** queries are questions in natural language
  - ▶ *“Who founded Intel?”*
  - ▶ *“For which label did Elvis record his first album?”*

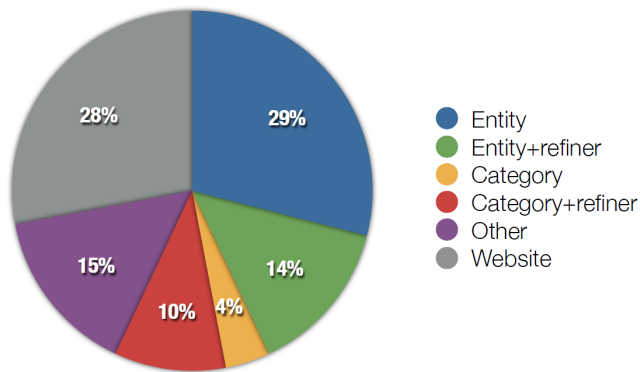
# DISTRIBUTION OF ENTITY WEB SEARCH QUERIES

[POUND ET AL. WWW'10]



# DISTRIBUTION OF ENTITY WEB SEARCH QUERIES

[LIN ET AL. WWW'11]



# RESEARCH CHALLENGES IN ERKG

1. How to design entity representations that capture the semantics of entity properties/relations and are effective for entity retrieval?
2. How to develop accurate and efficient entity retrieval models?

# ENTITY REPRESENTATION METHODS (DAY 1)

- Neumayer, Balog et al. On the Modeling of Entities for Ad-hoc Entity Search in the Web of Data, ECIR'12
- Neumayer, Balog et al. When Simple is (more than) Good Enough: Effective Semantic Search with (almost) no Semantics, ECIR'12
- Zhiltsov and Agichtein. Improving Entity Search over Linked Data by Modeling Latent Semantics, CIKM'13
- Zhiltsov, Kotov et al. Fielded Sequential Dependence Model for Ad-hoc Entity Retrieval in the Web of Data, SIGIR'15

## ENTITY RETRIEVAL MODELS (DAY 2)

- Classic unigram bag-of-words models for structured document retrieval, such as BM25F, Mixture of Language Models (MLM), Probabilistic Retrieval Model for Semi-structured Data (PRMS)
- Dali and Fortuna. Learning to Rank for Semantic Search, WWW'11
- Tonon, Demartini et al. Combining Inverted Indices and Structured Search for Ad-hoc Object Retrieval, SIGIR'12
- Sawant and Chakrabarti. Learning Joint Query Interpretation and Response Ranking, WWW'13
- Zhiltsov, Kotov et al. Fielded Sequential Dependence Model for Ad-hoc Entity Retrieval in the Web of Data, SIGIR'15
- Nikolaev, Kotov et al. Parameterized Fielded Term Dependence Models for Ad-hoc Entity Retrieval from Knowledge Graph, SIGIR'16

# OVERVIEW

Entities and Entity Retrieval

Knowledge Graphs

**Entity Representation**

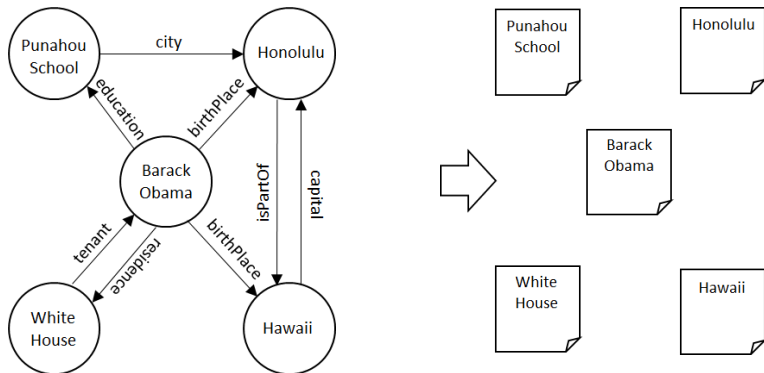
Entity Retrieval

Conclusion



# FROM ENTITY GRAPH TO ENTITY DOCUMENTS

Build a textual representation (i.e. “document”) for each entity by considering all triples, where it stands as subject (or object)



# LANGUAGE MODELING APPROACH

- Retrieval score of  $D$  is the likelihood of it being relevant to a given query  $Q$
- Query Likelihood retrieval model: retrieval score of  $D$  is the likelihood of generating  $Q$  from  $\Theta_D$ , the language model of  $D$

$$P(D|Q) \stackrel{\text{rank}}{=} \frac{P(Q|D)P(D)}{P(Q)} \propto P(Q|D)P(D), \text{ where}$$

$$P(Q|D) = \prod_{q_i \in Q} P(q_i | \theta_D)^{n(q_i, Q)}$$

# ENTITY LANGUAGE MODEL

If each entity is represented as an unstructured document  $E$ :

$$P(E|Q) \propto P(E)P(Q|\theta_E) = P(E) \prod_{q_i \in Q} P(q_i|\theta_E)^{n(q_i, Q)}$$

# STRUCTURED ENTITY DOCUMENTS (1)

- ▶ Entity descriptions are naturally structured, entities can be represented as fielded documents
- ▶ Entity documents can be ranked using conventional IR models
- ▶ In the simplest case, each predicate corresponds to one document field
- ▶ However, there are infinitely many predicates → optimization of field importance weights is computationally intractable

## STRUCTURED ENTITY DOCUMENTS (2)

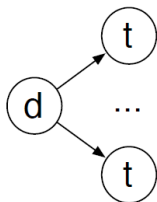
**Predicate folding:** group predicates together into a small set of predefined categories → entity documents with smaller number of fields

names	<ul style="list-style-type: none"> <li><b>foaf:name</b></li> <li><b>dbp:birthName</b></li> </ul>	<ul style="list-style-type: none"> <li>Barack Obama (en)</li> <li>Barack Hussein Obama II (en)</li> </ul>
attributes	<ul style="list-style-type: none"> <li><b>dbo:birthDate</b></li> <li><b>dbp:birthPlace</b></li> <li><b>dbo:office</b></li> </ul>	<ul style="list-style-type: none"> <li>1961-08-04 (xsd:date)</li> <li>Honolulu, Hawaii, U.S. (en)</li> <li>44<sup>th</sup> President of the United States</li> </ul>
outgoing relations	<ul style="list-style-type: none"> <li><b>dbo:party</b></li> <li><b>dbo:region</b></li> <li><b>dbo:predecessor</b></li> </ul>	<ul style="list-style-type: none"> <li>dbr:Democratic_Party_(United_States)</li> <li>dbr:Illinois</li> <li>dbr:Peter_Fitzgerald_(politician)</li> <li>dbr:George_W._Bush</li> <li>dbr:Alice_Palmer_(politician)</li> </ul>
incoming relations	<ul style="list-style-type: none"> <li>is <b>dbo:tenant</b> of</li> <li>is <b>dbo:president</b> of</li> </ul>	<ul style="list-style-type: none"> <li>dbr:White_House</li> <li>dbr:Joe_Biden</li> </ul>

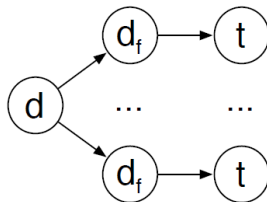
# PREDICATE FOLDING

- ▶ Grouping according to type (attributes, incoming/outgoing links)[[Pérez-Agüera et al. 2010](#)]
- ▶ Grouping according to importance (determined based on predicate popularity)[[Blanco et al. 2010](#)]

# MODEL COMPARISON



**Unstructured  
document model**



**Fielded  
document model**

## 2-FIELD ENTITY DOCUMENT

[NEUMAYER, BALOG ET AL., ECIR'12]

Each entity is represented as a two-field documents:

title

object values belonging to predicates ending with  
“name”, “label” or “title”

content

object values for 1000 most frequent predicates  
concatenated together into a flat text representation



## 3-FIELD ENTITY DOCUMENT

[ZHILTSOV AND AGICHTEIN, CIKM'13]

Each entity is represented as a three-field document:

names

literals of foaf:name, rdfs:label predicates along  
with tokens extracted from entity URIs

attributes

literals of all other predicates

outgoing links

names attributes of entities in the object position

## 5-FIELD ENTITY DOCUMENT

[ZHILTSOV, KOTOV ET AL., SIGIR'15]

Each entity is represented as a five-field document:

names

conventional names of the entities, such as the name of a person or the name of an organization

attributes

all entity properties, other than names

categories

classes or groups, to which the entity has been assigned

similar entity names

names of the entities that are very similar or identical to a given entity

related entity names

names of the entities that are part of the same RDF triple

## 5-FIELD ENTITY DOCUMENT EXAMPLE

Entity document for the entity *Barack Obama*.

Field	Content
names	barack obama barack hussein obama ii
attributes	44th current president united states birth place honolulu hawaii
categories	democratic party united states senator nobel peace prize laureate christian
similar entity names	barack obama jr barak hussein obama barack h obama ii
related entity names	spouse michelle obama illinois state predecessor george walker bush

# OVERVIEW

Entities and Entity Retrieval

Knowledge Graphs

Entity Representation

**Entity Retrieval**

Conclusion

# BM25F

[ROBERTSON AND ZARAGOZA, CIKM'04]

- ▶ Option 1: aggregation of BM25 scores across fields

$$P(E|Q) \stackrel{\text{rank}}{=} \sum_{q_i \in Q} \sum_{j=1}^F \log \frac{N}{df^j(q_i)} \frac{(k_1 + 1) \tilde{tf}^j(q_i)}{k_1((1 - b) + b \frac{|E^j|}{|E^j|_{\text{avg}}})}, 0 \leq b \leq 1$$

- ▶ Option 2 (more effective): field-specific length normalization

$$\tilde{tf}^j(q_i) = \sum_{j=1}^F w_j \frac{tf^j(q_i)}{B^j}$$

$$B^j = ((1 - b^j) + b^j \frac{|E^j|}{|E^j|_{\text{avg}}}), 0 \leq b^j \leq 1$$

$$P(E|Q) \stackrel{\text{rank}}{=} \sum_{q_i \in Q} \log \frac{N}{df^j(q_i)} \cdot \frac{(k_1 + 1) \tilde{tf}(q_i)}{k_1 + \tilde{tf}(q_i)}$$

# BM25F

[ROBERTSON AND ZARAGOZA, CIKM'04]

- Option 1: aggregation of BM25 scores across fields

$$P(E|Q) \stackrel{\text{rank}}{=} \sum_{q_i \in Q} \sum_{j=1}^F \log \frac{N}{df^j(q_i)} \frac{(k_1 + 1) \tilde{tf}^j(q_i)}{k_1((1 - b) + b \frac{|E^j|}{|E^j|_{\text{avg}}})}, 0 \leq b \leq 1$$

- Option 2 (more effective): field-specific length normalization

$$\tilde{tf}^j(q_i) = \sum_{j=1}^F w_j \frac{tf^j(q_i)}{B^j}$$

$$B^j = ((1 - b^j) + b^j \frac{|E^j|}{|E^j|_{\text{avg}}}), 0 \leq b^j \leq 1$$

$$P(E|Q) \stackrel{\text{rank}}{=} \sum_{q_i \in Q} \log \frac{N}{df^j(q_i)} \cdot \frac{(k_1 + 1) \tilde{tf}(q_i)}{k_1 + \tilde{tf}(q_i)}$$

# MIXTURE OF LANGUAGE MODELS

[OGILVIE AND CALLAN, SIGIR'03]

- ▶ Separate LM  $\theta_E^j$  is created for each field  $j$  of entity document  $E$
- ▶ Document LM is a linear combination of field LMs

$$P(Q|E) \stackrel{\text{rank}}{=} \prod_{q_i \in Q} P(q_i|\theta_E)^{tf(q_i)},$$

where

$$P(q_i|\theta_E) = \sum_{j=1}^F w_j P(q_i|\theta_E^j), \quad \sum_j w_j = 1$$

$$P(q_i|\theta_E^j) = \frac{tf_{q_i, E^j} + \mu_j \frac{cf_{q_i}^j}{|C_j|}}{|E^j| + \mu_j}$$

# SETTING FIELD WEIGHTS

- ▶ Heuristically: proportionate to the length of content in the field
- ▶ Empirically: by optimizing the target retrieval metric using training queries
- ▶ **Problems:**
  - ▶ Entities are sparse with respect to different fields (most entities have only a handful of predicates)
  - ▶ More fields in entity representations → more training data to optimize their weights



# PROBABILISTIC RETRIEVAL MODEL FOR SEMI-STRUCTURED DATA

[KIM, XUE AND CROFT, ECIR'09]

Extends Mixture of Language Models by dynamically determining the mapping of query terms onto entity document fields

$$P(q_i|\theta_E) = \sum_{j=1}^F w_j P(q_i|\theta_E^j), \quad \sum_j w_j = 1$$

$$P(q_i|\theta_E) = \sum_{j=1}^F P(E^j|q_i)P(q_i|\theta_E^j)$$

where

$$P(E^j|q_i) = \frac{P(q_i|E^j)P(E^j)}{\sum_{j=1}^F P(q_i|E^j)P(E^j)}$$

# PROBABILISTIC RETRIEVAL MODEL FOR SEMI-STRUCTURED DATA

[KIM, XUE AND CROFT, ECIR'09]

Extends Mixture of Language Models by dynamically determining the mapping of query terms onto entity document fields

$$P(q_i|\theta_E) = \sum_{j=1}^F w_j P(q_i|\theta_E^j), \quad \sum_j w_j = 1$$

$$P(q_i|\theta_E) = \sum_{j=1}^F P(E^j|q_i) P(q_i|\theta_E^j)$$

where

$$P(E^j|q_i) = \frac{P(q_i|E^j)P(E^j)}{\sum_{j=1}^F P(q_i|E^j)P(E^j)}$$

# PROBABILISTIC RETRIEVAL MODEL FOR SEMI-STRUCTURED DATA

[KIM, XUE AND CROFT, ECIR'09]

Extends Mixture of Language Models by dynamically determining the mapping of query terms onto entity document fields

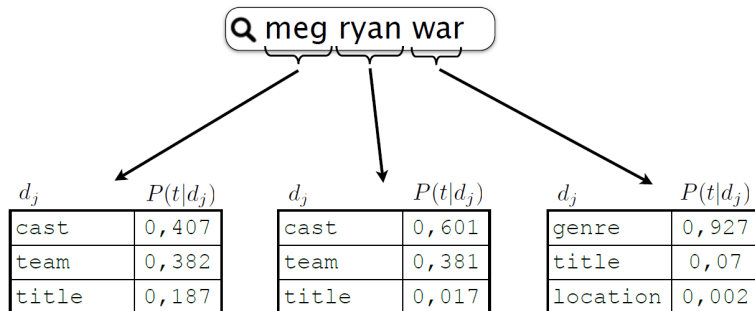
$$P(q_i|\theta_E) = \sum_{j=1}^F w_j P(q_i|\theta_E^j), \quad \sum_j w_j = 1$$

$$P(q_i|\theta_E) = \sum_{j=1}^F P(E^j|q_i)P(q_i|\theta_E^j)$$

where

$$P(E^j|q_i) = \frac{P(q_i|E^j)P(E^j)}{\sum_{j=1}^F P(q_i|E^j)P(E^j)}$$

## PRMS (EXAMPLE)



# HIERARCHICAL ENTITY MODEL (1)

[NEUMAYER, BALOG ET AL., ECIR'12]

Entity document fields are organized into a 2-level hierarchy:

- ▶ Predicate types are on the top level:

name

subject is  $E$ , object is literal and predicate comes from a predefined list (e.g. foaf:name or rdfs:label) or ends with "name", "label" or "title"

attributes

the subject is  $E$ , object is literal and the predicate is not of type name

outgoing links

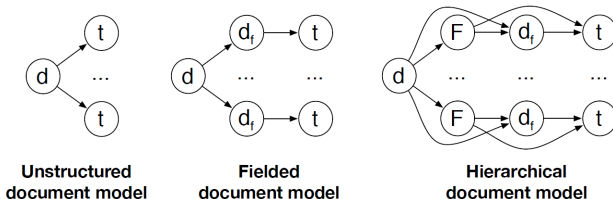
the subject is  $E$  and the object is a URI. URI is resolved by replacing it with entity name

incoming links

$E$  is an object, subject entity URI is resolved

- ▶ Individual predicates are at the bottom level

# HIERARCHICAL ENTITY MODEL (2)



$$\begin{aligned}
 P(q_i|\theta_E) &= \sum_{p_t} P(q_i|p_t, E)P(p_t|E) = \\
 &= \sum_{p_t} \left( \sum_{p \in p_t} P(q_i|p, p_t)P(p|p_t, E) \right) P(p_t|E)
 \end{aligned}$$

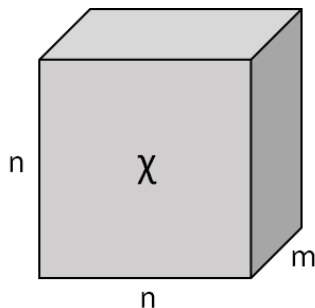
$P(q_i|p, p_t) = (1 - \lambda)P(q_i|p) + \lambda P(q_i|\theta_E^{p_t})$ , where  $P(q_i|p)$  ML estimate and  $P(q_i|\theta_E^{p_t})$  is Dirichlet-smoothed LM for predicate type  $p_t$

# LATENT DIMENSIONAL REPRESENTATION

[ZHILTSOV AND AGICHTEIN, CIKM'13]

- ▶ Compact representation of entities in low dimensional space by using a modified algorithm for tensor factorization
- ▶ Entities and entity-query pairs are represented with term-based and structural features

# KNOWLEDGE GRAPH AS TENSOR

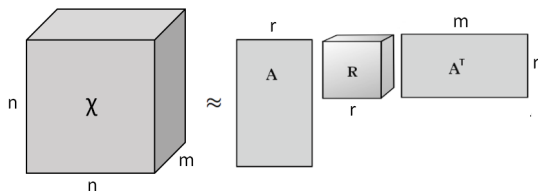


- ▶ For a knowledge graph with  $n$  distinct entities and  $m$  distinct predicates, we construct a tensor  $\mathcal{X}$  of size  $n \times n \times m$ , where  $\mathcal{X}_{ijk} = 1$ , if there is  $k$ -th predicate between  $i$ -th entity and  $j$ -th entity, and  $\mathcal{X}_{ijk} = 0$ , otherwise
- ▶ Each  $k$ -th frontal tensor slice  $\mathcal{X}_k$  is an adjacency matrix for the  $k$ -th predicate, which is sparse



# RESCAL TENSOR FACTORIZATION

[NIKEL, TRESP, ET AL., WWW'12]



- Given  $r$  is the number of latent factors, we factorize each  $X_k$  into the matrix product:

$$X_k = AR_kA^T, k = \overline{1, m},$$

where  $A$  is a dense  $n \times r$  matrix, a matrix of latent embeddings for entities, and  $R_k$  is an  $r \times r$  matrix of latent factors

- $A$  and  $R_k$  are solutions of the following optimization problem:

$$\min_{A, R} \frac{1}{2} \left( \sum_k \|X_k - AR_kA^T\|_F^2 \right) + \lambda \left( \|A\|_F^2 + \sum_k \|R_k\|_F^2 \right)$$

# RETRIEVAL METHOD

1. Retrieve initial set of entities using MLM
2. Re-rank the entities using Gradient Boosted Regression Tree (GBRT)

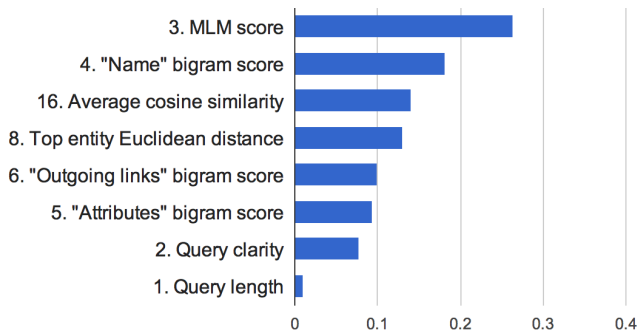
# FEATURES

#	Feature
<b>Term-based features</b>	
1	Query length
2	Query clarity
3	Uniformly weighted MLM score
4	Bigram relevance score for the "name" field
5	Bigram relevance score for the "attributes" field
6	Bigram relevance score for the "outgoing links" field
<b>Structural features</b>	
7	Top-3 entity cosine similarity, $\cos(\mathbf{e}, \mathbf{e}_{top})$
8	Top-3 entity Euclidean distance, $\ \mathbf{e} - \mathbf{e}_{top}\ $
9	Top-3 entity heat kernel, $e^{-\frac{\ \mathbf{e} - \mathbf{e}_{top}\ ^2}{\sigma}}$

# RESULTS

Features	Performance		
	NDCG	MAP	P@10
Term-based baseline	0.382	0.265	0.539
All features	<b>0.401 (+ 5.0%)*</b>	<b>0.276 (+ 4.2%)*</b>	<b>0.561 (+ 4.1%)*</b>

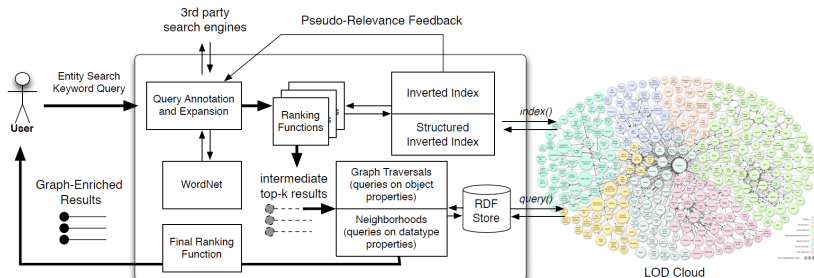
# FEATURE IMPORTANCE



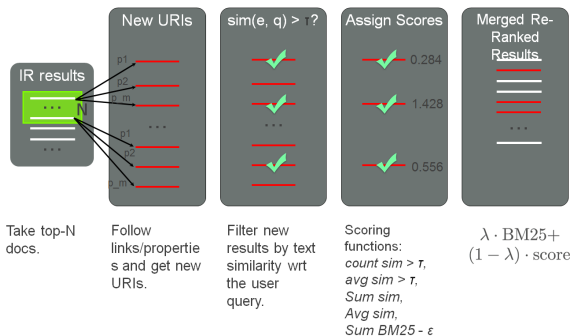
- ▶ Exploiting latent semantics of entities helps improve retrieval results (structural features improve NDCG and P@10)
- ▶ Most effective distance measures are cosine similarity and Euclidean distance
- ▶ However, the overall performance of the method is sensitive to top 3 retrieved results

# HYBRID IR AND DB ERKG METHODS

[TONON, DEMARTINI ET AL., SIGIR'12]



# HYBRID ERKG METHODS



Take top-N docs.

Follow links/properties and get new URIs.

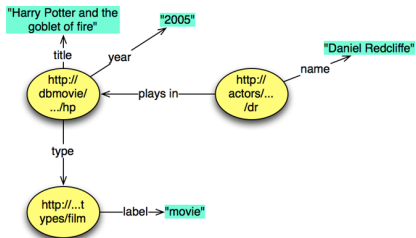
Filter new results by text similarity wrt the user query.

Scoring functions:  
*count sim* >  $\tau$ ,  
*avg sim* >  $\tau$ ,  
*Sum sim*,  
*Avg sim*,  
*Sum BM25 -  $\epsilon$*

$\lambda \cdot \text{BM25} + (1 - \lambda) \cdot \text{score}$

1. Retrieve an initial list of entities matching the query using standard retrieval function (BM25)
2. Expand the retrieved results by exploiting the structure of the knowledge graph (retrieved entities can be used as starting points for simple graph traversals, i.e. finding neighbors)
3. Filter out expanded results removing those with low similarity to the original query
4. Re-rank the results

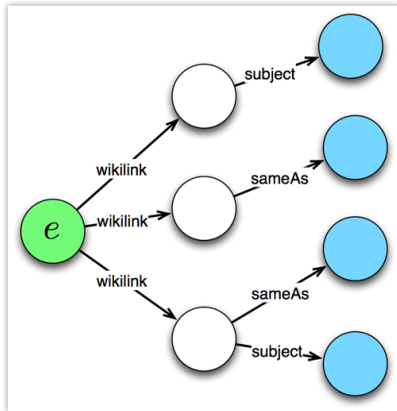
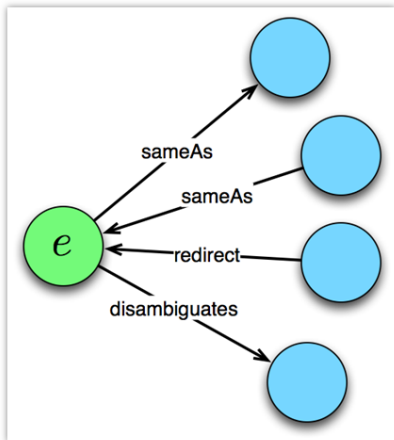
# RESULT EXPANSION STRATEGIES



- ▶ Follow predicates leading to other entities
- ▶ Follow datatype properties leading to additional entity attributes
- ▶ Explore just the neighborhood of a node and the neighbors of neighbors



# PREDICATES TO FOLLOW



# RESULTS

	2010 Collection			2011 Collection		
	MAP	P10	NDCG	MAP	P10	NDCG
BM25	0.2070	0.3348	0.5920	0.1484	0.2020	0.4267
SAMEAS	0.2293* (+11%)	0.3633* (+8%)	0.5932 (+0%)	0.1612 (+9%)	<b>0.2200</b> (+9%)	0.4433 (+4%)
S1_1	<b>0.2586*</b> (+25%)	<b>0.3848*</b> (+15%)	0.5965 (+1%)	0.1657 (+12%)	0.2140 (+6%)	0.4426 (+4%)
S1_2	0.2305* (+11%)	0.3217 (-4%)	0.5724* (-3%)	<b>0.1731</b> (+17%)	0.2180 (+8%)	<b>0.4532</b> (+6%)
S1_3	0.2306* (+11%)	0.3217 (-4%)	0.5721* (-4%)	0.1716 (+16%)	0.2140 (+6%)	0.4501 (+5%)
S2_1	0.2118 (+2%)	0.3370 (+1%)	0.5971 (+1%)	0.1550 (+4%)	0.2060 (+2%)	0.4376 (+3%)
S2_2	0.2118 (+2%)	0.3370 (+1%)	0.5965 (+1%)	0.1555 (+5%)	0.2080 (+3%)	0.4379 (+3%)
S2_3	0.2113 (+2%)	0.3402 (+2%)	<b>0.5978</b> (+1%)	0.1589 (+7%)	0.2120 (+5%)	0.4385 (+3%)

- ▶ The simple S1\_1 approach which exploits `<owl:sameAs>` links plus Wikipedia redirect and disambiguation information performs best obtaining 25% improvement of MAP over the BM25 baseline on the 2010 dataset

# LEARNING-TO-RANK METHOD FOR ENTITY RETRIEVAL

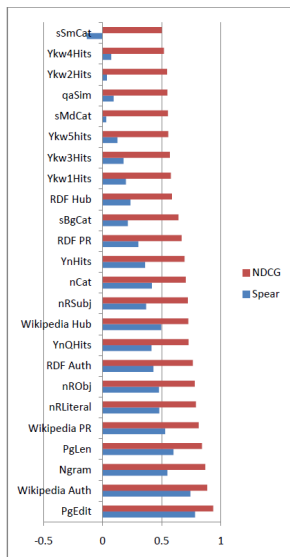
[DALI AND FORTUNA, WWW'11]

- ▶ Variety of features:
  - ▶ Popularity and importance of Wikipedia page: # of accesses from logs, # of edits, page length
  - ▶ RDF features: # of triples  $E$  is subject/object/subject and object is a literal, # of categories Wikipedia page for  $E$  belongs to, size of the biggest/smallest/median category
  - ▶ HITS scores and Pagerank of Wikipedia page and  $E$  in the RDF graph
  - ▶ # of hits from search engine API for the top 5 keywords from the abstract of Wikipedia page for  $E$
  - ▶ Count of entity name in Google N-grams
- ▶ RankSVM learning-to-rank method

# EVALUATION

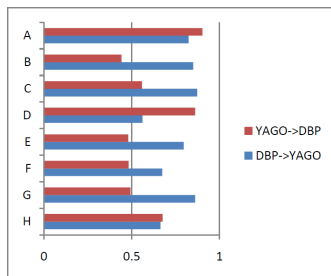
- ▶ Initial set of entities obtained using SPARQL queries
- ▶ 14 example queries for DBpedia and 27 example queries for Yago
- ▶ Example queries: “Which athlete was born in Philadelphia?”, “List of Schalke 04 players”, “Which countries have French as an official language?”, “Which objects are heavier than the Iosif Stalin tank?”

# FEATURE IMPORTANCE



- ▶ Features approximating the importance, hub and authority scores, PageRank of Wikipedia page are effective
- ▶ Google N-grams is effective proxy for entity popularity, cheaper than search engine API
- ▶ PageRank and HITS scores on RDF graph are not effective (outperformed by simpler RDF features)
- ▶ Feature combinations improve both robustness and accuracy of ranking

# TRANSFER LEARNING



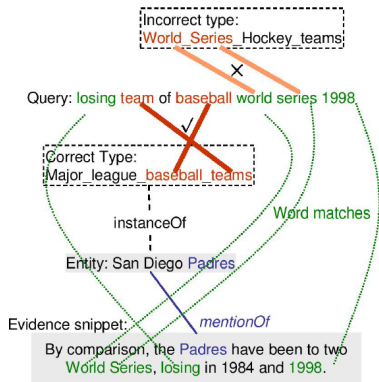
- ▶ Ranking model was trained on DBpedia questions and applied to Yago questions
- ▶ Only feature set A (all features) results in robust ranking model transfer
- ▶ In general, the ranking models for different knowledge graphs are non-transferable, unless they have been learned on large number of features
- ▶ The biggest inconsistencies occur on the models trained on graph based features → knowledge graphs preserve particularities reflecting their designer decisions

# JOINT TYPE DETECTION AND ENTITY RANKING

[SAWANT AND CHAKRABARTI, WWW'13]

- ▶ Method for answering “telegraphic” queries with target type
  - ▶ woodrow wilson president university
  - ▶ dolly clone institute
  - ▶ lead singer led zeppelin band
- ▶ Integrates type detection into ranking and considers multiple query interpretations
- ▶ Has generative and discriminative formulations

# METHOD

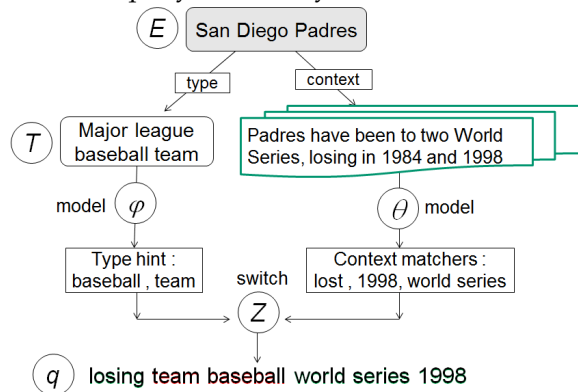


- ▶ All possible  $2^{|q|}$  query segmentations are considered
- ▶ Each query term is either a “type hint” or a “word matcher”



# GENERATIVE APPROACH

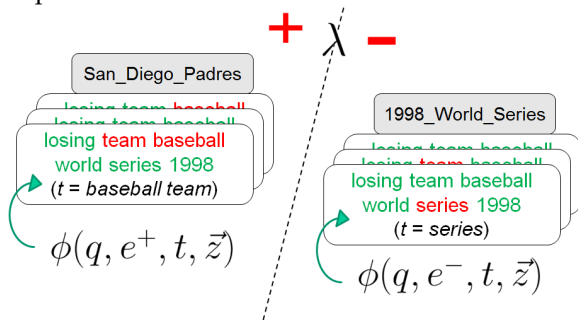
Generate query from entity



$$P(E|Q) \propto P(E) \sum_{t, \vec{z}} P(t|E) P(\vec{z}) P(h(\vec{q}, \vec{z})|t) P(s(\vec{q}, \vec{z})|E)$$

# DISCRIMINATIVE APPROACH

Separate correct and incorrect entities



$$\phi(q, e, t, \vec{z}) = \langle \phi_1(q, e), \phi_2(t, e), \phi_3(q, \vec{z}, t), \phi_4(q, \vec{z}, e) \rangle$$

# FIELDLED SEQUENTIAL DEPENDENCE MODEL

[ZHILTSOV, KOTOV ET AL., SIGIR'15]

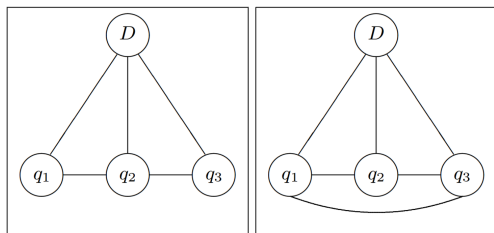
Previous research in ad-hoc IR has focused on two major directions:

- ▶ unigram bag-of-words retrieval models for multi-fielded documents
  - Ogilvie and Callan. Combining Document Representations for Known-item Search, SIGIR'03
  - Robertson et al. Simple BM25 Extension to Multiple Weighted Fields, CIKM'04
- ▶ retrieval models incorporating term dependencies
  - Metzler and Croft. A Markov Random Field Model for Term Dependencies, SIGIR'05
  - Huston and Croft. A Comparison of Retrieval Models using Term Dependencies, CIKM'14

**Goal:** to develop a retrieval model that captures both document structure and term dependencies

# SEQUENTIAL AND FULL DEPENDENCE MODELS

[METZLER AND CROFT, SIGIR'05]



Ranks w.r.t.  $P_{\Lambda}(D|Q) = \sum_{i \in \{T, U, O\}} \lambda_i f_i(Q, D)$

Potential function for unigrams is QL:

$$f_T(q_i, D) = \log P(q_i | \theta_D) = \log \frac{tf_{q_i, D} + \mu \frac{cf_{q_i}}{|C|}}{|D| + \mu}$$

SDM only considers two-word sequences in queries, FDM considers all two-word combinations.

# FSDM RANKING FUNCTION

FSDM incorporates document structure and term dependencies with the following ranking function:

$$P_{\Lambda}(D|Q) \stackrel{\text{rank}}{=} \lambda_T \sum_{q \in Q} \tilde{f}_T(q_i, D) + \\ \lambda_O \sum_{q \in Q} \tilde{f}_O(q_i, q_{i+1}, D) + \\ \lambda_U \sum_{q \in Q} \tilde{f}_U(q_i, q_{i+1}, D)$$

Separate MLMs for bigrams and unigrams give FSDM the flexibility to adjust the document scoring depending on the query type

MLM is a special case of FSDM, when  $\lambda_T = 1$ ,  $\lambda_O = 0$ ,  $\lambda_U = 0$

# FSDM RANKING FUNCTION

FSDM incorporates document structure and term dependencies with the following ranking function:

$$P_{\Lambda}(D|Q) \stackrel{\text{rank}}{=} \lambda_T \sum_{q \in Q} \tilde{f}_T(q_i, D) + \\ \lambda_O \sum_{q \in Q} \tilde{f}_O(q_i, q_{i+1}, D) + \\ \lambda_U \sum_{q \in Q} \tilde{f}_U(q_i, q_{i+1}, D)$$

Separate MLMs for bigrams and unigrams give FSDM the flexibility to adjust the document scoring depending on the query type

MLM is a special case of FSDM, when  $\lambda_T = 1$ ,  $\lambda_O = 0$ ,  $\lambda_U = 0$

# FSDM RANKING FUNCTION

FSDM incorporates document structure and term dependencies with the following ranking function:

$$P_{\Lambda}(D|Q) \stackrel{\text{rank}}{=} \lambda_T \sum_{q \in Q} \tilde{f}_T(q_i, D) + \\ \lambda_O \sum_{q \in Q} \tilde{f}_O(q_i, q_{i+1}, D) + \\ \lambda_U \sum_{q \in Q} \tilde{f}_U(q_i, q_{i+1}, D)$$

Separate MLMs for bigrams and unigrams give FSDM the flexibility to adjust the document scoring depending on the query type

MLM is a special case of FSDM, when  $\lambda_T = 1$ ,  $\lambda_O = 0$ ,  $\lambda_U = 0$

# FSDM RANKING FUNCTION

FSDM incorporates document structure and term dependencies with the following ranking function:

$$P_{\Lambda}(D|Q) \stackrel{\text{rank}}{=} \lambda_T \sum_{q \in Q} \tilde{f}_T(q_i, D) + \\ \lambda_O \sum_{q \in Q} \tilde{f}_O(q_i, q_{i+1}, D) + \\ \lambda_U \sum_{q \in Q} \tilde{f}_U(q_i, q_{i+1}, D)$$

Separate MLMs for bigrams and unigrams give FSDM the flexibility to adjust the document scoring depending on the query type

MLM is a special case of FSDM, when  $\lambda_T = 1$ ,  $\lambda_O = 0$ ,  $\lambda_U = 0$



## 5-FIELD ENTITY DOCUMENT

[ZHILTSOV, KOTOV ET AL., SIGIR'15]

Each entity is represented as a five-field document:

names

conventional names of the entities, such as the name of a person or the name of an organization

attributes

all entity properties, other than names

categories

classes or groups, to which the entity has been assigned

similar entity names

names of the entities that are very similar or identical to a given entity

related entity names

names of the entities that are part of the same RDF triple

# FSDM RANKING FUNCTION

Potential function for unigrams in case of FSDM:

$$\tilde{f}_T(q_i, D) = \log \sum_j w_j^T P(q_i | \theta_D^j) = \log \sum_j w_j^T \frac{tf_{q_i, D^j} + \mu_j \frac{cf_{q_i}^j}{|C_j|}}{|D^j| + \mu_j}$$

Example

apollo astronauts who walked on the moon

# FSDM RANKING FUNCTION

Potential function for unigrams in case of FSDM:

$$\tilde{f}_T(q_i, D) = \log \sum_j w_j^T P(q_i | \theta_D^j) = \log \sum_j w_j^T \frac{tf_{q_i, D^j} + \mu_j \frac{cf_{q_i}^j}{|C_j|}}{|D^j| + \mu_j}$$

Example

apollo astronauts who walked on the moon  
category

# FSDM RANKING FUNCTION

Potential function for unigrams in case of FSDM:

$$\tilde{f}_T(q_i, D) = \log \sum_j w_j^T P(q_i | \theta_D^j) = \log \sum_j w_j^T \frac{tf_{q_i, D^j} + \mu_j \frac{cf_{q_i}^j}{|C_j|}}{|D^j| + \mu_j}$$

Example

apollo astronauts who **walked on the moon**  
 category attribute

# PARAMETERS OF FSDM

Overall, FSDM has  $3 * F + 3$  free parameters:  $\langle w^T, w^O, w^U, \lambda \rangle$ .

## Properties of ranking function

1. Linearity with respect to  $\lambda$ .

We can apply any linear learning-to-rank algorithm to optimize the ranking function with respect to  $\lambda$ .

2. Linearity with respect to  $w$  of the arguments of monotonic  $\tilde{f}(\cdot)$  functions.

Optimization of the arguments as linear functions with respect to  $w$ , leads to optimization of each function  $\tilde{f}(\cdot)$ .

# OPTIMIZATION ALGORITHM

- 1:  $Q \leftarrow$  Training queries
- 2: **for**  $s \in \{T, O, U\}$  **do** // Optimize field weights of LMs independently
- 3:      $\lambda = e_s$
- 4:      $\hat{w}^s \leftarrow \text{CoordAsc}(Q, \lambda)$
- 5: **end for**
- 6:  $\hat{\lambda} \leftarrow \text{CoordAsc}(Q, \hat{w}_T, \hat{w}_O, \hat{w}_U)$  // Optimize  $\lambda$

The unit vectors  $e_T = (1, 0, 0)$ ,  $e_O = (0, 1, 0)$ ,  $e_U = (0, 0, 1)$  are the corresponding settings of the parameters  $\lambda$  in the formula of FSDM ranking function.

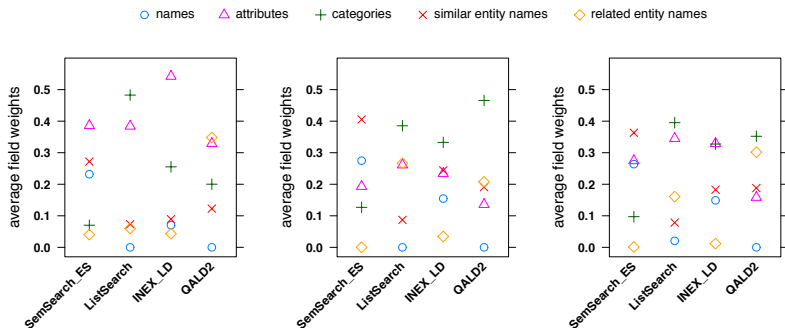
$\Rightarrow$  direct optimization w.r.t. target metric, e.g. MAP

## COLLECTION AND QUERY SETS

- ▶ DBpedia 3.7 was used as a knowledge
- ▶ Queries from Balog and Neumayer. A Test Collection for Entity Search in DBpedia, SIGIR'13.

Query set	Amount	Query types [Pound et al., 2010]
SemSearch ES	130	Entity
ListSearch	115	Type
INEX-LD	100	Entity, Type, Attribute, Relation
QALD-2	140	Entity, Type, Attribute, Relation

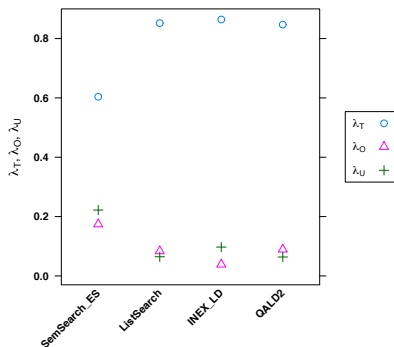
# TUNING FIELD WEIGHTS



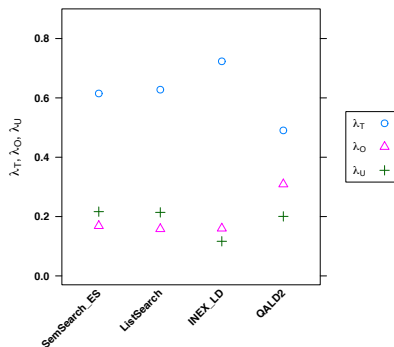
- ▶ *Attributes* field is consistently considered to be a very valuable for both unigrams and bigrams.
- ▶ The *names* field as well as the *similar entity names* field are highly important for queries aiming at finding named entities.
- ▶ Distinguishing *categories* from *related entity names* is particularly important for type queries.



# TUNING $\lambda$



(a) SDM



(b) FSDM

- Bigram matches are important for named entity queries.
- Transformation of SDM into FSDM increases the importance of bigram matches, which ultimately improves the retrieval performance

## EXPERIMENTAL RESULTS

Query set	Method	MAP	P@10	P@20	b-pref
SemSearch ES	MLM-CA	0.320	0.250	0.179	0.674
	SDM-CA	0.254*	0.202*	0.149*	0.671
	FSDM	<b>0.386<sub>†</sub>*</b>	<b>0.286<sub>†</sub>*</b>	<b>0.204<sub>†</sub>*</b>	<b>0.750<sub>†</sub>*</b>
ListSearch	MLM-CA	0.190	0.252	0.192	0.428
	SDM-CA	0.197	0.252	0.202	<b>0.471*</b>
	FSDM	<b>0.203</b>	<b>0.256</b>	<b>0.203</b>	0.466*
INEX-LD	MLM-CA	0.102	0.238	0.190	0.318
	SDM-CA	<b>0.117*</b>	0.258	0.199	0.335
	FSDM	0.111*	<b>0.263*</b>	<b>0.215<sub>†</sub>*</b>	<b>0.341*</b>
QALD-2	MLM-CA	0.152	0.103	0.084	0.373
	SDM-CA	0.184	0.106	0.090	0.465*
	FSDM	<b>0.195*</b>	<b>0.136<sub>†</sub>*</b>	<b>0.111*</b>	<b>0.466*</b>
All queries	MLM-CA	0.196	0.206	0.157	0.455
	SDM-CA	0.192	0.198	0.155	0.495*
	FSDM	<b>0.231<sub>†</sub>*</b>	<b>0.231<sub>†</sub>*</b>	<b>0.179<sub>†</sub>*</b>	<b>0.517<sub>†</sub>*</b>

# FSDM LIMITATION

In FSDM field weights are the same for all query concepts of the same type.

## Example

capitals in Europe which were host cities of summer Olympic games

# PARAMETRIC EXTENSION OF FSDM

$$w_{q_i,j}^T = \sum_k \alpha_{j,k}^U \phi_k(q_i, j)$$

- ▶  $\phi_k(q_i, j)$  is the the  $k$ -th feature value for unigram  $q_i$  in field  $j$ .
- ▶  $\alpha_{j,k}^U$  are feature weights that we learn.

$$\sum_j w_{q_i,j}^T = 1, w_{q_i,j}^T \geq 0, \alpha_{j,k}^U \geq 0, 0 \leq \phi_k(q_i, j) \leq 1$$

PFFDM is the same, but uses full dependence model.

# PARAMETRIC EXTENSION OF FSDM

$$w_{q_i,j}^T = \sum_k \alpha_{j,k}^U \phi_k(q_i, j)$$

- ▶  $\phi_k(q_i, j)$  is the the  $k$ -th feature value for unigram  $q_i$  in field  $j$ .
- ▶  $\alpha_{j,k}^U$  are feature weights that we learn.

$$\sum_j w_{q_i,j}^T = 1, w_{q_i,j}^T \geq 0, \alpha_{j,k}^U \geq 0, 0 \leq \phi_k(q_i, j) \leq 1$$

PFFDM is the same, but uses full dependence model.

# PARAMETRIC EXTENSION OF FSDM

$$w_{q_i,j}^T = \sum_k \alpha_{j,k}^U \phi_k(q_i, j)$$

- ▶  $\phi_k(q_i, j)$  is the the  $k$ -th feature value for unigram  $q_i$  in field  $j$ .
- ▶  $\alpha_{j,k}^U$  are feature weights that we learn.

$$\sum_j w_{q_i,j}^T = 1, w_{q_i,j}^T \geq 0, \alpha_{j,k}^U \geq 0, 0 \leq \phi_k(q_i, j) \leq 1$$

PFDDM is the same, but uses full dependence model.

# PARAMETRIC EXTENSION OF FSDM

$$w_{q_i,j}^T = \sum_k \alpha_{j,k}^U \phi_k(q_i, j)$$

- ▶  $\phi_k(q_i, j)$  is the the  $k$ -th feature value for unigram  $q_i$  in field  $j$ .
- ▶  $\alpha_{j,k}^U$  are feature weights that we learn.

$$\sum_j w_{q_i,j}^T = 1, w_{q_i,j}^T \geq 0, \alpha_{j,k}^U \geq 0, 0 \leq \phi_k(q_i, j) \leq 1$$

PPFDM is the same, but uses full dependence model.

# PARAMETRIC EXTENSION OF FSDM

$$w_{q_i,j}^T = \sum_k \alpha_{j,k}^U \phi_k(q_i, j)$$

- ▶  $\phi_k(q_i, j)$  is the the  $k$ -th feature value for unigram  $q_i$  in field  $j$ .
- ▶  $\alpha_{j,k}^U$  are feature weights that we learn.

$$\sum_j w_{q_i,j}^T = 1, w_{q_i,j}^T \geq 0, \alpha_{j,k}^U \geq 0, 0 \leq \phi_k(q_i, j) \leq 1$$

PFFDM is the same, but uses full dependence model.



# FEATURES

Source	Feature	Description	CT
Collection statistics	$FP(\kappa, j)$	Posterior probability $P(E_j w)$ .	UG BG
	$TS(\kappa, j)$	Top SDM score on $j$ -th field when $\kappa$ is used as a query.	BG
Stanford POS Tagger	$NNP(\kappa)$	Is concept $\kappa$ a proper noun?	UG
	$NNS(\kappa)$	Is $\kappa$ a plural non-proper noun?	UG BG
	$JJS(\kappa)$	Is $\kappa$ a superlative adjective?	UG
Stanford Parser	$NPP(\kappa)$	Is $\kappa$ part of a noun phrase?	BG
	$NNO(\kappa)$	Is $\kappa$ the only singular non-proper noun in a noun phrase?	UG
	$INT$	Intercept feature (= 1).	UG BG

## FEATURES

Source	Feature	Description	CT
Collection statistics	$FP(\kappa, j)$	Posterior probability $P(E_j w)$ .	UG BG
	$TS(\kappa, j)$	Top SDM score on $j$ -th field when $\kappa$ is used as a query.	BG
Stanford POS Tagger	$NNP(\kappa)$	Is concept $\kappa$ a proper noun?	UG
	$NNS(\kappa)$	Is $\kappa$ a plural non-proper noun?	UG BG
	$JJS(\kappa)$	Is $\kappa$ a superlative adjective?	UG
Stanford Parser	$NPP(\kappa)$	Is $\kappa$ part of a noun phrase?	BG
	$NNO(\kappa)$	Is $\kappa$ the only singular non-proper noun in a noun phrase?	UG
	$INT$	Intercept feature (= 1).	UG BG

# PARAMETERS OF PFSDM

Both PFSDM and PFFDM have  $F * U + F * B + 3$  free parameters:  
 $\langle \hat{\alpha}^U, \hat{\alpha}^B, \hat{\lambda} \rangle$ .

We perform direct optimization w.r.t. target metric (e.g. MAP) using coordinate ascent.

# COLLECTIONS

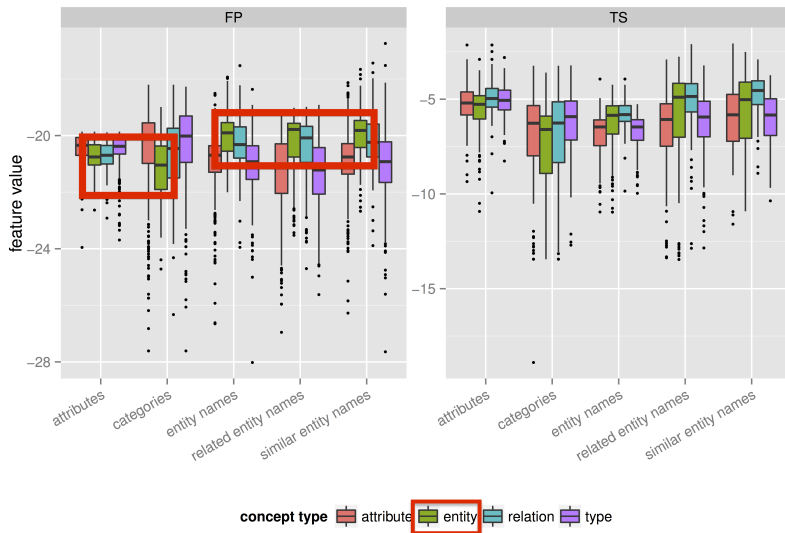
## 1. DBPedia 3.7

- ▶ Structured version of on-line encyclopedia Wikipedia
- ▶ Provides the descriptions of over 3.5 million entities belonging to 320 classes

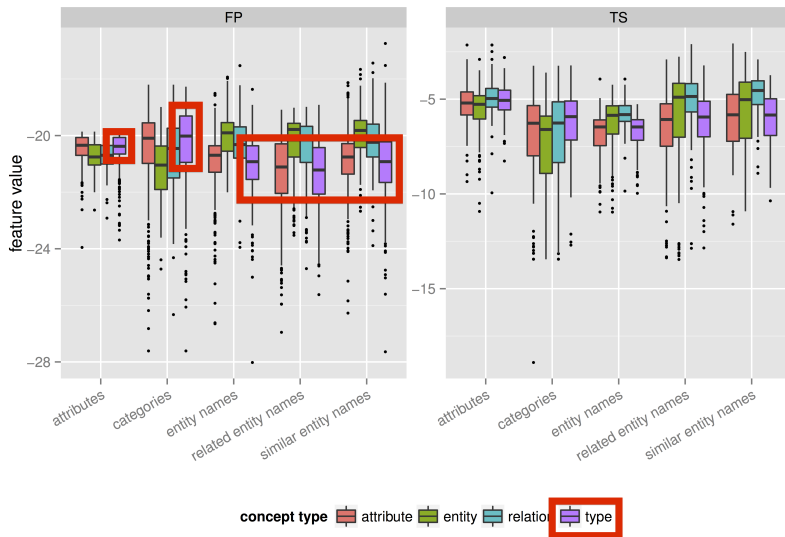
## 2. BTC-2009

- ▶ Contains entities from multiple knowledge bases.
- ▶ Consists of 1.14 billion RDF triples.

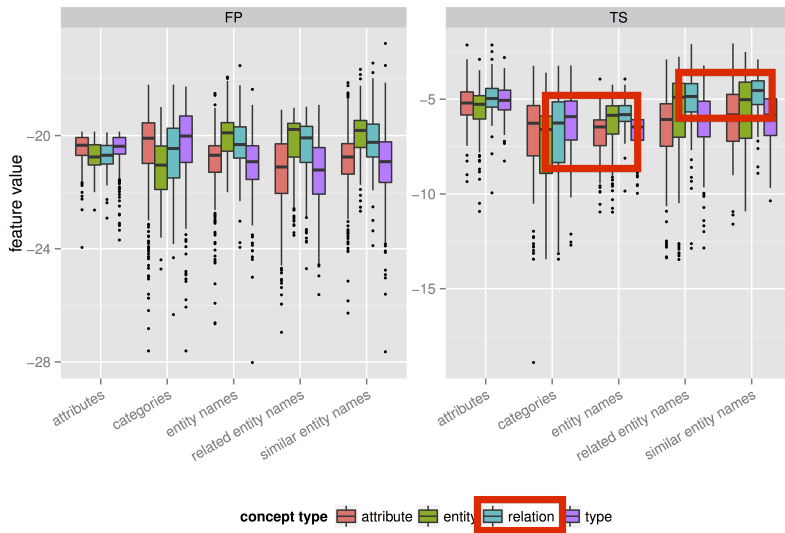
# REAL-VALUED FEATURES ANALYSIS



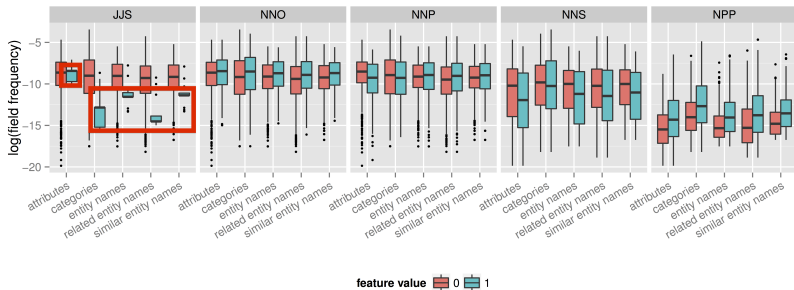
# REAL-VALUED FEATURES ANALYSIS



# REAL-VALUED FEATURES ANALYSIS

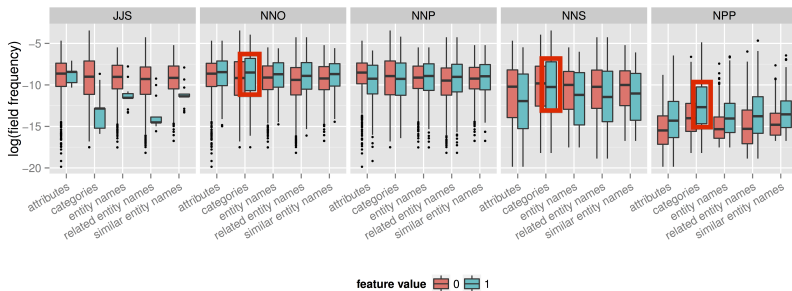


# NLP-BASED FEATURES ANALYSIS

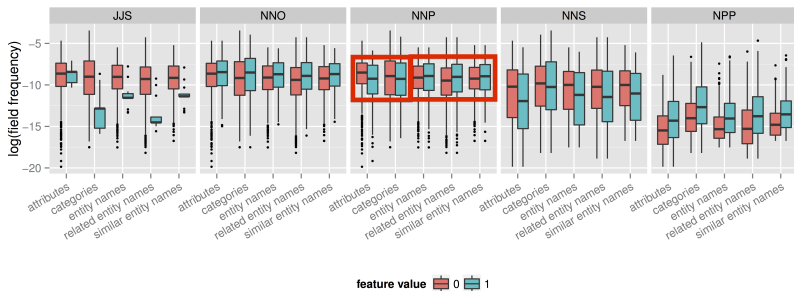




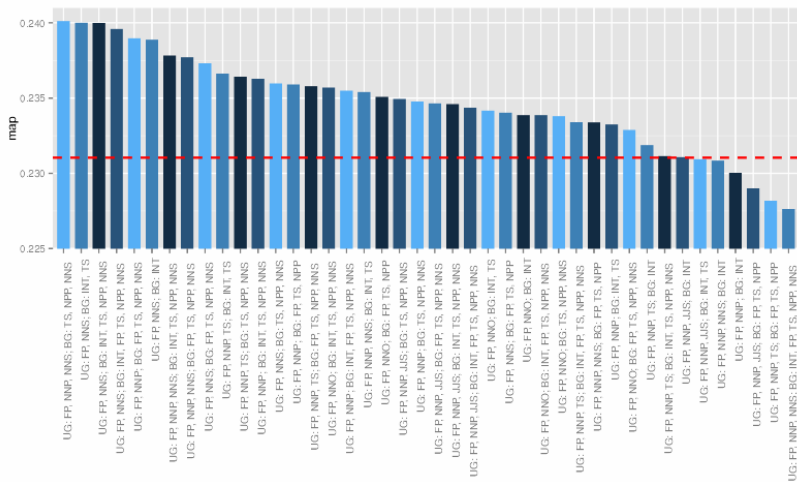
# NLP-BASED FEATURES ANALYSIS



# NLP-BASED FEATURES ANALYSIS



# FEATURE EFFECTIVENESS



## DBPEDIA RESULTS (USING BEST FEATURES COMBINATION)

Query set	Method	MAP	P@10	P@20	b-pref
SemSearch ES	PRMS	0.230	0.177	0.549	0.317
	FSDM	0.386	0.286	0.737	0.476
	PFSDM	0.394*	0.286*	0.757*	0.494* <sub>†</sub>
	FFDM	0.389*	0.286*	0.734*	0.479*
	PFFDM	0.380*	0.286*	0.739*	0.477*
ListSearch	PRMS	0.111	0.154	0.355	0.176
	FSDM	0.203	0.256	0.447	0.274
	PFSDM	0.201*	0.253*	0.443*	0.278*
	FFDM	0.226* <sub>†</sub>	0.282* <sub>†</sub>	0.499* <sub>†</sub>	0.313* <sub>†</sub>
	PFFDM	0.228* <sub>†</sub>	0.286* <sub>†</sub>	0.487* <sub>†</sub>	0.302* <sub>†</sub>
INEX-LD	PRMS	0.064	0.145	0.409	0.216
	FSDM	0.111	0.263	0.546	0.322
	PFSDM	0.116*	0.259*	0.579*	0.341*
	FFDM	0.122* <sub>†</sub>	0.273*	0.560*	0.345* <sub>†</sub>
	PFFDM	0.121* <sub>†</sub>	0.274*	0.556*	0.343*
QALD-2	PRMS	0.120	0.079	0.188	0.147
	FSDM	0.195	0.136	0.283	0.229
	PFSDM	0.218* <sub>†</sub>	0.140*	0.308*	0.253* <sub>†</sub>
	FFDM	0.200*	0.139*	0.292*	0.237*
	PFFDM	0.219* <sub>†</sub>	0.147*	0.310*	0.267* <sub>†</sub>
All queries	PRMS	0.136	0.136	0.370	0.214
	FSDM	0.231	0.231	0.498	0.325
	PFSDM	0.240* <sub>†</sub>	0.231*	0.516* <sub>†</sub>	0.342* <sub>†</sub>
	FFDM	0.241* <sub>†</sub>	0.240* <sub>†</sub>	0.515* <sub>†</sub>	0.342* <sub>†</sub>
	PFFDM	0.244* <sub>†</sub>	0.244* <sub>†</sub>	0.518* <sub>†</sub>	0.347* <sub>†</sub>

# DBPEDIA RESULTS (USING BEST FEATURES COMBINATION)

Query set	Method	MAP	P@10	P@20	b-pref
SemSearch ES	PRMS	0.230	0.177	0.549	0.317
	FSDM	0.386	0.286	0.737	0.476
	PFSDM	0.394*	0.286*	0.757*	0.494* <sub>†</sub>
	FFDM	0.389*	0.286*	0.734*	0.479*
	PFFDM	0.380*	0.286*	0.739*	0.477*
ListSearch	PRMS	0.111	0.154	0.355	0.176
	FSDM	0.203	0.256	0.447	0.274
	PFSDM	0.201*	0.253*	0.443*	0.278*
	FFDM	0.226* <sub>†</sub>	0.282* <sub>†</sub>	0.499* <sub>†</sub>	0.313* <sub>†</sub>
	PFFDM	0.228* <sub>†</sub>	0.286* <sub>†</sub>	0.487* <sub>†</sub>	0.302* <sub>†</sub>
INEX-LD	PRMS	0.064	0.145	0.409	0.216
	FSDM	0.111	0.263	0.546	0.322
	PFSDM	0.116*	0.259*	0.579*	0.341*
	FFDM	0.122* <sub>†</sub>	0.273*	0.560*	0.345* <sub>†</sub>
	PFFDM	0.121* <sub>†</sub>	0.274*	0.556*	0.343*
QALD-2	PRMS	0.120	0.079	0.188	0.147
	FSDM	0.195	0.136	0.283	0.229
	PFSDM	0.218* <sub>†</sub>	0.140*	0.308*	0.253* <sub>†</sub>
	FFDM	0.200*	0.139*	0.292*	0.237*
	PFFDM	0.219* <sub>†</sub>	0.147*	0.310*	0.267* <sub>†</sub>
All queries	PRMS	0.136	0.136	0.370	0.214
	FSDM	0.231	0.231	0.498	0.325
	PFSDM	0.240* <sub>†</sub>	0.231*	0.516* <sub>†</sub>	0.342* <sub>†</sub>
	FFDM	0.241* <sub>†</sub>	0.240* <sub>†</sub>	0.515* <sub>†</sub>	0.342* <sub>†</sub>
	PFFDM	0.244* <sub>†</sub>	0.244* <sub>†</sub>	0.518* <sub>†</sub>	0.347* <sub>†</sub>

# DBPEDIA RESULTS (USING BEST FEATURES COMBINATION)

Query set	Method	MAP	P@10	P@20	b-pref
SemSearch ES	PRMS	0.230	0.177	0.549	0.317
	FSDM	0.386	0.286	0.737	0.476
	PFSDM	0.394*	0.286*	0.757*	0.494* <sub>†</sub>
	FFDM	0.389*	0.286*	0.734*	0.479*
	PFFDM	0.380*	0.286*	0.739*	0.477*
ListSearch	PRMS	0.111	0.154	0.355	0.176
	FSDM	0.203	0.256	0.447	0.274
	PFSDM	0.201*	0.253*	0.443*	0.278*
	FFDM	0.226* <sub>†</sub>	0.282* <sub>†</sub>	0.499* <sub>†</sub>	0.313* <sub>†</sub>
	PFFDM	0.228* <sub>†</sub>	0.286* <sub>†</sub>	0.487* <sub>†</sub>	0.302* <sub>†</sub>
INEX-LD	PRMS	0.064	0.145	0.409	0.216
	FSDM	0.111	0.263	0.546	0.322
	PFSDM	0.116*	0.259*	0.579*	0.341*
	FFDM	0.122* <sub>†</sub>	0.273*	0.560*	0.345* <sub>†</sub>
	PFFDM	0.121* <sub>†</sub>	0.274*	0.556*	0.343*
QALD-2	PRMS	0.120	0.079	0.188	0.147
	FSDM	0.195	0.136	0.283	0.229
	PFSDM	0.218* <sub>†</sub>	0.140*	0.308*	0.253* <sub>†</sub>
	FFDM	0.200*	0.139*	0.292*	0.237*
	PFFDM	0.219* <sub>†</sub>	0.147*	0.310*	0.267* <sub>†</sub>
All queries	PRMS	0.136	0.136	0.370	0.214
	FSDM	0.231	0.231	0.498	0.325
	PFSDM	0.240* <sub>†</sub>	0.231*	0.516* <sub>†</sub>	0.342* <sub>†</sub>
	FFDM	0.241* <sub>†</sub>	0.240* <sub>†</sub>	0.515* <sub>†</sub>	0.342* <sub>†</sub>
	PFFDM	0.244* <sub>†</sub>	0.244* <sub>†</sub>	0.518* <sub>†</sub>	0.347* <sub>†</sub>

# DBPEDIA RESULTS (USING BEST FEATURES COMBINATION)

Query set	Method	MAP	P@10	P@20	b-pref
SemSearch ES	PRMS	0.230	0.177	0.549	0.317
	FSDM	0.386	0.286	0.737	0.476
	PFSDM	0.394*	0.286*	0.757*	0.494* <sub>†</sub>
	FFDM	0.389*	0.286*	0.734*	0.479*
	PFFDM	0.380*	0.286*	0.739*	0.477*
ListSearch	PRMS	0.111	0.154	0.355	0.176
	FSDM	0.203	0.256	0.447	0.274
	PFSDM	0.201*	0.253*	0.443*	0.278*
	FFDM	0.226* <sub>†</sub>	0.282* <sub>†</sub>	0.499* <sub>†</sub>	0.313* <sub>†</sub>
	PFFDM	0.228* <sub>†</sub>	0.286* <sub>†</sub>	0.487* <sub>†</sub>	0.302* <sub>†</sub>
INEX-LD	PRMS	0.064	0.145	0.409	0.216
	FSDM	0.111	0.263	0.546	0.322
	PFSDM	0.116*	0.259*	0.579*	0.341*
	FFDM	0.122* <sub>†</sub>	0.273*	0.560*	0.345* <sub>†</sub>
	PFFDM	0.121* <sub>†</sub>	0.274*	0.556*	0.343*
QALD-2	PRMS	0.120	0.079	0.188	0.147
	FSDM	0.195	0.136	0.283	0.229
	PFSDM	0.218* <sub>†</sub>	0.140*	0.308*	0.253* <sub>†</sub>
	FFDM	0.200* <sub>†</sub>	0.139*	0.292*	0.237*
	PFFDM	0.219* <sub>†</sub>	0.147*	0.310*	0.267* <sub>†</sub>
All queries	PRMS	0.136	0.136	0.370	0.214
	FSDM	0.231	0.231	0.498	0.325
	PFSDM	0.240* <sub>†</sub>	0.231*	0.516* <sub>†</sub>	0.342* <sub>†</sub>
	FFDM	0.241* <sub>†</sub>	0.240* <sub>†</sub>	0.515* <sub>†</sub>	0.342* <sub>†</sub>
	PFFDM	0.244* <sub>†</sub>	0.244* <sub>†</sub>	0.518* <sub>†</sub>	0.347* <sub>†</sub>

# DBPEDIA RESULTS (USING BEST FEATURES COMBINATION)

Query set	Method	MAP	P@10	P@20	b-pref
SemSearch ES	PRMS	0.230	0.177	0.549	0.317
	FSDM	0.386	0.286	0.737	0.476
	PFSDM	0.394*	0.286*	0.757*	0.494* <sub>†</sub>
	FFDM	0.389*	0.286*	0.734*	0.479*
	PFFDM	0.380*	0.286*	0.739*	0.477*
ListSearch	PRMS	0.111	0.154	0.355	0.176
	FSDM	0.203	0.256	0.447	0.274
	PFSDM	0.201*	0.253*	0.443*	0.278*
	FFDM	0.226* <sub>†</sub>	0.282* <sub>†</sub>	0.499* <sub>†</sub>	0.313* <sub>†</sub>
	PFFDM	0.228* <sub>†</sub>	0.286* <sub>†</sub>	0.487* <sub>†</sub>	0.302* <sub>†</sub>
INEX-LD	PRMS	0.064	0.145	0.409	0.216
	FSDM	0.111	0.263	0.546	0.322
	PFSDM	0.116*	0.259*	0.579*	0.341*
	FFDM	0.122* <sub>†</sub>	0.273*	0.560* <sub>†</sub>	0.345* <sub>†</sub>
	PFFDM	0.121* <sub>†</sub>	0.274*	0.556*	0.343*
QALD-2	PRMS	0.120	0.079	0.188	0.147
	FSDM	0.195	0.136	0.283	0.229
	PFSDM	0.218* <sub>†</sub>	0.140*	0.308*	0.253* <sub>†</sub>
	FFDM	0.200*	0.139*	0.292*	0.237*
	PFFDM	0.219* <sub>†</sub>	0.147*	0.310*	0.267* <sub>†</sub>
All queries	PRMS	0.136	0.136	0.370	0.214
	FSDM	0.231	0.231	0.498	0.325
	PFSDM	0.240* <sub>†</sub>	0.231*	0.516* <sub>†</sub>	0.342* <sub>†</sub>
	FFDM	0.241* <sub>†</sub>	0.240* <sub>†</sub>	0.515* <sub>†</sub>	0.342* <sub>†</sub>
	PFFDM	0.244* <sub>†</sub>	0.244* <sub>†</sub>	0.518* <sub>†</sub>	0.347* <sub>†</sub>



# BTC2009 RESULTS

<b>Method</b>	<b>MAP</b>	<b>P@10</b>	<b>P@20</b>	<b>b-pref</b>
PRMS	0.098	0.198	0.545	0.269
FSDM	0.171	0.323	0.631	0.358
PFSDM	0.182* †	0.335*	0.657* †	0.371*
FFDM	0.180* †	0.330* †	0.647* †	0.373* †
PFFDM	0.187*	0.342* †	0.650*	0.377*

# FUTURE WORK

We hypothesize that FSDM and PFDSM can be effective in other structured information retrieval scenarios, such as product and social graph search, and leave verification of this hypothesis to industry or research community.

# OVERVIEW

Entities and Entity Retrieval

Knowledge Graphs

Entity Representation

Entity Retrieval

**Conclusion**

- ▶ Code and runs are available at:

`github.com/teanalab`

- ▶ Send me an email at `kotov@wayne.edu`, if you have any questions about this tutorial or would like to collaborate on future projects.

# Thank you!