Graph-based RDF Data Managment

Lei Zou

Peking University Institute of Computer Science and Technology

◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ 三臣 = のへぐ

RDF and Semantic Web

- RDF is a language for the conceptual modeling of information about web resources
- A building block of semantic web
 - Facilitates exchange of information
 - Search engines can retrieve more relevant information

- Facilitates data integration (mashes)
- Machine understandable
 - Understand the information on the web and the interrelationships among them

What's Sematic Web: A Simple Example (RDFa)

The traditional Web (HTML) only considers the display of the content.

How is the page displayed, such as which font and the format of the pictures ?

```
<html>
   <font size="3" color="red"> Lei Zou </font>
   \langle hr \rangle
   Email:<a href="mailto: zoulei@pku.edu.cn">
zoulei@pku.edu.cn </a>
   <font size="3" color="black">Publications: </font>
   <div>
      Lei Zou, Jinhui Mo, Lei Chen, M. Tamer Ozsu,
Dongyan Zhao, gStore: Answering SPARQL Queries Via
Subgraph Matching, VLDB, 2011
   \langle div \rangle
</html>
```

What's Sematic Web: A Simple Example (RDFa)

Sematic Web considers the sematics of the content. What does the content in the page mean? e.g., What are the mean of "zoulei@pku.edu.cn" and "VLDB" ?

```
<html>
<div resource="#me" typeof="Person" >
<font size="3" color="red"> <span property= http://xmlns.com/foaf/0.1/name> Lei
Zou </span> </font>
<hr/>
<a property=" http://xmlns.com/foaf/0.1/mbox" href= "mailto: zoulei@pku.edu.cn "
> zoulei@pku.edu.cn </a>
<font size="3" color="black">Publications: </font>
\langle p \rangle
<div resource="www.vldb.org/pvldb/vol4/p482-zou.pdf">
   <span property=" http://purl.org/dc/terms/contributor"> Lei Zou </span>.
   <span property="http://purl.org/dc/terms/contributor"> Jinghui Mo </span>,
   <span property=" http://purl.org/dc/terms/contributor"> Lei Chen </span>,
   <span property=" http://purl.org/dc/terms/contributor"> M. Tamer Özsu</span>.
   <span property=" http://purl.org/dc/terms/contributor"> Dongyan Zhao</span>,
   <span property=" http://purl.org/dc/terms/title"> gStore: Answering SPAROL
Queries Via Subgraph Matching </span>,
   <span property=" http://purl.org/dc/terms/Publisher"> VLDB </span>
   <span property=" http://purl.org/dc/terms/Date">2011</span>
\langle div \rangle
</html>
```

URL:	HTML			
<html> <div resourc<br=""><font #me"="" 3"="" color="red" size="
Zou </span:

</td><td>e=" typeof="
"> <s > </s =" http://xmins.cor ku.edu.cn "3" color="black">f "3" color="black">f ce="www.vldb.org/p</div></html>	Person" > span property= http://xmlns.com/foaf/0.1/name> Lei n/foaf/0.1/mbox" href= "mailto: zoulei@pku.edu.cn " Publications: wldb/vol4/p482-zou.pdf">	PREVIEW	Examples	
laximum 15	00 characters allow	ved. Over-length text will be truncated		

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

Preview

Extracted structured data	
rdfa-node	
relationship:	
name:	mbox
value:	zoulei@pku.edu.cn
href:	mailto:%20zoulei@pku.edu.cn
property:	
name:	Lei Zou
contributor:	Lei Zou
contributor:	Jinghui Mo
contributor:	Lei Chen
contributor:	M. Tamer Özsu
contributor:	Dongyan Zhao
title:	gStore: Answering SPARQL Queries Via Subgraph Matching
Publisher:	VLDB
Date	2011

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 善臣 - のへで



Enhanced search result preview

Disclaimer: this preview is only shown as a example of what a search engine might display. It is to the discretion of each search engine provider to decide whether your page will be displayed as an enhanced search result or not in their search results pages.

Bob Smith

Enter structured-data.org/examples/google-rs/people.md.html Abuquinque, NM - engineer, ACME Corp. an actual search result may display other content relating to your search terms here.

Raw structured data extracted from the page:

rdfitype	vad:Person	
vadiaddrosa	rdfitype vmdiAddre	
	wmd:locality	Alboquerque
	vmdiregion	SR
vmd:affiliation	ACHE Corp	
vndinano	Bob Smith	
vndinickname	Smithy	
vmd:title	engineer	
vmdsurl	http://www.asheple.com	

Parser statistics

Statements	12
Templates	http://data-vocabulary.org/Person http://data-vocabulary.org/Address

What's Sematic Web: Facebook Social Graph



Sac

What's Sematic Web: From Two Perspectives

Expressiveness	↑	
More Semantic; More Powerful Reasoning	RDF, RDFS,OWL,OWL Full, OWL 2 	Open Linked Data, Web-scale Triple Store, Semantic Wiki

How to get more data ? Scalability How to manage the Web-scale Semantic Data? (=)

What's Sematic Web: From Two Perspectives



How to get more data ? Scalability How to manage the Web-scale Semantic Data? $\langle a \rangle = \langle a \rangle \langle a \rangle$

596

Some Interesting Products

IBM Watson \$1,400 \$200 \$0 ken WATSON Jude 98% sad song make it better

▲□▶ ▲□▶ ▲目▶ ▲目▶ = 目 - のへぐ

Some Interesting Products

EVI— acquired by Amazon on October 2012.



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

Some Interesting Products

Google Knowledge Graph



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

RDF Uses

► Yago and DBPedia extract facts from Wikipedia & represent as RDF → structural queries

- Communities build RDF data
 - ▶ E.g., biologists: Bio2RDF and Uniprot RDF
- Web data integration
 - Linked Data Cloud

▶ ...

- ... are growing and fast
 - Linked data cloud currently consists of 325 datasets with >25B triples

◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ 三臣 = のへぐ

Size almost doubling every year

- ... are growing and fast
 - Linked data cloud currently consists of 325 datasets with >25B triples
 - Size almost doubling every year





▲□▶ ▲□▶ ▲□▶ ▲□▶ □ □ のへで

Linking Open Data cloud diagram, by Richard Cyganiak and Anja Jentzsch. http://lod-cloud.net/

- ... are growing and fast
 - Linked data cloud currently consists of 325 datasets with >25B triples
 - Size almost doubling every year





Linking Open Data cloud diagram, by Richard Cyganiak and Anja Jentzsch. http://lod-cloud.net/

◆□ > ◆母 > ◆臣 > ◆臣 > ○臣 - のへで

- ... are growing and fast
 - Linked data cloud currently consists of 325 datasets with >25B triples
 - Size almost doubling every year





Linking Open Data cloud diagram, by Richard Cyganiak and Anja Jentzsch. http://lod-cloud.net/

・ロト・西ト・山田・山田・山市・山市・山市・山市

- ... are growing and fast
 - Linked data cloud currently consists of 325 datasets with >25B triples
 - Size almost doubling every year



April '14: 1091 datasets, ??? triples

Max Schmachtenberg, Christian Bizer, and Heiko Paulheim: Adoption of Linked Data Best Practices in Different Topical Domains. In *Proc. ISWC*, 2014.

Outline

RDF Introduction

gStore: a graph-based SPARQL query engine Answering SPARQL queries using graph pattern matching [Zou et al., PVLDB 2011, VLDB J 2014]

gAnswer: Natural Language Question Answering over RDF A Data Driven Approach [Zou et al., SIGMOD 2014; Zheng et al., SIGMOD 2015]

Outline

RDF Introduction

gStore: a graph-based SPARQL query engine Answering SPARQL queries using graph pattern matching [Zou et al., PVLDB 2011, VLDB J 2014]

gAnswer: Natural Language Question Answering over RDF A Data Driven Approach [Zou et al., SIGMOD 2014; Zheng et al., SIGMOD 2015]

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

Everything is an uniquely named resource



◆□ ▶ ◆□ ▶ ◆ 三 ▶ ◆ 三 ● ● ● ●

http://en.wikipedia.org/wiki/Abraham_Lincoln

- Everything is an uniquely named resource
- Namespaces can be used to scope the names

xmlns:y=http://en.wikipedia.org/wiki y:Abraham_Lincoln



◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ 三臣 = のへぐ

- Everything is an uniquely named resource
- Namespaces can be used to scope the names
- Properties of resources can be defined

xmlns:y=http://en.wikipedia.org/wiki y:Abraham_Lincoln



Abraham_Lincoln:hasName "Abraham Lincoln" Abraham_Lincoln:BornOnDate: "1809-02-12" Abraham_Lincoln:DiedOnDate: "1865-04-15"

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQ@

- Everything is an uniquely named resource
- Namespaces can be used to scope the names
- Properties of resources can be defined
- Relationships with other resources can be defined

xmlns:y=http://en.wikipedia.org/wiki y:Abraham_Lincoln



Abraham_Lincoln:hasName "Abraham Lincoln" Abraham_Lincoln:BornOnDate: "1809-02-12" Abraham_Lincoln:DiedOnDate: "1865-04-15"

Abraham_Lincoln:DiedIn



y:Washington_DC

- Everything is an uniquely named resource
- Namespaces can be used to scope the names
- Properties of resources can be defined
- Relationships with other resources can be defined
- Resources can be contributed by different people/groups and can be located anywhere in the web
 - Integrated web "database"

xmlns:y=http://en.wikipedia.org/wiki y:Abraham_Lincoln



Abraham_Lincoln:hasName "Abraham Lincoln" Abraham_Lincoln:BornOnDate: "1809-02-12" Abraham_Lincoln:DiedOnDate: "1865-04-15"

Abraham_Lincoln:DiedIn



y:Washington_DC

RDF Data Model

 Triple: Subject, Predicate (Property), Object (s, p, o)
 Subject: the entity that is described (URI or blank node)
 Predicate: a feature of the entity (URI)
 Object: value of the feature (URI, blank node or literal)



► $(s, p, o) \in (U \cup B) \times U \times (U \cup B \cup L)$

U: set of URIs

- B: set of blank nodes
- L: set of literals

Set of RDF triples is called an RDF graph

Subject	Predicate	Object
Abraham_Lincoln	hasName	"Abraham Lincoln"
Abraham_Lincoln	BornOnDate	"1809-02-12"
Abraham_Lincoln	DiedOnDate	"1865-04-15"

RDF Example Instance

Prefix: y=http://en.wikipedia.org/wiki Subject Predicate Object v: Abraham_Lincoln basName "Abraham Lincoln" Literal v: Abraham_Lincoln BornOnDate "1809-02-12" v: Abraham_Lincoln DiedOnDate "1865-04-15" UR v:Abraham_Lincoln hornIn v:Hodgenville_KY v: Abraham Lincoln DiedIn y: Washington_DC y:Abraham_Lincoln title "President" y:Abraham_Lincoln gender "Male" v: Washington_DC hasName "Washington D.C." URI y:Washington_DC foundingYear "1790" y:Hodgenville_KY hasName "Hodgenville" y:United_States hasName "United States" y:United_States hasCapital < y:Washington_DC y:United_States foundingYear "1776" y:Reese_Witherspoon bornOnDate "1976-03-22" y:Reese_Witherspoon bornIn y:New_Orleans_LA y:Reese_Witherspoon hasName "Reese Witherspoon" y:Reese_Witherspoon gender "Female" y:Reese_Witherspoon title "Actress" v:New_Orleans_LA foundingYear "1718" v:New_Orleans_LA locatedIn v:United_States v:Franklin_Roosevelt hasName "Franklin D. Roosevelt" v:Franklin_Roosevelt bornIn v:Hvde_Park_NY v:Franklin_Roosevelt title "President" v:Franklin_Roosevelt "Male" gender v:Hvde_Park_NY foundingYear "1810" v:Hvde_Park_NY locatedIn v:United_States v:Marilvn_Monroe "Female" gender v:Marilvn_Monroe hasName "Marilyn Monroe" v:Marilvn_Monroe hornOnDate "1926-07-01" v:Marilvn_Monroe diedOnDate "1962-08-05"

Sac

RDF Graph



◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ○ □ ○ ○ ○ ○

RDF Query Model

- Query Model SPARQL Protocol and RDF Query Language
- ► Given U (set of URIs), L (set of literals), and V (set of variables), a SPARQL expression is defined recursively:
 - an atomic triple pattern, which is an element of

 $(U \cup V) \times (U \cup V) \times (U \cup V \cup L)$

- ?x hasName "Abraham Lincoln"
- ▶ *P* FILTER *R*, where *P* is a graph pattern expression and *R* is a built-in SPARQL condition (i.e., analogous to a SQL predicate)

?x price ?p FILTER(?p < 30)</p>

 P1 AND/OPT/UNION P2, where P1 and P2 are graph pattern expressions

Example:

SELECT ?name

WHERE $\{$

```
?m <bornIn> ?city. ?m <hasName> ?name.
?m<bornOnDate> ?bd. ?city <foundingYear> ``1718``.
FILTER(regex(str(?bd),``1976``))
```

SPARQL Queries



Naïve Triple Store Design

SELECT ?name
WHERE {
 ?m <bornIn > ?city. ?m <hasName> ?name.
 ?m<bornOnDate> ?bd. ?city <foundingYear> ''1718''.
FILTER(regex(str(?bd), ''1976''))

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQ@

Subject	Property	Object
y:Abraham_Lincoln	hasName	"Abraham Lincoln"
y:Abraham_Lincoln	bornOnDate	"1809-02-12"
y:Abraham_Lincoln	diedOnDate	"1865-04-15"
y:Abraham_Lincoln	bornIn	y:Hodgenville_KY
y:Abraham_Lincoln	diedIn	y:Washington_DC
y:Abraham_Lincoln	title	"President"
y:Abraham_Lincoln	gender	"Male"
y:Washington_DC	hasName	"Washington D.C."
y:Washington_DC	foundingYear	"1790"
y:Hodgenville_KY	hasName	"Hodgenville"
y:United_States	hasName	"United States"
y:United_States	hasCapital	y:Washington_DC
y:United_States	foundingYear	"1776"
y:Reese_Witherspoon	bornOnDate	"1976-03-22"
y:Reese_Witherspoon	bornIn	y:New_Orleans_LA
y:Reese_Witherspoon	hasName	"Reese Witherspoon"
y:Reese_Witherspoon	gender	"Female"
y:Reese_Witherspoon	title	"Actress"
y:New_Orleans_LA	foundingYear	"1718"
y:New_Orleans_LA	locatedIn	y:United_States
y:Franklin_Roosevelt	hasName	"Franklin D. Roo-
		sevelt"
y:Franklin_Roosevelt	bornIn	y:Hyde_Park_NY
y:Franklin_Roosevelt	title	"President"
y:Franklin_Roosevelt	gender	"Male"
y:Hyde_Park_NY	foundingYear	"1810"
y:Hyde_Park_NY	locatedIn	y:United_States
y:Marilyn_Monroe	gender	"Female"
y:Marilyn_Monroe	hasName	"Marilyn Monroe"
y:Marilyn_Monroe	bornOnDate	"1926-07-01"
y:Marilyn_Monroe	diedOnDate	"1962-08-05"

}

Naïve Triple Store Design

SELECT ?name
WHERE {
 ?m <bornln > ?city. ?m <hasName> ?name.
 ?m<bornOnDate> ?bd. ?city <foundingYear > 1''1718''.
 FILTER(regex(str(?bd), ''1976''))

Subject	Property	Object
y:Abraham_Lincoln	hasName	"Abraham Lincoln"
y:Abraham_Lincoln	bornOnDate	"1809-02-12"
y:Abraham_Lincoln	diedOnDate	"1865-04-15"
y:Abraham_Lincoln	bornIn	y:Hodgenville_KY
y:Abraham_Lincoln	diedIn	y:Washington_DC
y:Abraham_Lincoln	title	"President"
y:Abraham_Lincoln	gender	"Male"
y:Washington_DC	hasName	"Washington D.C."
y:Washington_DC	foundingYear	"1790"
y:Hodgenville_KY	hasName	"Hodgenville"
y:United_States	hasName	"United States"
y:United_States	hasCapital	y:Washington_DC
y:United_States	foundingYear	"1776"
y:Reese_Witherspoon	bornOnDate	"1976-03-22"
y:Reese_Witherspoon	bornIn	y:New_Orleans_LA
y:Reese_Witherspoon	hasName	"Reese Witherspoon"
y:Reese_Witherspoon	gender	"Female"
y:Reese_Witherspoon	title	"Actress"
y:New_Orleans_LA	foundingYear	"1718"
y:New_Orleans_LA	locatedIn	y:United_States
y:Franklin_Roosevelt	hasName	"Franklin D. Roo-
		sevelt"
y:Franklin_Roosevelt	bornIn	y:Hyde_Park_NY
y:Franklin_Roosevelt	title	"President"
y:Franklin_Roosevelt	gender	"Male"
y:Hyde_Park_NY	foundingYear	"1810"
y:Hyde_Park_NY	locatedIn	y:United_States
y:Marilyn_Monroe	gender	"Female"
y:Marilyn_Monroe	hasName	"Marilyn Monroe"
y:Marilyn_Monroe	bornOnDate	"1926-07-01"
y:Marilyn_Monroe	diedOnDate	"1962-08-05"

SELECT T2.object				
FROM T as T1, T as T2, T as T3,				
T as T4				
WHERE T1. property="bornin"				
AND T2.property="hasName"				
AND T3. property="bornOnDate"				
AND T1.subject=T2.subject				
AND T2.subject=T3.subject				
AND T4. propety="foundingYear"				
AND T1.object=T4.subject				
AND T4.object="1718"				
AND T3.object LIKE '%1976%'				

Naïve Triple Store Design

SELECT ?name WHERE {

?m <bornIn > ?city. ?m <hasName> ?name. ?m<bornOnDate> ?bd. ?city <foundingYear> FILTER(regex(str(?bd),''1976''))

Subject	Property	Object
y:Abraham_Lincoln	hasName	"Abraham Lincoln"
y:Abraham_Lincoln	bornOnDate	"1809-02-12"
y:Abraham_Lincoln	diedOnDate	"1865-04-15"
y:Abraham_Lincoln	bornIn	y:Hodgenville_KY
y:Abraham_Lincoln	diedIn	y:Washington_DC
y:Abraham_Lincoln	title	"President"
y:Abraham_Lincoln	gender	"Male"
y:Washington_DC	hasName	"Washington D.C."
y:Washington_DC	foundingYear	"1790"
y:Hodgenville_KY	hasName	"Hodgenville"
y:United_States	hasName	"United States"
y:United_States	hasCapital	y:Washington_DC
y:United_States	foundingYear	"1776"
y:Reese_Witherspoon	bornOnDate	"1976-03-22"
y:Reese_Witherspoon	bornIn	y:New_Orleans_LA
y:Reese_Witherspoon	hasName	"Reese Witherspoon"
y:Reese_Witherspoon	gender	"Female"
y:Reese_Witherspoon	title	"Actress"
y:New_Orleans_LA	foundingYear	"1718"
y:New_Orleans_LA	locatedIn	y:United_States
y:Franklin_Roosevelt	hasName	"Franklin D. Roo-
		sevelt"
y:Franklin_Roosevelt	bornIn	y:Hyde_Park_NY
y:Franklin_Roosevelt	title	"President"
y:Franklin_Roosevelt	gender	"Male"
y:Hyde_Park_NY	foundingYear	"1810"
y:Hyde_Park_NY	locatedIn	y:United_States
y:Marilyn_Monroe	gender	"Female"
y:Marilyn_Monroe	hasName	"Marilyn Monroe"
y:Marilyn_Monroe	bornOnDate	"1926-07-01"
y:Marilyn_Monroe	diedOnDate	"1962-08-05"

Too many self-joins!

SELECT T2.object FROM T as T1, T as T2, T as T3, T as T4 WHERE T1.property="bornln" AND T2.property="bornOnDate" AND T3.property="bornOnDate" AND T1.subject=T2.subject AND T2.subject=T3.subject AND T4.propety="foundingYear" AND T1.object=T4.subject AND T4.object="1718" AND T3.object LIKE '%1976%'

Existing Solutions

- 1. Property table
 - \blacktriangleright Each class of objects go to a different table \Rightarrow similar to normalized relations
 - Eliminates some of the joins
- 2. Vertically partitioned tables
 - For each property, build a two-column table, containing both subject and object, ordered by subjects
 - Can use merge join (faster)
 - Good for subject-subject joins but does not help with subject-object joins
- 3. Exhaustive indexing
 - Create indexes for each permutation of the three columns
 - Query components become range queries over individual relations with merge-join to combine
 - Excessive space usage
Property Tables

Subjee y:Abraham_ Reese With

- Grouping by entities; Jena [Wilkinson et al., SWDB 03] ,FlexTable [Wang et al., DASFAA 10], DB2-RDF [Bornea et al., SIGMOD 13]
- Clustered property table: group together the properties that tend to occur in the same (or similar) subjects
- Property-class table: cluster the subjects with the same type of property into one property table

Sı	ubject	Property	Obje	ect				
y:/	Abraham_Lincoln	hasName	"Abr	ahan	n Lincoln"			
y:/	Abraham_Lincoln	bornOnDate	"180	9-02	-12"			
y:/	Abraham_Lincoln	diedOnDate	"186	5-04	-15"			
y:\	Washington_DC	hasName	"Wa	shing	gton D.C."			
y:\	Washington_DC	foundingYea	r "179	0"				
	1							
	\				Subject		hasName	foundingYear
	hasName I	oornOnDate died	IOnDate]	y:Washington_I	DC "	Washington D.C."	1790
ncoln	"Abraham Lincoln"	1809-02-12 180	5-04-15		y:Hyde_Park_N	IY	"Hyde Park"	1810'
rspoon	"Reese Witherspoon"	1976-03-22						

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 善臣 - のへで

Property Tables

- Grouping by entities; Jena [Wilkinson et al., SWDB 03] ,FlexTable [Wang et al., DASFAA 10], DB2-RDF [Bornea et al., SIGMOD 13]
- Clustered property table: group together the properties that

Advantages

- Fewer joins
- If the data is structured, we have a relational system similar to normalized relations



Property Tables

- Grouping by entities; Jena [Wilkinson et al., SWDB 03] ,FlexTable [Wang et al., DASFAA 10], DB2-RDF [Bornea et al., SIGMOD 13]
- Clustered property table: group together the properties that

Advantages

- Fewer joins
- If the data is structured, we have a relational system similar to normalized relations

Disadvantages

- Potentially a lot of NULLs
- Clustering is not trivial
- Multi-valued properties are complicated

Subject	hasName	bornOnDate	diedOnDate
y:Abraham_Lincoln	"Abraham Lincoln"	1809-02-12	1865-04-15
y:Reese_Witherspoon	"Reese Witherspoon"	1976-03-22	

Jubject	naarvanne	rounding rear
y:Washington_DC	"Washington D.C."	1790
y:Hyde_Park_NY	"Hyde Park"	1810'

▲□▶ ▲圖▶ ▲臣▶ ▲臣▶ 三臣 - のへで

Binary Tables

- Grouping by properties: For each property, build a two-column table, containing both subject and object, ordered by subjects [Abadi et al., VLDB 07]
- Also called vertical partitioned tables
- n two column tables (n is the number of unique properties in the data)

Subject	Property	Object		
y:Abraham_Lincoln	hasName	"Abraham Lincoln"		
y:Abraham_Lincoln	bornOnDate	"1809-02-12"		
y:Abraham_Lincoln	diedOnDate	"1865-04-15"		
y:Washington_DC	hasName	"Washington D.C."		
y:Washington_DC	foundingYear	"1790"		

hasName

Subject	Object
y:Abraham_Lincoln	"Abraham Lincoln"
y:Washington_DC	"Washington D.C."

bornOnDate				
Subject	Object			
y:Abraham_Lincoln	1809-02-12			
y:Reese_Witherspoon	1976-03-22			

foundingYear

Subject	Object
y:Washington_DC	1790
y:Hyde_Park_NY	1810

◆□▶ ◆圖▶ ◆臣▶ ◆臣▶ ─臣 ─の�?

Binary Tables

- Grouping by properties: For each property, build a two-column table, containing both subject and object, ordered by subjects
 Advantages
 - Supports multi-valued properties
 - No NULLs
 - No clustering
 - Read only needed attributes (i.e. less I/O)
 - Good performance for subject-subject joins

y:washington_DC	nasivame	vvasnington D.C.
y:Washington_DC	foundingYear	"1790"

hasName

Subject	Object
y:Abraham_Lincoln	"Abraham Lincoln"
y:Washington_DC	"Washington D.C."

bornOnDate				
Subject	Object			
y:Abraham_Lincoln	1809-02-12			
y:Reese_Witherspoon	1976-03-22			

foundingYear

Subject	Object
y:Washington_DC	1790
y:Hyde_Park_NY	1810

Binary Tables

- Grouping by properties: For each property, build a two-column table, containing both subject and object, ordered by subjects
 Advantages
 - Supports multi-valued properties
 - No NULLs
 - No clustering
 - Read only needed attributes (i.e. less I/O)
 - Good performance for subject-subject joins

Disadvantages

- Not useful for subject-object joins
- Expensive inserts

Subject	Object]	Subject	Object		Subject	UDJELL	
y:Abraham_Lincoln	"Abraham Lincoln"	1	y:Abraham_Lincoln	1809-02-	12	y:Washington_DC	1/90	
y:Washington_DC	"Washington D.C."	1	y:Reese_Witherspoon	1976-03-	22	y.Hyue_Fark_INT	1010	

Exhaustive Indexing

- RDF-3X [Neumann and Weikum, PVLDB 08], Hexastore [Weiss et al., PVLDB 08]
- Strings are mapped to ids using a mapping table

Original triple table

Subject	Property	Object
y:Abraham_Lincoln	hasName	"Abraham Lincoln"
y:Abraham_Lincoln	bornOnDate	"1809-02-12"
y:Abraham_Lincoln	diedOnDate	"1865-04-15"
y:Washington_DC	hasName	"Washington D.C."
y:Washington_DC	foundingYear	"1790"

Encoded triple table

Subject	Property	Object	
0	1	2	
0	3	4	
0	5	6	
7	1	8	
7	9	10	

Mapping table

ID	Value
0	y:Abraham_Lincoln
1	hasName
2	"Abraham Lincoln"
3	bornOnDate
4	"1809-02-12"
5	diedOnDate
6	"1865-04-15"
7	y:Washington_DC
8	"Washington D.C."
9	foundingYear
10	"1790"

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Exhaustive Indexing

- RDF-3X [Neumann and Weikum, PVLDB 08], Hexastore [Weiss et al., PVLDB 08]
- Strings are mapped to ids using a mapping table
- Triples are indexed in a clustered B+ tree in lexicographic order



Exhaustive Indexing

- RDF-3X [Neumann and Weikum, PVLDB 08], Hexastore [Weiss et al., PVLDB 08]
- Strings are mapped to ids using a mapping table
- Triples are indexed in a clustered B+ tree in lexicographic order
- Create indexes for permutations of the three columns: SPO, SOP, PSO, POS, OPS, OSP



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Exhaustive Indexing–Query Execution

- Each triple pattern can be answered by a range query
- Joins between triple patterns computed using merge join
- Join order is easy due to extensive indexing

Subject	Property	Object
0	1	2
0	3	4
0	5	6
7	1	8
7	9	10
:	:	

ID	Value
0	y:Abraham_Lincoln
1	hasName
2	"Abraham Lincoln"
3	bornOnDate
4	"1809-02-12"
5	diedOnDate
6	"1865-04-15"
7	y:Washington_DC
8	"Washington D.C."
9	foundingYear
10	"1790"

Exhaustive Indexing–Query Execution

- Each triple pattern can be answered by a range query
- Joins between triple patterns computed using merge join
- Join order is easy due to extensive indexing

Advantages

- Eliminates some of the joins they become range queries
- Merge join is easy and fast

7	1	8
7	9	10
:	:	:

3	bornOnDate
4	"1809-02-12"
5	diedOnDate
6	"1865-04-15"
7	y:Washington_DC
8	"Washington D.C."
9	foundingYear
10	"1790"

Exhaustive Indexing–Query Execution

- Each triple pattern can be answered by a range query
- Joins between triple patterns computed using merge join
- Join order is easy due to extensive indexing

Advantages

- Eliminates some of the joins they become range queries
- Merge join is easy and fast

Disadvantages

Space usage

6	"1865-04-15"
7	y:Washington_DC
8	"Washington D.C."
9	foundingYear
10	"1790"

Outline

RDF Introduction

gStore: a graph-based SPARQL query engine Answering SPARQL queries using graph pattern matching [Zou et al., PVLDB 2011, VLDB J 2014]

gAnswer: Natural Language Question Answering over RDF A Data Driven Approach [Zou et al., SIGMOD 2014; Zheng et al., SIGMOD 2015]

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQ@

gStore – General Idea

- We work directly on the RDF graph and the SPARQL query graph
 - Answering SPARQL query \equiv subgraph matching
 - Subgraph matching is computationally expensive
- Use a signature-based encoding of each entity and class vertex to speed up matching
- Filter-and-evaluate
 - Use a false positive algorithm to prune nodes and obtain a set of candidates; then do more detailed evaluation on those
- We develop an index (VS*-tree) over the data signature graph (has light maintenance load) for efficient pruning

0. Start with RDF Graph G



0. Start with RDF Graph G



gStore



◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ● ○ ○ ○ ○

gStore

C++ API

Ē1	Harristan "Salaration at "
21	<pre>#loclodecistreamy</pre>
	// before can this mampin, you must start up the STORY served at fills user manned ./get the mainfine args, char * args()
1	// initialize the Object server's IP address and pert- ourcoednmeeter grf'IIV.s.s.', (30)g:
	// Dalls a new desidence by a form film: // onto that the scalaring path is reliabed to gammag. go.ball4f*ds_lime(v*, *cample/vif_sign/comm_lt_commet.m)*1;
	(/ then you can examine the difficult party and the distribution. MAINING party = ***********************************
	stdilout of snown of stdilendil
	// infand this detabane. me.undedf*de_liner:*\$;
12	7/ slop, you was load some maint skyldaws tirectly and then youty. 06.1004(*0) [Junif/31] answer a 0. servicestly]
	stditopat ee answer ee stditendir
	return 1/

JAVA API

import pprc.GatoreConnector)

```
() Merer hat the example, you must start up the Disnet solver at first lush command J_DISTURCE_
philo: these beachThroughe
philo: exact you maintering[] areas
```

```
// initializes the effort server's it astrona and port.
GetoreContentor gt - new GetoreContentor (*17).5.5.2*, (2013).
() build a new metabane by a MDF file.
// note that the relative path is related to gentyer-
ge builds "do nimeto", "excepte feit artginforme in offenes siva :
If then you can essente StMAGL many in this Calabiant,
String spengi . "select on above "
       · *rs mitrays, university and the ends
      . "in diame elements."
       · "te mitabreditres ip. "
     + "N ghitresheitt 79) *
       · "In minime chilifetifenerili."
      · "To diversifier Ty."
       · "In minner departmently."
String answer a gr-parylipargil;
System.out.printinianewers!
If unload this metabase.
gerundenadenin Linnet 1783
17 mins, you can load nome maint matamany strently and then carey.
ad-16849*/in.010901**12
annues - oc. merrdeparelb /
System.out.println&answerp1
```

gStore

Queries: gubject subLocation predict object objLocation (36.1312200154285.41 Drawitt) Drawiff wasBornin * 20 temporal:Start>=1850-XX-XX End <=1900-XX-XX Put # Location Drawitt Drawith type_star wordnet_cit/ Draw it! Put # Location Draw it! temporal:Start End + -Time Line 1850-11908-XX-XX -3000 0 1000 1500 1600 1700 1750 1800 1850 1950 1950 2000 2050 9999 -9999 Use exists examples: altd -Show ! . answers in one page Query ----卫星 混合地图 1965 員大利 保险利益 2.11 HEX No DRE HERP 1 IL MAR ma han -----西说牙 ites, surt

▲□▶ ▲□▶ ▲目▶ ▲目▶ 三目 - のへぐ

Peer Review Comments



Charu Aggarwal, ACM Fellow, IBM T. J. Watson Researcher

In the realm of RDF, SPARQL is widely used as the query processing language. However, writing of a SPARQL query is often too challenging, because it requires the exact knowledge of structure, node labels and types. <u>gStore [43]</u>, which is the first study that considers a subgraph matching-based query answering technique in RDF data, allows approximate node label matching, but adheres to strict structural matches. In contrast, our NeMa framework permits both structural and node label mismatches,

> —NeMa: Fast Graph Search with Label Similarity, Proc. of VLDB: 181-192 (2013)

Outline

RDF Introduction

gStore: a graph-based SPARQL query engine Answering SPARQL queries using graph pattern matching [Zou et al., PVLDB 2011, VLDB J 2014]

gAnswer: Natural Language Question Answering over RDF A Data Driven Approach [Zou et al., SIGMOD 2014; Zheng et al., SIGMOD 2015]

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQ@

gAnswer: Natural Language Question Answering Over Knowledge Graph–A Graph Data Driven Approach

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

gAnswer: Natural Language Question Answering Over Knowledge Graph–A Graph Data Driven Approach

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

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで



Running Example

Question: Who was married to an actor that play in Philadelphia ?

Subject	Property	Object
Antonio_Banderas	type	actor
Antonio_Banderas	spouse	Melanie_Griffith
Antonio_Banderas	starring	Philadelphia_(film)
Philadelphia_(film)	type	film
Jonathan_Demme	director	film
Philadelphia	type	city
Aaron_McKie	bornIn	Philadelphia
James_Anderson	playForTeam	Philadelphia_76ers
Constantin_Stanislavski	create	An_Actor_Prepares
Philadelphia_76ers	type	Basketball_team
An_Actor_Prepares	type	Book

Running Example

Question: Who was married to an actor that play in Philadelphia ?

Subject	Property	Object	
Antonio_Banderas	type	actor	
Antonio_Banderas	spouse	Melanie_Griffith	
Antonio_Banderas	starring	Philadelphia_(film)	
Philadelphia_(film)	type	film	
Jonathan_Demme	director	film	
Philadelphia	type	city	
Aaron_McKie	bornIn	Philadelphia	
James_Anderson	playForTeam	Philadelphia_7	Melanie Griffith
Constantin_Stanislavski	create	An_Actor_Prepa	
Philadelphia_76ers	type	Basketball_team	
An_Actor_Prepares	type	Book	

Existing Solutions



SELECT ?y Translate NL Question to structured queries WHERE { ?x starring Philadelphia_(film). ?x type actor. ?x spouse ?y . } Query Processing Melanie Griffith



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで





▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 - のへ(で)



▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 - のへ(で)



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへ(?)

Experiments: Datasets

RDF repository: DBPedia

	DBpedia
Number of Entities	5.2 million
Number of Triples	60 million
Number of Predicates	1643
Size of RDF Graphs (in GB)	6.1

Table : Statistics of RDF Graph

Relation Phrase Dictionary: Patty

Table : Statistics of Relation Phrase Dataset

	wordnet-wikipedia	freebase-wikipedia
Number of Textual Patterns	350,568	1,631,530
Number of Entity Pairs	3,862,304	15,802,947
Average Entity Pair	11	9
Number For Each Pattern		

Experiments: Online

Benchmark: QALD-3, 99 Natural Language Questions

Table :	Evaluating	QALD-3	Testing	Questions	(on	DBpedia))
---------	------------	--------	---------	-----------	-----	----------	---

	Processed	Right	Partially	Recall	Precision	F-1
Our	76	32	11	0.40	0.40	0.40
Method						
squall2sparql	96	77	13	0.85	0.89	0.87
CASIA	52	29	8	0.36	0.35	0.36
Scalewelis	70	1	38	0.33	0.33	0.33
RTV	55	30	4	0.34	0.32	0.33
Intui2	99	28	4	0.32	0.32	0.32
SWIP	21	14	2	0.15	0.16	0.16
DEANNA	27	21	0	0.21	0.21	0.21

Experiments: Online

ID	Questions	Response Time (in ms)
Q2	Who was the successor of John F. Kennedy?	1699
Q3	Who is the mayor of Berlin?	677
Q14	Give me all members of Prodigy?	811
Q17	Give me all cars that are produced in Germany ?	297
Q19	Give me all people that were born in Vienna and died in Berlin ?	2557
Q20	How tall is Michael Jordan ?	942
Q21	What is the capital of Canada ?	1342
Q22	Who is the governor of Wyoming ?	796
Q24	Who was the father of Queen Elizabeth II?	538
Q27	Sean Parnell is the governor of which U.S. state ?	1210
Q28	Give me all movies directed by Francis Ford Coppola.	577
Q30	What is the birth name of Angela Merkel ?	250
Q35	Who developed Minecraft ?.	2565
Q39	Give me all companies in Munich.	1312
Q41	Who founded Intel?	1105
Q42	Who is the husband of Amanda Palmer ?	1418
Q44	Which cities does the Weser flow through ?	1139
Q45	Which countries are connected by the Rhine ?	736
Q54	What are the nicknames of San Francisco ?	321
Q58	What is the time zone of Salt Lake City ?	316
Q63	Give me all Argentine films.	427
Q70	Is Michelle Obama the wife of Barack Obama ?	316
Q74	When did Michael Jackson die ?	258
Q76	List the children of Margaret Thatcher.	1139
Q77	Who was called Scarface?	719
Q81	Which books by Kerouac were published by Viking Press ?	796
Q83	How high is the Mount Everest ?	635
Q84	Who created the comic Captain America ?	589
Q86	What is the largest city in Australia ?	1419
Q89	In which city was the former Dutch queen Juliana buried ?	1700
Q98	Which country does the creator of Miffy come from ?	2121
Q100	Who produces Orangina ?	367

Experiments: Online



Figure : Online Running Time Comparison

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで
Experiments: Online

QUESTION ANSWERING OVER LINKED DATA QALD-4: Evaluation results

Workshop website: http://www.sc.cit-ec.uni-bielefeld.de/qald

	Total	Processed	Right	Partially	Recall	Precision	F-measure
Year	50	40	34	6	0.71	0.72	0.72
gAnswer	50	25	16	4	0.37	0.37	0.37
CASIA	50	26	15	4	0.40	0.32	0.36
Intui3	50	33	10	4	0.25	0.23	0.24
ISOFT	50	28	10	3	0.26	0.21	0.23
RO FII	50	50	6	0	0.12	0.12	0.12

Results¹ for Task 1: Multilingual question answering over DBpedia

Figure : QALD-4 Results

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

An Example: using gAnswer+gStore in CDBLP





 作者投究
 高级投究

 室小峰 云计算 论文
 词义检察

 常見監察:学术交前ScholarSoace的设的作者列表,学术主项自动中或高衡EavScholar

亚关注学术空间微信号

THE THE ALL ALL A A A ALL THE

为	时为子不会	KI	化义收求消沉		
+	软件学报	+	计算机学报	+	计算机研究与发展
	电子学报	4	中国图象图形学报		中文信息学报
+	计算机科学	+	小型微型计算机系统	+	计算机科学与探索
*	中国科学F辑	•	计算机辅助设计与图形	学学报	
÷	计算机工程与科学	•	计算机工程(暂停收录)		
	中国数据库学术会谈	IND	3C)		



Sac

An Example: using gAnswer+gStore in CDBLP

学 WA	术空间Se MDM, Database	cholarSp Group, Renmin	ace n University of	[:] China			
	首页	检索	期刊	分类	机构	会议	关于
金索	言词: 孟小峭	• 云计算	论文		WHERE (? ?paper :T ?paper :K FILTER(re	?author :Name itle ?title. eywords ?k. gex(str(?k)," 元	``孟小峰''. 计算"))
年間 注: C 里。	家到结果4条 DBLP语义搜索功能图	B查询部分使用了:	北京大学邹磊老师	提供的gAnswer系的	统,在此表示感谢。	了解更多相关信	息请点击这
4	慈祥 马友忠 孟小	峰:一种云环境	下的大数据Top-K	查询方法, 软件学	报, 2014 (04): 813	-826	
3	史英杰 孟小峰:	云数据管理系统	中查询技术研究综	述. 计算机学报, 2	2013 (02): -		
z	孟小峰 葱祥: 2	大数据管理:概念、	技术与挑战. 计算	机研究与发展, 20	13 (01): -		
1	權峰 孟小峰 徐建	建良: 云计算中面	而隐私保护的查询	则处理技术研究.计	+算机科学与探索,:	2012 (05):	

An Example: using gAnswer+gStore in CDBLP

	空间 主 om.cn 全球	办:中国知 最大的数字图刊	网络	1	充值。 建主动 送卡」	PG REAL	1卡 数手编奏 预数 2番 广告服务 Eng
首页 期刊全文庫	学位论文库	会议论文库	年基全文库	学术百科	184	CNICC	200 (200) 1000
基础科学 工程科技印	1 工程科技印刷	1 医新卫生科机	E 位息科技 1	改重科技 1 哲学	与人文科学(社	会科学[辑]	社会科学印辑) 经济管
输入关键词		28.238	援助刊	da Bi	i Aan Tari 🔒	An to B	2. 久主宏
高型推索 间: Top-Kg	迪 索针算 "副用用平	自检察(点击这里	说素更多	-4-0 E	a var lød -	AU WAL	
<软件学报> 2014年04期	一种云	环境下的大数	数据Top-K查注 E 工小峰	向方法	ta.	人收藏設務	如两广告投》 010-6298299 244
【携娶】: Top-K查询在搜	索引擎、电子商务	等领域有着广泛	的应用.Top-K查到	向从海量数据中返	回最符合用户图	求的前K	和笑烟刊
个结果,主要目的是消除信息 4apReduce的特点,从数据3 具有良好的性质和3"展任。 【作者单位】: 中国人员2 【关键词】: 可p-K索切: 【基金】: 国家自然科学学 点专资将将形态全(20130004 (分录号]: TP311.13	3.过载带来的负面 3.分、数据符选等 大学信息学程; 云计算 Mar Redi 3.(613/3050,91 130001)	影响、大數指背景 方面对云环境下 1008 224008) 国家高生	下的Top-K童调,8 的大数描Top-K童 支术研究发展计划	6数据管理和分析 询问题进行深入和 (863)(2013AA01	逐方面带来断的 研究,实验结果表 3204) 高等学校(挑战,结合 。 明, 该方法 。 。 等士学科 。	电子技术应用 指安电子科技大学学校 材料工程 水声译丛 计算机工程与应用 一种科技大学学报(自然) 计算机工程 系统工程与电子技术
【正文快照】:							计关机构

Conclusions

 Graph Database is a Possible Way for RDF Knowledge Base Management.

<□> <圖> < ≧> < ≧> < ≧> < ≧ < つへぐ

Conclusions

 Graph Database is a Possible Way for RDF Knowledge Base Management.

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Subgraph Matching is a Strong Tool.

Conclusions

- Graph Database is a Possible Way for RDF Knowledge Base Management.
- Subgraph Matching is a Strong Tool.
- Using RDF repository, how to Provide Knowledge Services for Applications and Common Users?

▲□▶ ▲□▶ ▲□▶ ▲□▶ = ● のへで

Thank you!









▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ = 直 - 釣�?