

Making the Semantic Web Easier to Use

Tim Finin

University of Maryland, Baltimore County

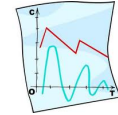
Joint work with Lushan Han, Varish Mulwad, Anupam Joshi

Presented at the Workshop on Data Engineering Meets the Semantic Web of the IEEE International Conference on Data Engineering, 1 April 2012

UMBC
ebiquity

<http://ebiq.org/r/339>

Overview



- Linked Open Data 101
- Two ongoing UMBC dissertations
 - Varish Mulwad, Generating linked data from tables
 - Lushan Han, Querying linked data with a quasi-NL interface

2/49

Linked Open Data (LOD)



- Linked **data** is just RDF data, typically just the instances (ABox), not schema (TBox)
- RDF data is a graph of triples
 - URI URI string
dbr:Barack_Obama dbo:spouse "Michelle Obama"
 - URI URI URI
dbr:Barack_Obama dbo:spouse dbpedia:Michelle_Obama
- Best **linked** data practice prefers the 2nd pattern, using nodes rather than strings for "entities"
- Liked **open** data is just linked data freely accessible on the Web along with any required ontologies

3/49

Semantic Web

Use Semantic Web Technology to publish shared data & knowledge


Semantic web technologies allow machines to share data and knowledge using common web language and protocols.

~ 1997

Semantic Web beginning

Semantic Web => Linked Open Data

Use Semantic Web Technology to publish shared data & knowledge




2007

Data is inter-linked to support integration and fusion of knowledge

LOD beginning

Semantic Web => Linked Open Data

Use Semantic Web Technology to publish shared data & knowledge



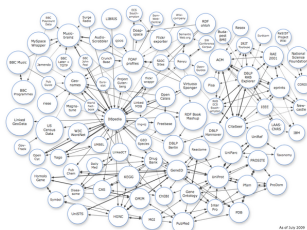
2008

Data is inter-linked to support integration and fusion of knowledge

LOD growing

Semantic Web => Linked Open Data

Use Semantic Web Technology to publish shared data & knowledge



2009

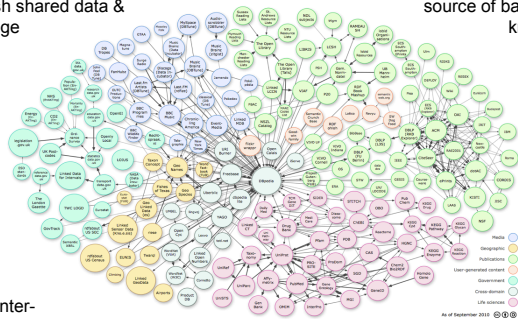
Data is inter-linked to support integration and fusion of knowledge

... and growing

Linked Open Data

Use Semantic Web Technology to publish shared data & knowledge

LOD is the new Cyc: a common source of background knowledge



2010

Data is inter-linked to support integration and fusion of knowledge

...growing faster

Linked Open Data

Use Semantic Web Technology to publish shared data & knowledge **LOD is the new Cyc: a common source of background knowledge**

Data is inter-linked to support integration and fusion of knowledge

2011: 31B facts in 295 datasets interlinked by 504M assertions on ckan.net

Exploiting LOD not (yet) Easy

- Publishing or using LOD data has inherent difficulties for the potential user
 - It's difficult to explore LOD data and to query it for answers
 - It's challenging to publish data using appropriate LOD vocabularies & link it to existing data
- Problem: $O(10^4)$ schema terms, $O(10^{11})$ instances
- I'll describe two ongoing research projects that are addressing these problems

10/49

Generating Linked Data by Inferring the Semantics of Tables

Research with Varish Mulwad

<http://ebiq.org/j/96>

Goal: Table => LOD*

Name	Team	Position	Height
Michael Jordan	Chicago	Shooting guard	1.98
Allen Iverson	Philadelphia	Point guard	1.83
Yao Ming	Houston	Center	2.29
Tim Duncan	San Antonio	Power forward	2.11

* DBpedia

12/49

Goal: Table => LOD*

Name	Team	Position	Height
Michael Jordan	Chicago	Shooting guard	1.98
Allen Iverson	Philadelphia	Point guard	1.83
Yao Ming	Houston		
Tim Duncan	San Antonio		

RDF Linked Data

@prefix dbpedia: <http://dbpedia.org/resource/> .
 @prefix dbo: <http://dbpedia.org/ontology/> .
 @prefix yago: <http://dbpedia.org/class/yago/> .

"Name"@en is rdfs:label of dbo:BasketballPlayer .
 "Team"@en is rdfs:label of yago:NationalBasketballAssociationTeams .

"Michael Jordan"@en is rdfs:label of dbpedia:Michael Jordan .
 dbpedia:Michael Jordan a dbo:BasketballPlayer .

"Chicago Bulls"@en is rdfs:label of dbpedia:Chicago Bulls .
 dbpedia:Chicago Bulls a yago:NationalBasketballAssociationTeams .

All this in a completely automated way

* DBpedia 13/49

Tables are everywhere !! ... yet ...

The web – **154 million** high quality relational tables

When Assessed	Major	Non-major	Minor	Total
n	416	1,027	2,022	3,465
Age (mean)	59.62	57.42	64.22	59.42
Diabetes duration (years)	22.23	17.02	14.12	18.42
HbA1c (mean)	8.1	8.1	8.1	8.1
HbA1c (SD)	1.5	1.5	1.5	1.5
Timing (years since diagnosis)	100.1	107.3	101.1	102.8

	Omeprazole Oral Suspension (n = 178)	Intravenous Cimetidine (n = 181)	Confidence Interval for the Difference in Rates, %
Clinically significant bleeding, n (%)	7 (3.9)	10 (5.5)	-100.0, 2.8*
Any overt bleeding, n (%)	34 (19.1)	58 (32.0)	-21.9, -4.0**
Inadequate pH control, n (%)	32 (18.0)	105 (58.0)	

14/49

Evidence-based medicine

Evidence-based medicine judges the efficacy of treatments or tests by meta-analyses of clinical trials. Key information is often found in tables in articles

	Omeprazole Oral Suspension (n = 178)	Intravenous Cimetidine (n = 181)	Confidence Interval for the Difference in Rates, %
Clinically significant bleeding, n (%)	7 (3.9)	10 (5.5)	-100.0, 2.8*
Any overt bleeding, n (%)	34 (19.1)	58 (32.0)	-21.9, -4.0**
Inadequate pH control, n (%)	32 (18.0)	105 (58.0)	

of Clinical trials published in 2008

Trial Day	Omeprazole Oral Suspension, %	Intravenous Cimetidine, %	P Value
1	2.4 (4/166)	11.5 (20/174)	<.01
2	4.6 (13/279)	16.3 (30/185)	<.01
3	2.8 (4/143)	17.8 (28/157)	<.01
4	4.0 (5/124)	13.1 (19/145)	<.01
5	2.8 (3/109)	15.5 (16/103)	<.01
6	2.2 (2/99)	20.5 (19/92)	<.01
7	1.4 (1/72)	17.9 (14/78)	<.01
8	3.0 (2/66)	24.3 (17/70)	<.01
9	3.8 (2/53)	32.1 (19/59)	<.01
10	4.7 (2/43)	33.3 (17/51)	<.01
11	5.0 (2/40)	39.4 (14/44)	<.01
12	6.0 (2/33)	25.6 (10/39)	<.01
13	6.0 (2/33)	27.3 (10/37)	<.01
14	3.7 (1/27)	28.6 (9/31)	.02

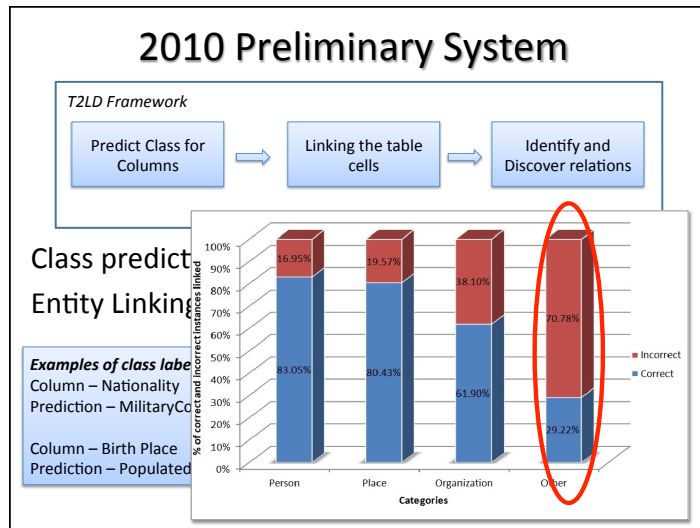
Characteristic	417 Cases	1998 Control Subjects	P
Age, mean (SD), years	70.9 (11.2)	69.0 (9.6)	.01
White, %	6.1	23.7 (11)	<.001
Time enrolled in (GEE) mean (SD) years	22.4 (12.7)	36.5	<.001
Postmenopausal hormone therapy, %	28.1 (9.9)	27.8 (9.3)	<.001
Body mass index, mean (SD) kg/m ²	21.2	2.2	<.001
Body mass index in prior 3 months, %	5.2	12.2	<.001
Hospitalization in prior 3 months, %	35.6	19.8	<.001
Major fracture in prior 3 months, %	31.5	0.1	
Vascular diseases ¹ , %			
Myocardial infarction			
Stroke			
Vascular disease ² , %			
Myocardial infarction			
Stroke			
Vascular disease ³ , %			
Myocardial infarction			
Stroke			

Figure: Evidence-Based Medicine - the Essential Role of Systematic Review 15/49

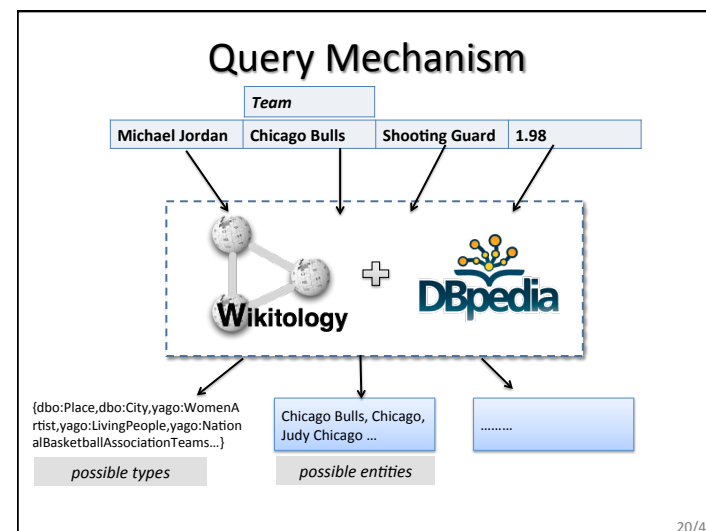
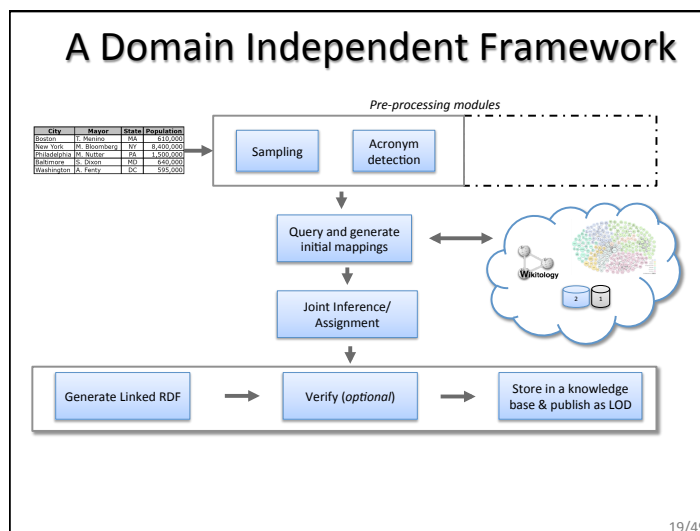
~ 400,000 datasets 🤖

~ < 1 % in RDF 😊

16/49



- ### Sources of Errors
-
- The *sequential* approach let errors percolate from one phase to the next
 - The system was biased toward predicting overly general classes over more appropriate specific ones
 - Heuristics largely drive the system
 - Although we consider multiple sources of evidence, we did not joint assignment
- 18/49



Ranking the candidates

- $C_i = \text{"State"} ; L_{C_i} = \text{AdministrativeRegion}$
 - String in column header
 - Class from an ontology
- $f_1 = [\text{Levenshtein distance}(C_i, L_{C_i}), \text{Dice Score}(C_i, L_{C_i}), \text{Semantic Similarity}(C_i, L_{C_i}), \text{InformationGain}(L_{C_i})]$
 - String similarity metrics
- $\psi_1 = \exp(w_1^T f_1(C_i, L_{C_i}))$

21/49

Ranking the candidates

- $R_{ij} = \text{"Baltimore"} ; E_{ij} = \text{Baltimore_Maryland}$
 - String in table cell
 - Entity from the knowledge base (KB)
- $f_2 = [\text{Levenshtein distance}(R_{ij}, E_{ij}), \text{Dice Score}(R_{ij}, E_{ij}), \text{PageRank}(E_{ij}), \text{KBScore}(E_{ij}), \text{PageLength}(E_{ij})]$
 - String similarity metrics
 - Popularity metrics
- $\psi_2 = \exp(w_2^T f_2(R_{ij}, E_{ij}))$

22/49

Joint Inference over evidence in a table

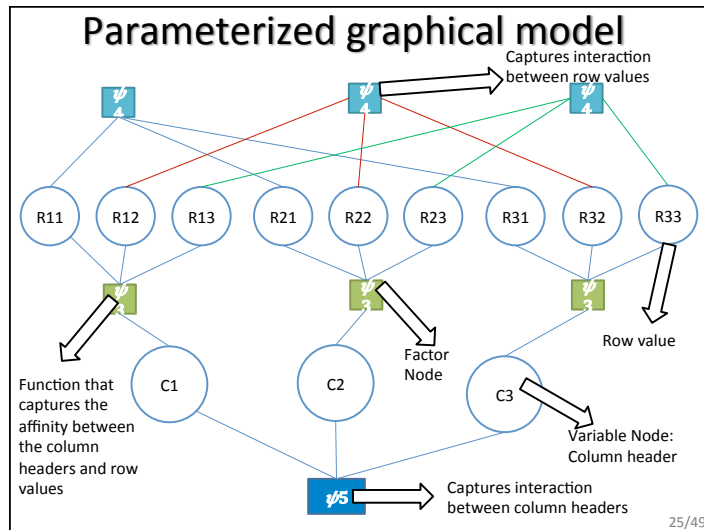
✓ Probabilistic Graphical Models

23/49

A graphical model for tables

Joint inference over evidence in a table

24/49



Challenge: Interpreting Literals

Many columns have literals, e.g., numbers

Population	
Population?	690,000
Profit in \$K ?	345,000
	510,020
	120,000

Age	
Age in years?	75
Percent?	65
	50
	25

- Predict properties based on cell values
- Cyc had hand coded rules: *humans don't live past 120*
- We extract *value distributions* from LOD resources
 - Differ for subclasses: *age of people vs. political leaders vs. athletes*
 - Represent as *measurements*: value + units
- Metric: possibility/probability of values given distribution

26/49

Other Challenges

- Using table *captions* and other text is associated documents to provide context
- **Size** of some data.gov tables (> 400K rows!) makes using full graphical model impractical
 - Sample table and run model on the subset
- Achieving acceptable accuracy may require **human input**
 - 100% accuracy unattainable automatically
 - How best to let humans offer advice and/or correct interpretations?

27/49

PMI as an association measure

We use **pointwise mutual information** (pmi) to measure the association between two RDF resources (nodes)

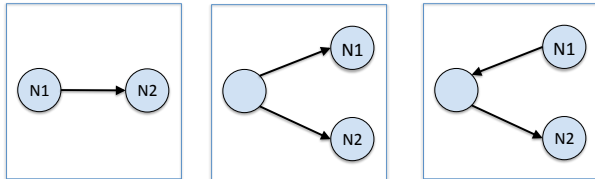
$$pmi(x; y) \equiv \log \frac{p(x, y)}{p(x)p(y)} = \log \frac{p(x|y)}{p(x)} = \log \frac{p(y|x)}{p(y)}$$

pmi is used for word association by comparing how often two words occur together in text to their expected co-occurrence if independent

28/49

PMI for RDF instances

- For text, the co-occurrence context is usually a window of some number of words (e.g, 50)
- For RDF instances, we count three graph patterns as instances of the co-occurrence of N1 and N2



- Other graph patterns can be added, but we've not evaluated their utility or cost to compute.

29/49

PMI for RDF types

- We also want to measure the association strength between *RDF types*, e.g., a `dbo:Actor` associated with a `dbo:Film` vs. a `dbo:Place`
- We can also measure the association of an RDF property and types, e.g. `dbo:author` used with a `dbo:Film` vs. a `dbo:Book`
- Such simple statistics can be efficiently computed for large RDF collections in parallel

PREFIX `dbo:` <<http://dbpedia.org/ontology/>>

30/49

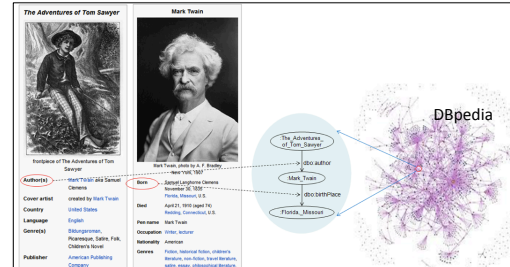
GoRelations: Intuitive Query System for Linked Data

Research with Lushan Han

<http://ebiq.org/j/93>

Dbpedia is the Stereotypical LOD

- DBpedia is an important example of Linked Open Data
 - Extracts structured data from Infoboxes in Wikipedia
 - Stores in RDF using custom ontologies Yago terms
- The major integration point for the entire LOD cloud
- Explorable as HTML, but harder to query in SPARQL



32/49

Browsing DBpedia's Mark Twain

Querying LOD is Much Harder

- Querying DBpedia requires a lot of a user
 - Understand the RDF model
 - Master SPARQL, a formal query language
 - Understand ontology terms: 320 classes & 1600 properties !
 - Know instance URIs (>1M entities !)
 - Term heterogeneity (Place vs. PopulatedPlace)
- Querying large LOD sets overwhelming
- Natural language query systems still a research goal

34/49

Goal

- Allow a user with a basic understanding of RDF to query DBpedia and ultimately distributed LOD collections
 - To explore what data is in the system
 - To get answers to question
 - To create SPARQL queries for reuse or adaptation
- Desiderata
 - Easy to learn and to use
 - Good accuracy (e.g., precision and recall)
 - Fast

35/49

Key Idea

Structured keyword queries

Reduce problem complexity by:

- User enters a *simple graph*, and
- Annotates the nodes and arcs with *words and phrases*

36/49

Structured Keyword Queries



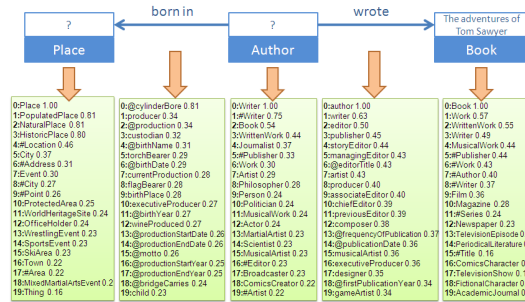
- Nodes denote entities and links binary relations
- Entities described by two unrestricted terms: *name* or value and *type* or concept
- Result entities marked with ? and those not with *
- A compromise between a natural language Q&A system and SPARQL
 - Users provide compositional structure of the question
 - Free to use their own terms in annotating the structure

37/49

Translation – Step One

finding semantically similar ontology terms

For each concept or relation in the graph, generate the *k* most semantically similar candidate ontology classes or properties

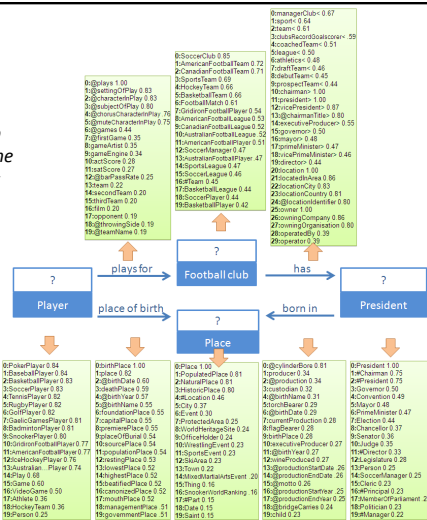


Similarity metric is distributional similarity, LSA, and WordNet.

38/49

Another Example

Football players who were born in the same place as their team's president



39/49

Translation – Step Two

disambiguation algorithm

- To assemble the best interpretation we rely on *statistics of the data*
- Primary measure is *pointwise mutual information* (PMI) between RDF terms in the LOD collection
 - This measures the degree to which two RDF terms occur together in the knowledge base
- In a reasonable interpretation, *ontology terms associate* in the way that their corresponding *user terms* connect in the structured keyword query

40/49

Translation – Step Two disambiguation algorithm

Three aspects are combined to derive an *overall goodness measure* for each candidate interpretation

Joint disambiguation

$$\operatorname{argmax}_{p_1 \dots p_m, c_1 \dots c_n \in H} \operatorname{goodness}(G) = \operatorname{argmax}_{p_1 \dots p_m, c_1 \dots c_n \in H} \sum_{i=1}^m \operatorname{goodness}(L_i) \quad (1)$$

Resolving direction

If $[\overline{\operatorname{PMI}}(c(O_i), p(R_i)) + \overline{\operatorname{PMI}}(p(R_i), c(S_i))] - [\overline{\operatorname{PMI}}(c(S_i), p(R_i)) + \overline{\operatorname{PMI}}(p(R_i), c(O_i))] > \alpha$

Then $S_i' = O_i, O_i' = S_i$

Else $S_i' = S_i, O_i' = O_i$ (2)

Link reason-ability

$$\operatorname{goodness}(L_i) = \max(\overline{\operatorname{PMI}}(c(S_i'), p(R_i)) \cdot \operatorname{sim}(S_i', c(S_i')) \cdot \operatorname{sim}(R_i, p(R_i)) + \overline{\operatorname{PMI}}(p(R_i), c(O_i')) \cdot \operatorname{sim}(O_i', c(O_i')) \cdot \operatorname{sim}(R_i, p(R_i)), \beta) + \operatorname{PMI}(c(S_i'), c(O_i')) \cdot \operatorname{sim}(S_i', c(S_i')) \cdot \operatorname{sim}(O_i', c(O_i')) \quad (3)$$

41/49

Example of Translation result

Concepts: Place => Place, Author => Writer, Book => Book
Properties: born in => birthPlace, wrote => author (inverse direction)

← born in

?

Place

← wrote →

?

Author


→

The adventures of Tom Sawyer

Book

<p>0:Place 1.0</p> <p>1:PopulatedPlace 0.5</p> <p>2:HistoricPlace 0.5</p> <p>3:#Part 0.38</p> <p>4:#Office 0.35</p> <p>5:#Residence 0.32</p> <p>6:#Role 0.30</p> <p>7:Town 0.30</p> <p>8:#Point 0.29</p> <p>9:PersonFunction 0.27</p> <p>10:City 0.25</p> <p>11:#Location 0.21</p> <p>12:#Address 0.21</p> <p>13:Village 0.16</p> <p>14:Country 0.14</p>	<p>0:@cylinderBore 0.53</p> <p>1:torchBearer 0.26</p> <p>2:flagBearer 0.26</p> <p>3:crosses 0.24</p> <p>4:@name 0.20</p> <p>5:#motto 0.19</p> <p>6:child 0.17</p> <p>7:@birthYear 0.15</p> <p>8:birthPlace 0.15</p> <p>9:@birthName 0.15</p> <p>10:@birthDate 0.15</p> <p>11:family 0.15</p> <p>12:@weight 0.13</p> <p>13:@year 0.12</p> <p>14:origin 0.11</p> <p>15:@teamName 0.11</p> <p>16:@meshName 0.11</p> <p>17:@iupacName 0.11</p>	<p>0:Writer 0.76</p> <p>1:#Author 0.5</p> <p>2:Book 0.36</p> <p>3:Artist 0.22</p> <p>4:Philosopher 0.18</p> <p>5:Actor 0.18</p> <p>6:#Editor 0.15</p> <p>7:Language 0.11</p> <p>8:MusicalArtist 0.11</p> <p>9:Work 0.10</p> <p>10:#Publisher 0.10</p>	<p>0:writer 0.48</p> <p>1:@pages 0.42</p> <p>2:page 0.42</p> <p>3:author 0.39</p> <p>4:@note 0.34</p> <p>5:@notes 0.34</p> <p>6:@volume 0.26</p> <p>7:@name 0.23</p> <p>8:@grayPage 0.22</p> <p>9:spokenIn 0.21</p> <p>10:@address 0.20</p> <p>11:editor 0.19</p> <p>12:@date 0.18</p> <p>13:artist 0.17</p> <p>14:language 0.16</p> <p>15:@year 0.15</p> <p>16:@title 0.15</p> <p>17:recordedIn 0.13</p>
--	--	---	---

SPARQL Generation



The translation of a semantic graph query to SPARQL is straightforward given the mappings

Concepts

- Place => Place
- Author => Writer
- Book => Book

Relations

- born in => birthPlace
- wrote => author


→

```
PREFIX dbo: <http://dbpedia.org/ontology/>

SELECT DISTINCT ?x, ?y WHERE {
  ?o a dbo:Book .
  ?o rdfs:label ?label0 .
  ?label0 bif:contains "'The adventures of Tom Sawyer'" .
  ?x a dbo:Writer .
  ?y a dbo:Place .
  {?o dbo:author ?x}
  {?x dbo:birthPlace ?y} .
}
```

43/49

Evaluation




- 33 test questions from 2011 *Workshop on Question Answering over Linked Data* answerable using DBpedia
- Three human subjects unfamiliar with DBpedia translated the test questions into semantic graph queries
- Compared with two top natural language QA systems: [PowerAqua](#) and [True Knowledge](#)

		Prec.	Recall	F
GoRelations	regular	0.687	0.722	0.704
	concise	0.736	0.803	0.768
PowerAqua	1st triple	0.372	0.483	0.420
	all triples	0.334	0.483	0.395
	merged	0.255	0.291	0.272
True Knowledge		0.469	0.535	0.500

44/49

URL	height	date of birth
http://dbpedia.org/resource/Jack_Shapiro	1.5494	1907-03-22
http://dbpedia.org/resource/Josh_Betts	1.5748	1982-08-25
http://dbpedia.org/resource/Nate_Abrams	1.6256	1897-12-25
http://dbpedia.org/resource/Herbert_Clow	1.6256	1899-05-07
http://dbpedia.org/resource/Trindon_Holiday	1.651	1966-04-27
http://dbpedia.org/resource/Earl_Warweg	1.6764	1897-01-11
http://dbpedia.org/resource/Jacuzzi_Rodgers	1.6764	1994-02-06
http://dbpedia.org/resource/John_Brailer	1.6764	1876-12-27
http://dbpedia.org/resource/Menz_Lindsey	1.6764	1897-07-25
http://dbpedia.org/resource/Stefan_Logan	1.6764	1981-06-02
http://dbpedia.org/resource/Cory_Ross	1.6764	1982-09-22
http://dbpedia.org/resource/Darren_Sproles	1.6764	1983-06-20
http://dbpedia.org/resource/Jacuzzi_Rodgers	1.6764	1994-02-06
http://dbpedia.org/resource/Eddie_Schater	1.6764	1902-01-26
http://dbpedia.org/resource/John_Barret_%28American_football%29	1.6764	1899-02-25
http://dbpedia.org/resource/Dick_Thornton_%28American_football%29	1.6764	1908-02-04
http://dbpedia.org/resource/Walter_French_%28baseball%29	1.7018	1899-07-12
http://dbpedia.org/resource/Al_Loeb	1.7018	1890-03-11
http://dbpedia.org/resource/Amvann_Carter	1.7018	1981-09-09
http://dbpedia.org/resource/Terry_Cutley	1.7018	1984-06-22
http://dbpedia.org/resource/Kendall_Hunter	1.7018	1988-09-16
http://dbpedia.org/resource/LaRod_Stephens-Howling	1.7018	1987-04-26
http://dbpedia.org/resource/Jayson_Foster	1.7018	1985-07-22
http://dbpedia.org/resource/Mark_McMillan	1.7018	1970-04-29
http://dbpedia.org/resource/Al_Cornwee	1.7018	1906-07-16
http://dbpedia.org/resource/Bodie_Weldon	1.7018	1895-11-07
http://dbpedia.org/resource/Bill_Stobbs	1.7018	1896-05-28
http://dbpedia.org/resource/Johnny_Bryan	1.7018	1897-02-28
http://dbpedia.org/resource/Andrew_Hawkins	1.7018	1986-03-10
http://dbpedia.org/resource/Chad_Owens	1.7018	1987-04-03
http://dbpedia.org/resource/Mark_Higgs	1.7018	1987-04-03
http://dbpedia.org/resource/Eamonn_McLennox	1.7018	1984-08-17
http://dbpedia.org/resource/Danny_Woodhead	1.7018	1984-08-17
http://dbpedia.org/resource/Gauntlett_Woolfe	1.7018	1984-08-17
http://dbpedia.org/resource/Shiand_Williams	1.7018	1980-02-10

Current challenges



- Baseline system works well for DBpedia
- Current challenges we are addressing are
 - Adding direct entity matching
 - Relaxing the need for type information
 - Testing on other LOD collections and extending to a set of distributed LOD collections
 - Developing a better Web interface
 - Allowing user feedback and advice
- See <http://ebiq.org/93> for more information & try our alpha version at <http://ebiq.org/GOR>

Final Conclusions

- Linked Data is an emerging paradigm for sharing structured and semi-structured data
 - Backed by machine-understandable semantics
 - Based on successful Web languages and protocols
- Generating and exploring Linked Data resources can be challenging
 - Schemas are large, too many URIs
- New tools for mapping tables to Linked Data and translating structured natural language queries help reduce the barriers

49/49

<http://ebiq.org/>