D2R2: Disk-oriented Deductive Reasoning in a RISC-style RDF Engine

RuleML

Nov. 03, 2011

Mohamed Yahya <u>Martin Theobald</u> {myahya,mtb}@mpi-inf.mpg.de

Max-Planck Institute for Informatics, Saarbrücken, Germany

Background: Information Extraction

Stanford University - V Eile Edit View History	Subject		Predicate		Object		
W Stanford University - Wikipedia, the free enc +			Stanford University		type		Private University
WIKIPEDIA The Free Encyclopedia	Stanford University From Wikipedia, the free encyclopedia			hasPresident		J.L.Hennessy	
Main page Contents Featured content	"Stanford" redirects here. For other uses, see Stanford (disambiguation). The Leland Stanford Junior University, commonly referred to as Stanford University or Stanford, private research university on an 8,180-acre (3,310 ha) campus located near Palo Alto, California. It is in the previous clarp of the standard st				has	Students	15,319
Current events Random article Donate to Wikipedia	northwest of San Jose and 37 miles (60 km Leland Stanford, a Californian railroad type Letand Stanford, Jr., who gled of typenol type	;		foundedBy		L.Stanford	
 Interaction Help About Wikipedia Community portal 	as a coeducational and nondenominationa death and after much of the campus was o War II, Provost Frederick Terman supporte local industry in what would become know			foundedIn		1891	
Recent changes Contact Wikipedia	accelerator, was one of the original four Al university in computer science, mathemat faculty, staff, and alumni have won the Nol						
 Print/export Languages العروبة वारवा विद्यारा Ennapyckaa Catala Catala Catala Česky Dansk Deutsch Español Español	winners for a single institution. Stainford companies including Cisco Systems, Goo Rambus, Silicon Graphics, Sun Microsyst The university is organized into seven sch Earth Sciences as well as professional sc Stanford has a student body of approxima is a founding member of the Association o research funding and \$13.8 billion in endow Stanford competes in 34 varsity sports an Pacific-12 Conference. Stanford's athletic 1995. ^[9] In the 2008 Summer Olympics in more than any other university in the Unite Contents [hde] 1 History 1.1 Origins 1.1.1 Coeducation 1.1.2 Early finances 1.2 20th century 1.2.1 Football 1.2.2 Hoover Institution 1.3 Post 1945 1.3.1 Biology 1.3.2 High tech 1.3.3 Physics	cuty and adumin have founded many prominent technology gle, Hewlett-Packard, LinkedIn, Netscape Communications, ems, Varian Associates, and Yahool. ^[7] bools including academic schools of Humanities and Sciences a hools of Business, Education, Engineering, Law, and Medicine. Ietly 6,900 undergraduate ^[6] and 8,400 graduate students. ^[6] Sta f American Universities and in 2010 managed US\$1.15 billion in wment support, with \$21.4 billion in consolidated net assets. ^[6] d is one of two private universities in the NCAA Division I-A program has won the NACDA Directors' Cup every year since Beijing, Stanford athletes won 25 medals, including 8 gold med id States. ^[10]	Seal of S Motto Motto in English Motto in English Established Type Endowment Provost Academic staff Students Undergraduates Location Campus Athletic nickname Colors	tanford University Die Luft der Freiheit weh (German) ^[1] The wind of freedom blows ^[1] 1891 ^[2] Private US \$13.8 billion ^[2] John L. Hennessy John Echemendy 1,910 ^[4] 15,319 6,878 ^[6] 8,441 ^[6] Stanford, California, United States Suburban, 8,180 acres (33.1 km ²) ^[6] Cardinal Cardinal	t .::		

D2R2: Disk-oriented Deductive Reasoning in a RISC-style RDF Engine

YAGO Knowledge Base

- Combine knowledge from WordNet & Wikipedia
- Additional
 Gazetteers

 (geonames.org)
- Part of the Linked-Data cloud



• YAGO2: >30 Mio RDF facts in base relations, >95% precision



Motivation (1/2)





Motivation (2/2)







Motivation (2/2)



Outline

- Recursive query processing:
 - QSQR
 - Chaining & dynamic subquery scheduling
- Adding deductive reasoning to an RDFengine
- Experiments



Outline

 Recursive query processing:
 – QSQR



QSQR

- Query(q):
 _ Q = ? ← likes(John,?y)

- Rules(D):
 - $r_1: \underline{likes}(?x,?y) \leftarrow \underline{friend}(?x,?y)$
 - r_2 : <u>friend(?x,?y</u>) \leftarrow marriedTo(?x,?y)
 - r_3 : <u>friend(?x,?y</u>) \leftarrow <u>likes(?x,?z</u>) \land praises(?z,?y)

• Facts(D):

- f₁: marriedTo(John,Mary)
- f₂: praises(Mary,Tom)



Algorithm 1: QSQR(D, q)Input: A Datalog program D and an intensional query qInput: The global ans_ and input_ relations Output: All answers for q 1 begin 2 Set all ans_ relations to be empty Set (R^{γ}, J) to be the generalized query corresponding to q 3 4 repeat Set all input_ relations to be empty 5 Call QSQR_EVAL_GENERALIZED($D, (R^{\gamma}, J)$) 6 until Until no ans_ relation has changed in the last iteration; 7 return All answers for g by performing a selection on ans \mathbb{R}^{γ} using. 8

> [Abiteboul, Hull, Vianu: Foundations of Databases, 1995]

<john,?></john,?>	<john,?></john,?>				
ans_friend	ans_likes				
<john, mary=""></john,>	<john, mary=""></john,>				
<john, tom=""></john,>	<john, tom=""></john,>				

input likes input friend

RuleML, 03.11.2011

D2R2: Disk-oriented Deductive Reasoning in a RISC-style RDF Engine

Outline

- Recursive query processing:
 - QSQR
 - Chaining & dynamic subquery scheduling



QSQR + Extensions

- **1. Chaining**: make use of the underlying engine's index structures & ability to perform joins
- 2. Dynamic scheduling: next set of atoms to be evaluated in a subquery is determined after current set is evaluated

QSQR + Extensions – Chaining

 $I_1(?x,?y) \leftarrow I_2(?x,?y) \wedge E_1(?x,?z) \wedge E_2(?y,?z)$

Two ways to evaluate: $? \leftarrow I_1(a,?y)$



1. Grouping: better utilization of query optimizer a. $E_1(a,?z) \land E_2(?y,?z)$ b. $I_2(a,?y)$ for every ?y More choices: (SMJ, HJ, NLJ) Merging, Hashing

QSQR + Extensions – Dynamic Scheduling

 $\mathbf{I}_{1}(?x,?y) \leftarrow \mathbf{I}_{2}(?x,?y) \wedge \mathbf{E}_{1}(?x,?z) \wedge \mathbf{E}_{2}(?y,?z)$

Just because the rule places I₂ first, doesn't mean it has to be evaluated first:

Declarative paradigm: What vs. How

Extend QSQR to support *dynamically* deciding which (set of) literals to evaluate next.



Outline

- Recursive query processing:
 - QSQR
 - Chaining & dynamic subquery scheduling
- Adding deductive reasoning to an RDFengine



RDF-3X Engine [T. Neumann et al.: VLDB'08, SIGMOD'09]

- No-tuning RISC, versioning, online updates, transactions
- Aggressive indexing (15 SPO permutations & projections)
- Fast DP-based join-order optimization for up to 20-30 joins!
- Aggressive sideways information passing



RDF-3X (1/2)

 15 subsets & permutations of SPO attributes stored in exhaustive B⁺-tree indexes:



- Disk-oriented statistics, relying on the indexes above.
- Index-only tables: B⁺-tree leafs contain data & statistics.



RDF-3X (2/2)

• Operators see integers:



Dictionary



- Sideways information passing (SIP):
 - Operators communicate across query plan to skip pages.



RDF-3X Integration (1/2)

Our query patterns do not always match what RDF-3X expects.

- **Small subqueries** (single literal/atom) are very common:
 - » Bypass RDF-3X query optimizer & employ own precompiled plans. I₁(?x,?y) ← I₂(?x,?y) ∧ E₂(?y,?z) ∧ I₃(?z,?y)
- **Predicates** are always given in our context:
 - » Reduce the number of plans RDF-3X considers: only POS, PSO, PO and PS needed (others still used for statistics).

RDF-3X Integration (2/2)

- During recursion, same pages can be requested multiple times:
 - » Add caching on top of compressed indexes.





Outline

- Recursive query processing:
 - QSQR
 - Chaining & dynamic subquery scheduling
- Adding deductive reasoning to an RDFengine
- Experiments



Experiments (1/5)

- Datasets:
 - YAGO1 22 Mio RDF facts (3.8 GB)
 - \rightarrow 16 partly recursive rules
 - LUBM Scale factor 1: 100k RDF facts (5.7 MB)
 - → single recursive rule (*subOrganizationOf*)
- D2R2 implemented in C++ on top of RDF-3X
- Other systems:
 - IRIS Reasoner (incl. PostgreSQL 8.4 backend)
 - Jena Semantic Web Framework (incl. TDB 0.8.7 backend)

Experiments (2/5)



YAGO setting, no rules

- QE1 = ? \leftarrow bornIn(Tipper Gore; ?y);
- QE3 = ? ← isMarriedTo(Al Gore; ?y); bornIn(?y; ?z);
- QE6 = ? ← actedIn(Schwarzenegger; ?y); actedIn(?x; ?y); bornIn(?x; ?z);

Experiments (3/5)



R1: bornIn(?x; ?z) \leftarrow isCitizenOf(?x; ?y); locatedIn(?z; ?y); livesIn(?x; ?z);

R2: bornIn(?c; ?y) \leftarrow <u>livesIn(</u>?x; ?y); <u>livesIn(</u>?z; ?y); <u>isMarriedT</u>o(?x; ?z); hasChild(?x; ?c); hasChild(?z; ?c);

R3: <u>livesIn(?y; ?z)</u> \leftarrow <u>isMarriedT</u> o(?x; ?y); <u>livesIn(?x; ?z);</u>

R4: <u>isMarriedTo(</u>?x; ?y) ← hasChild(?x; ?z); hasChild(?y; ?z); notEquals(?x; ?y);

•••



RuleML, 03.11.2011

D2R2: Disk-oriented Deductive Reasoning in a RISC-style RDF Engine



Experiments (4/5)



- Q13: ? ← ancestor(Edward IV of England; ?x);

```
R1: ancestor(?x,?y) \leftarrow parent(?x,?y);
R2: ancestor(?x,?y) \leftarrow ancestor(?x,?z); parent(?z,?y);
```

...

D2R2: Disk-oriented Deductive Reasoning in a RISC-style RDF Engine



Experiments (5/5)



- R2: <u>subOrganizationOf(</u>?x; ?z) ← <u>subOrganizationOf(</u>?x; ?y); <u>subOrganizationOf(</u>?y; ?z);
- R3: hasAlumnus(?x; ?y) ← degreeF rom(?y; ?x);
- ...

• QL1: ? ← type(?x; GraduateStudent); takesCourse(?x; http://www:Department0:University0:edu=GraduateCourse0);

Conclusions

- Setting
 - Large disk-resident RDF collections & deductive reasoning
- Recursive queries
 - QSQR with extensions: *chaining* & *dynamic* scheduling
- Facts store
 - RDF-3X with modifications for our setting
- <u>Current work</u>
 - Distributed reasoning
 - Lineage tracing & probabilistic inference ("soft rules")



Thank You!

Questions?

RuleML, 03.11.2011

D2R2: Disk-oriented Deductive Reasoning in a RISC-style RDF Engine

