Scalable RDF Data Management



... with a Touch of Uncertainty

Martin Theobald

University of Luxembourg

Faculty of Sciences, Technology & Communication

Joint work with:

 Hernán Blanco, Maximilian Dylla, Sairam Gurajada, Maarten Van den Heuvel, Iris Miliaraki, Dat Ba Nguyen











Jeffrey Ullman

From Wikipedia, the free encyclopedia

Jeffrey David Ullman (born November 22, 1942) is a renowned computer scientist. His textbooks on compilers (various editions are popularly known as the Dragon Book), theory of computation (also known as the <u>Cinderella book</u>), data structures, and databases are regarded as standards in their fields.

Contents [hide]	
1 Early life & Career	
2 Books	
3 References	
4 External links	
Early life & Car	eer
Ullman received a Ba	chelor of Science

Ullman received a Bachelor of Science degree in Engineering Mathematics from Columbia University in 1963 and his Ph.D. in Electrical Engineering from Princeton University in 1966. He then worked for several years at <u>Bell Labs.</u> From 1969 to 1979 he was a professor at Princeton. Since 1979 he has been a professor at Stanford University, where he is currently the Stanford W.

	Jeffrey D. Ullman
Born	November 22, 1942 (age 69)
Citizenship	American
Nationality	American
Institutions	Stanford University
Alma mater	Columbia University, Princeton University
Doctoral advisor	Arthur Bernstein, Archie McKellar
Doctoral students	Surajit Chaudhuri, Kevin Karplus, David Maier, Harry Mairson, Alberto O. Mendelzon, Jeffrey F. Naughton, Anand Rajaraman, Yehoshua Sagiy, Mibelis Vannakakis
Known for	database theory, database systems, formal language theory
Notable awards	Fellow of the Association for Computing Machinery (1994), ACM SIGMOD Contributions Award (1996), ACM SIGMOD Best Paper Award (1996), Karl V. Karlstrom outstanding educator award (1998), Knuth Prize (2000), ACM SIGMOD Edgar F. Codd Innovations Award (2006), ACM SIGMOD Test of Time Award (2006), IEEE John von Neumann Medal (2010)

Ascherman Professor of Computer Science (Emeritus). In 1995 he was inducted as a Fellow of the Association for Computing Machinery and in 2000 he was awarded the Knuth Prize. Ullman is also the co-recipient (with John Hopcroft) of the 2010 IEEE John von Neumann Medal, "For laying the foundations for the fields of automata and language theory and many seminal contributions to theoretical computer science."^[1]

[edit]

Ullman's research interests include database theory, data integration, data mining, and education using the information infrastructure. He is one of the founders of the field of database theory, and was the doctoral advisor of an entire generation of students who later became leading database theorists in their own right. <u>He was the Ph.D.</u> advisor of <u>Sergey Brin</u>, one of the co-founders of <u>Google</u>, and served on Google's technical advisory board. He is currently the <u>CEO of Gradiance</u>.

Books

Information Extraction

DBpedia/YAGO et al.

bornOn(Jeff, 09/22/42) gradFrom(Jeff, Columbia) gradFrom(Jeff, Princeton) hasAdvisor(Jeff, Arthur) hasAdvisor(Surajit, Jeff) knownFor(Jeff, Theory)

>120 M facts for YAGO3 (from Wikipedia infoboxes)

New fact candidates

author(Jeff, Drag_Book) [0.6] author(Jeff, Cind_Book) [0.8] worksAt(Jeff, Bell_Labs) [0.5] hasAdvisor(Sergej, Jeff) [0.7] type(Jeff, ACM_Fellow) [0.5] type(Jeff, CEO) [0.3]

>100's M additional facts (from Wikipedia free-text)

[edit]

Database Systems: The Complete Book (with H. Garcia-Molina and J. Widom), Prentice-Hall, Englewood Cliffs, NJ,

Linked-Open-Data Cloud



 Web Apps Examples Web Apps Examples ussuming "University of Luxembourg" is a university Use "Luxembourg" as a country instead uput interpretation: Université du Luxembourg year founded tesult: 1848 time from today: The year 1848 was 169 years ago. 	⊃⊄ Randor Open code <i>(</i>
assuming "University of Luxembourg" is a university Use "Luxembourg" as a country instead apput interpretation: Université du Luxembourg year founded esult: 1848 ime from today: The year 1848 was 169 years ago. roperties: 1848 is a leap year.	Open code 🧲
nput interpretation: Université du Luxembourg year founded tesult: 1848 ime from today: The year 1848 was 169 years ago. troperties: 1848 is a leap year.	Open code 🤇
Université du Luxembourg year founded esult: 1848 ime from today: The year 1848 was 169 years ago.	Open code 🤇
lesult: 1848 ime from today: The year 1848 was 169 years ago. roperties: 1848 is a leap year.	G
1848 ime from today: The year 1848 was 169 years ago. roperties: 1848 is a leap year.	ē
ime from today: The year 1848 was 169 years ago. roperties: 1848 is a leap year.	G
1848 is a leap year.	
Saturday, January 1, 1848 to Sunday, December 31, 1848	More calendars
lotable events in 1848: January 24: California Gold Rush	More
February 2: Treaty of Guadalupe Hidalgo	
February 21: Communist Manifesto published	
July 4: Marx and Engels publish their "Communist Manifesto"	
December 2: Franz Josef I becomes Emperor of Austria and King of Hungary	
alendar:	
January February	

27 28 29

23

30 31

24 25 26 27 28 29

Wolfram Alpha

The "Computational Knowledge Engine"

- Fully implemented in Wolfram-Mathematica
- 10 trillion+ facts
- 50,000+ algorithms and statistical analyses
- 5,000+ templates for visualization and layouts
- 1,000+ domain-specific linguistic analyses

http://www.wolframalpha.com/

IBM Watson: Deep Question Answering

- William Wilkinson's "An Account of the Principalities of Wallachia and Moldavia" inspired this author's most famous novel
- This town is known as "Sin City" & its downtown is "Glitter Gulch"
- As of 2010, this is the only former Yugoslav republic in the EU
- 99 cents got me a 4-pack of Ytterlig coasters from this Swedish chain
- U.S. City: largest airport is named for a World War II Hero; its second largest for a World War II Battle



Question classification & decomposition



Knowledge back-ends

D. Ferrucci et al.: Building Watson: An Overview of the DeepQA Project. Al Magazine, Fall 2010.

https://www.ibm.com/watson/

RDF-Centered Research Topics

Information Extraction

[SIGMOD'09, WebDB'10, PODS'10, WSDM'11, CIKM'12, CLEF/INEX'11/'12, LDOW'14, TACL'16]

Uncertain RDF Data & Probabilistic Databases

[ICDE'08, VLDB-J'08, SSDBM'10, BTW'11, CIKM'11, ICDE'13, PVLDB'14, VLDB PhD Workshop'15]

 Scalable RDF Indexing & SPARQL Query Processing

[SIGMOD'14, SWIM'14, SIGMOD'16]

"David played for manu, real, and la galaxy. His wife posh performed with the spice girls."









"David played for manu, real, and la galaxy. His wife posh performed with the spice girls."



 J-NERD jointly recognizes and disambiguates named entities with respect to a background knowledge base such as YAGO.









- Probability distribution over possible tokens x and combined NER/D labels y
- Probabilistic inference: find the most likely labels y, given the observed tokens x
- Viterbi algorithm (dynamic programming) for fast and exact inference

$$p(\mathbf{x}, \mathbf{y}) = \frac{1}{Z} \prod_{A} \mathcal{F}_A(\mathbf{x}_A, \mathbf{y}_A)$$

$$\mathbf{y}^* = \arg \max_{\mathbf{y}} p(\mathbf{y} \,|\, \mathbf{x})$$



- Probability distribution over possible tokens x and combined NER/D labels y
- Probabilistic inference: find the most likely labels y, given the observed tokens x

$$p(\mathbf{x}, \mathbf{y}) = \frac{1}{Z} \prod_{A} \mathcal{F}_{A}(\mathbf{x}_{A}, \mathbf{y}_{A})$$

- $\mathbf{y}^* = \arg \max_{\mathbf{y}} p(\mathbf{y} \,|\, \mathbf{x})$
- General factor graphs: MCMC-style sampling for approximate inference



Evaluation on the CoNLL newswire collection with YAGO2 groundtruth annotations (1,244 labeled articles)

_	Method	Prec	Rec	\mathbf{F}_1
_	P-NERD	80.1	75.1	77.5
_	J-NERD	81.9	75.8	78.7
_	AIDA-light	78.7	76.1	77.3
_	TagMe	64.6	43.2	51.8
) _	SpotLight	71.1	47.9	57.3
,				

Ultimate PhD Challenge (I)

"Paris Hilton stays in the Hilton in Paris."

	🗎 gate.d5.mpi-inf.mpg.de	C	0 1
Max-Planck-Institut für Informatik: AIDA		AIDA Web interface (aida)	+
Disambiguation Method: prior prior+sim prior+sim+coherence	Paris Hilton stays in th	ne Hilton in Paris.	
Parameters			Discusto
Prior-Similarity-Coherence balancing ratio:			Disambiguate
	Input Type:TEXT Overall runtin	me: 2032 ms	
Ambiguity degree 5	Paris Hilton [Paris Hilton] Stay	rs in the Hilton [Paris Hilton]	in Paris [Paris Hilton].
Coherence Measure:	Run Information Gra	ph Removal Steps	
Entities Type Filters:	• 0: Paris Hilton (solve	ed by local sim. only)	
Stanford NER Manual You can manually tag the mentions by putting them between [[and]]. HTML Tables are automatcially disambiguated in the manual mode.	26: Hilton 36: Paris chunkid: 3A8E68D0C750C09CC098FA92743F/	AC571494247460213_singlechunk	
Fast Mode: Enabled Examples YAGOTypes		Types tag cloud	d Focused Types tag cloud

- All of the current NED tools (incl. AIDA, J-NERD, Spotlight, TagMe) get this sentence wrong!
- Humans (usually) get it right, though.

RDF-Centered Research Topics

Information Extraction

[SIGMOD'09, WebDB'10, PODS'10, WSDM'11, CIKM'12, CLEF/INEX'11/'12, LDOW'14, TACL'16]

Uncertain RDF Data & Probabilistic Databases

[ICDE'08, VLDB-J'08, SSDBM'10, BTW'11, CIKM'11, ICDE'13, PVLDB'14, VLDB PhD Workshop'15]

 Scalable RDF Indexing & SPARQL Query Processing

[SIGMOD'14, SWIM'14, SIGMOD'16]

Probabilistic Database

A probabilistic database D^p (compactly) encodes a probability distribution over a finite set of deterministic database instances D_i .

(I) C	0 ₁ : 0.42	D	₂ : 0.18	C	0 ₃ : 0.28	D ₄ : 0.12
works	At(sub, obj)	works	At(sub, obj)	works	At(sub, obj)	worksAt(sub, obj)
Jeff	Stanford	Jeff	Stanford	Jeff	Princeton	
Jeff	Princeton					

Special Cases:

(II) **D**^p tuple-independent

worksAt(sub, obj)		р
Jeff	Stanford	0.6
Jeff	Princeton	0.7

(III) **D**^p block-independent

worksAt(sub, obj)		р
Jeff	Stanford	0.6
	Princeton	0.4

Note:

(I) and (II) here are equivalent;(II) and (III) not!

Query Answering Problem: ("Marginal Probabilities" of Query Answers) Run query **Q** against each instance D_i ; for each answer tuple t_j , $P(t_i)$ is the sum of the probabilities of all instances D_i where t_i exists.

Flashback: Stanford Trio System

[Widom: CIDR'05]

Uncertainty-Lineage Databases (ULDBs)

- 1. Alternatives
- 2. '?' (Maybe) Annotations
- 3. Confidence values
- 4. Lineage

Trio's Data Model

[Widom: CIDR'05]

1. Alternatives: uncertainty about value



Trio's Data Model

[Widom: CIDR'05]

- 1. Alternatives
- 2. '?' (Maybe): uncertainty about presence



Trio's Data Model

[Widom: CIDR'05]

1. Alternatives

2. '?' (Maybe) Annotations

3. Confidences: weighted uncertainty



So Far: Data Model is <u>Not</u> Closed

[Widom: CIDR'05]



Drives (*person*, *car*)

Jimmy, Toyota || Jimmy, Mazda

Billy, Honda || Frank, Honda

Hank, Honda

Suspects = **T**_{person}(Saw ⋈ Drives)



Example with Lineage



$$\begin{split} \lambda(31) &= (11,2) \land (21,2) \\ \lambda(32,1) &= (11,1) \land (22,1); \ \lambda(32,2) = (11,1) \land (22,2) \\ \lambda(33) &= (11,1) \land 23 \end{split}$$

Example with Lineage

ID	Saw	(witness, <i>car</i>)
11	Cathy	Honda Mazda

ID	Drives (<i>person</i> , <i>car</i>)
21	Jimmy, Toyota Jimmy, Mazda
22	Billy, Honda Frank, Honda
23	Hank, Honda

Suspects = **T**_{person}(Saw ⋈ Drives)



Operational Semantics



Completeness: any (finite) set of possible instances can be represented

(will be coming back to this subtlety again later...)

Summary on Trio's Data Model

Uncertainty-Lineage Databases (ULDBs)

- 1. Alternatives
- 2. '?' (Maybe) Annotations
- 3. Confidence values
- 4. Lineage

Theorem: ULDBs are **closed** and **complete**.

Formally studied properties like *minimization*, *equivalence*, *approximation* and *membership* **based on lineage**. [Benjelloun, Das Sarma, Halevy, Widom, Theobald: VLDB-J. 2008]



.. back to Wikipedia

rack Obama citizenship conspiracy th... 🕂

📓 https://en.wikipedia.org/wiki/Barack_C 🏫 🔻 C 🛛 💈 🕶 Google

Born in Kenya [edit source | edit beta]

Some opponents of Obama's presidential eligibility claim that he was born in Kenva, and was therefore not born a United States citizen.

_ D X

•

م

Whether Obama having been born outside the U.S. would have invalidated his U.S. citizenship at birth is debated. Andrew Malcolm, of the *Los Angeles Times*, has argued that Obama would still be eligible for the presidency, irrespective of where he was born, because his mother was an American citizen, saying that Obama's mother "could have been on Mars when wee Barry emerged and he'd still be American."^[59] A contrary view is promoted by UCLA Law Professor Eugene Volokh, who has said that in the hypothetical scenario that Obama was born outside the U.S., he would *not* be a natural-born citizen, since the then-applicable law would have required Obama's mother to have been in the U.S. at least "five years after the age of 14", but Ann Dunham was three months shy of her 19th birthday when Obama was born.^[60]

Obama's paternal step-grandmother's version of events [edit source | edit beta]

An incorrect but popularly reported claim is that his father's stepmother, Sarah Obama, told Anabaptist Bishop Ron McRae in a recorded transatlantic telephone conversation that she was present when Obama was born in Kenya.

bornIn(Barack, Hawaii)

min Baraon, Ronya

Deductive Database:

Datalog, Core of SQL &

(**Soft**) Deduction Rules vs. (**Hard**) Consistency Constraints

- People may live in more than one | RDF/S, OWL2-RL, etc. livesIn(x,y) \Leftarrow marriedTo(x,z) \land livesIn(z,y) [0.5]
- ▶ People are not born in different place of a different data bornIn(x,y) ∧ bornIn(x,z) ⇒ y=z

 More General FOL

 born0n(x,y) ∧ born0n(x,z) ⇒ y=z

 Constraints:
- People are not married to more that (at the same time, in most countries?)
 marriedTo(x,y,t₁) ∧ marriedTo(x,z,t₂) ∧ yrz
 ⇒ disjoint(t₁,t₂)

Deductive Grounding w/ Lineage

(SLD Resolution in Datalog/Prolog)

[Yahya, Theobald: RuleML'11, Dylla, Miliaraki, Theobald: ICDE'13]



Rules

hasAdvisor(x,y) <

worksAt(y,z)

 \Rightarrow graduatedFrom(x,z)

graduatedFrom(x,y) ^
graduatedFrom(x,z)

⇒ y=z

Base Facts

graduatedFrom(Surajit, Princeton) [0.7] graduatedFrom(Surajit, Stanford) [0.6] graduatedFrom(David, Princeton) [0.9] hasAdvisor(Surajit, Jeff) [0.8] hasAdvisor(David, Jeff) [0.7] worksAt(Jeff, Stanford) [0.9] type(Princeton, University) [1.0] type(Stanford, University) [1.0] type(Jeff, Computer_Scientist) [1.0] type(Surajit, Computer_Scientist) [1.0]

Lineage & Possible Worlds



[Das Sarma, Theobald, Widom: ICDE'08, Dylla, Miliaraki, Theobald: ICDE'13]

1) Deductive Grounding

- Top-down Datalog evaluation
- Plus tracing the lineage of individual query answers

2) Lineage DAGs

- Grounded soft & hard rules
- Base facts with confidences

3) Probabilistic Inference

→ Compute marginals:

P(Q): sum up the probabilities of all possible worlds that entail the query answers

P(Q|H): drop "impossible worlds"

	P(Q ₂ H)=0.2664 / <mark>0.412</mark> = 0.6466	0.0784 / <mark>0.412</mark> = 0.1903	P(Q ₁ H)=	4 I 4	=0.078	P(Q ₁) P(Q ₂)
	P(W)	Q ₂ :	D:0.9	C:0.8	B:0.6	A:0.7
\sum	0.7x0.6x0.8x0.9 = 0.3024	0	1	1	1	1
	0.7x0.6x0.8x0.1 = 0.0336	0	0	1	1	1
	= 0.0756	0	1	0	1	1
	= 0.0084	0	0	0	1	1
	= 0.2016	0	1	1	0	1
$\overline{)}$	= 0.0224	0	0	1	0	1
> 0.0784	= 0.0504	0	1	0	0	1
	= 0.0056	0	0	0	0	1
)	0.3x0.6x0.8x0.9 = 0.1296	1	1	1	1	0
	0.3x0.6x0.8x0.1 = 0.0144	1	0	1	1	0
≻ 0.2664	0.3x0.6x0.2x0.9 = 0.0324	1	1	0	1	0
	0.3x0.6x0.2x0.1 = 0.0036	1	0	0	1	0
) > 1.0	0.3x0.4x0.8x0.9 = 0.0864	1	1	1	0	0
0.41	= 0.0096	0	0	1	0	0
	= 0.0216	0	1	0	0	0
	= 0.0024		0	•	0	0

Dichotomy of Queries

[Suciu & Dalvi: SIGMOD'05 Tutorial on "Foundations of Probabilistic Answers to Queries"]

A probabilistic database D^p (compactly) encodes a probability distribution over a finite set of deterministic database instances D_i .

Is there any professor who works at a university that is located in CA?
Q() :- isProfessor(pers), worksAt(pers,uni), located(uni, CA)



<u>Theorem</u>: The query answering problem for the above join query over a tuple-independent probabilistic database is **#P-hard**.

Inference in Probabilistic Databases

Safe query plans [Dalvi & Suciu: VLDB-J'07+J-ACM'12]

Can propagate confidences along with relational operators.

Read-once functions [Sen et al.: PVLDB'10; Olteanu & Huang: SUM'08]

 Can factorize Boolean formula (in polynomial time) into read-once form, where every variable occurs at most once.

Knowledge compilation [Olteanu et al.: ICDE'10; ICDT'11; VLDB-J'13]

Can compile Boolean formula into a decision diagram (OBDD/SDD), such that inference resolves to *independent-and* and *independent-or* operations over the decomposed formula.

Top-*k* **pruning** [Ré, Davli & Suciu: ICDE'07; Karp, Luby & Madras: J-Alg.'89; Olteanu & Wen: ICDE'12]

- Can return top-k answers based on *lower* and *upper bounds*, even without knowing their exact marginal probabilities.
- <u>Multi-Simulation</u>: run multiple Markov-Chain-Monte-Carlo (MCMC) simulations in parallel.

Top-k Ranking by Marginal Probabilities



Bounds for First-Order Formulas

[Dylla, Miliaraki, Theobald: ICDE'13]

Theorem 1:

Given a (partially grounded) **first-order lineage** formula Φ :

 $\Phi(Q_2) = B \lor \exists y \text{ gradFrom}(S,y)$

 Lower bound P_{low} (for all query answers that can be obtained from grounding Φ): Substitute ∃y gradFrom(S,y) with *false* (or *true* if negated).

 $\mathsf{P}_{\mathsf{low}}(\mathsf{Q}_2) = \mathsf{P}(\mathsf{B} \lor \mathsf{false}) = \mathsf{P}(\mathsf{B}) = \mathbf{0.6}$

 Upper bound P_{up} (for all query answers that can be obtained from grounding Φ): Substitute ∃y gradFrom(S,y) with *true* (or *false* if negated).

```
P_{up}(Q_2) = P(B \lor true) = P(true) = 1.0
```

Proof: (sketch)

Substitution of a subformula with *false* reduces the number of models *(possible worlds)* that satisfy Φ ; substitution with *true* increases them.

Convergence of Bounds

[Dylla,Miliaraki,Theobald: ICDE'13]

Theorem 2:

Let Φ_1, \ldots, Φ_n be a series of first-order lineage formulas obtained from grounding Φ via SLD resolution, and let ϕ be the *propositional lineage* formula of an answer obtained from this grounding procedure.

Then rewriting each Φ_i according to Theorem 1 into $P_{i,low}$ and $P_{i,up}$ creates a **monotonic series** of lower and upper bounds that **converges to P(\phi)**.

$$\begin{array}{l} 0 = \mathsf{P}(\mathsf{false}) \leq \mathsf{P}(\mathsf{B} \lor \mathsf{false}) = 0.6 \leq \mathsf{P}(\mathsf{B} \lor (\mathsf{C} \land \mathsf{D})) = 0.888 \\ \\ \leq \mathsf{P}(\mathsf{B} \lor \mathsf{true}) = \ \mathsf{P}(\mathsf{true}) = 1 \end{array}$$

Proof: (sketch, via induction)

Substitution of *true* with a formula reduces the number of models that satisfy Φ ; substitution of *false* with a formula increases this number.

Top-k Stopping Condition

[Fagin et al.'01; Balke,Kießling'02; Dylla,Miliaraki,Theobald: ICDE'13]

"Fagin's Algorithm"

Maintain two disjoint queues:

Top-k sorted by P_{low} and *Candidates* sorted by P_{up}

Return the top-k queue at the t-th grounding step when:

 $\mathsf{P}_{i,\mathsf{low}}(\mathsf{Q}_k) \mid_{\mathsf{Q}_k \in \mathsf{Top-}k} \mathsf{P}_{i,\mathsf{up}}(\mathsf{Q}_j) \mid_{\mathsf{Q}_j \in \mathsf{Candidates}}$



Temporal-Probabilistic Database

[Wang,Yahya,Theobald: MUD'10; Dylla,Miliaraki,Theobald: PVLDB'13]



Example using the Allen predicate overlaps

Inference in Temporal-Probabilistic Databases





Inference in Temporal-Probabilistic Databases

[Wang,Yahya,Theobald: MUD'10; Dylla,Miliaraki,Theobald: PVLDB'13]

Derived Facts

teamMates(Beckham, teamMates(Beckham, Ronaldo, T₄) Zidane, T₅)

> teamMates(Ronaldo, Zidane, T₆)

Non-independent Independent

Closed and complete representation model (incl. lineage)

Temporal alignment is linear in the number of input intervals

Probabilistic inference per interval remains #P-hard

 Inference requires lineage decompositions, top-k pruning, or Monte Carlo approximations (Luby-Karp for DNF, MCMC-style sampling)

Ultimate PhD Challenge (II)



RDF-Centered Research Topics

Information Extraction

[SIGMOD'09, WebDB'10, PODS'10, WSDM'11, CIKM'12, CLEF/INEX'11/'12, LDOW'14, TACL'16]

Uncertain RDF Data & Probabilistic Databases

[ICDE'08, VLDB-J'08, SSDBM'10, BTW'11, CIKM'11, ICDE'13, PVLDB'14, VLDB PhD Workshop'15]

Scalable RDF Indexing & SPARQL Query Processing [SIGMOD'14, SWIM'14, SIGMOD'16]

RDF & SPARQL



Data complexity of core SPARQL: *polynomial* **Combined data & query complexity**: *exponential*

(same as SQL w/o recursion)

TriAD Architecture

RDF

Indexing



TriAD Architecture

SPARQL Query

Processing





 \rightarrow TriAD follows a very classical master-slave architecture; however with a direct (asynchronous) communication among all slaves at query time.

Locality-Based Graph Summarization: METIS



Min-k-Cut

For a desired amount of k evenly sized partitions, assign each node in the RDF data graph to exactly one partition, such that the number of cut edges among those partitions is minimized.

Summary Graph



- Drop all nodes and edges inside the partitions
- Keep only inter-partition edges
- Introduce self-loop edges for intra-partition edges

Querying the Summary Graph

Global Dictionary:	
Barack_Obama	$\rightarrow P_{I}$
USA	$\rightarrow P_{I}$
Lady_Gaga	$\rightarrow P_2$
Peace_Nobel_Prize	$\rightarrow P_4$

...



Querying the Summary Graph



- Summary graph guarantees no false negatives (i.e., "missed results"); the subsequent processing of the query against the pruned data graph also ensures no false positives.
- Facilitates join-ahead pruning by skipping over irrelevant partitions.

Example Query Plan



- A copy of the same query plan is shipped to all slaves:
 - > DIS operators (leafs) are augmented with **locality** and **pruning information**.
 - 6 SPO permutations allow the usage of **DMJ op's at the first level of joins**.

Distributed & Multithreaded Query Execution



Experiments

TriAD is implemented in C++ using GCC 4.4, Boost-1.5 & MPICH2.

All experiments were run on a proprietary cluster with 32 x 48 GB RAM, 2 quad-core XENON CPUs and a 1GBit Ethernet connection.

LUBM – Lehigh University Benchmark

Scale Factor 160: 28 Mio RDF triples \rightarrow 16 GB data \rightarrow 3 GB index Scale Factor 10240: 1.8 Bio RDF triples \rightarrow 730 GB data \rightarrow 150 GB index

BTC – Billion Triples Challenge (2012)

DBpedia/Yago/Freebase: 1.4 Bio RDF triples \rightarrow 231 GB data \rightarrow 130 GB index

• WSDTS – Waterloo SPARQL Diversity Test Suite

Scale Factor 1000: 109 Mio RDF triples \rightarrow 15 GB data \rightarrow 9.1 GB index

<u>9 Competitors:</u> RDF-3x (MPII), MonetDB (U-Amsterdam), BitMat (Rensselaer Polytech), TripleBit (U-Huazhong/U-Georgia), Hadoop-RDF-3x (Yale), Apache Hadoop / Spark (UC Berkeley), SHARD (open-source), Trinity.RDF (MSR)

Benchmark Results

	TriAD	TriAD-SG	Trinity.RDF	SHARD	H-RDF-3X		4store		RDF-3X		BitMat	
		(200K)			(cold)	(warm)	(cold)	(warm)	(cold)	(warm)	(cold)	(warm)
Q1	7,631	2,146	12,648	6.9E5	2.3e6	1.7e5	aborted	aborted	1.9e6	1.8e6	17,339	11,295
Q2	1,663	2,025	6,018	2.1e5	5.3E5	4,095	1.1E5	15,113	6.3E5	17,835	2.4e5	1.8e5
Q3	4,290	1,647	8,735	4.7e5	2.2e6	1.3e5	aborted	aborted	1.7e6	1.7e6	8,429	2,679
Q4	2.1	1.3	5	3.9E5	166	1	1,903	12	243	3	aborted	aborted
Q5	0.5	0.7	4	97,545	85	1	2,429	12	99	1	472	338
Q6	69	1.4	9	1.8E5	5.8e5	23,440	3,572	9	913	287	7,796	5,377
Q7	14,895	16,863	31,214	3.9E5	2.3e6	2.1e5	aborted	aborted	6.5E5	46,262	71,157	36,905
Geo-												
Mean	249	106	450	3.0E5	91,378	2,406	-	-	31,345	2,991	-	-

LUBM-10240: Query Processing Times in Milliseconds (ms)

	#Results	TriAD	TriAD-SG (200K)	H-R (cold)	ADF-3X (warm)	RDF (cold) (F-3X warm)	_			
Q1	1	1.5	0.3	49	6	297	4	-			
Q2 Q3 Q4	1 1 0	61 1 0.6				#Slaves	(Geo	L1-L5 Mean)	S1-S7 (GeoMean)	F1-F5 (GeoMean)	C1-C3 (GeoMean)
Q5 Q6 Q7	5 0 0	$ \begin{array}{cccc} 5 & 51 \\ 0 & 0.5 \\ 0 & 50 \end{array} $	Tr Tr	TriAD TriAD-SG(75K) TriAD SHARD RDF-3X (cold)		1		2 8	2 4	94 35	494 767
Q8 Geo Mean		128 7.4	SF RI			5 5 1		2 3.2E5 10,066	3 5.8E5 167	29 7.1E5 1,749	270 7.7E5 6.610
BTC: Query Pro			/ Prc RI M	RDF-3X (warm) MonetDB (cold) MonetDB (warm)		1 1 1		18 3530 171	2 10,459 744	41 timeout timeout	354 timeout timeout

WSDTS-1000: Query Processing Times (ms)

Ultimate PhD Challenge (III)

From Map & Reduce



• over Synchronous Dataflows to Asynchronous Dataflows!



Summary

Information Extraction

- Natural-Language
 Processing & Understanding
- Named-Entity Recognition
 & Disambiguation
- Extraction of N-Ary Relations
- Knowledge-Graph Construction, Integration & Maintenance

Uncertain Data

- Probabilistic & Temporal Data(base) Models
- Data Integration & Cleaning
- Model- & Dissociationbased Bounds
- Scalable Probabilistic
 Inference

Big Data

- Big Data Analytics
- Distributed Graph Engines
- Real-Time Dataflows & Stream Processing
- Message Passing & Asynchronous Protocols

References

- Sairam Gurajada, Martin Theobald: Distributed Set Reachability. SIGMOD 2016
- Dat Ba Nguyen, Martin Theobald, Gerhard Weikum: Joint Named Entity Resolution and Disambiguation with Rich Linguistic Features. TACL Vol. 4, 2016
- Hernán Blanco: Scaling Probabilistic Databases. VLDB PhD Workshop 2015
- Dat Ba Nguyen, Johannes Hoffart, Martin Theobald, Gerhard Weikum: AIDA-light: High-Throughput Named-Entity Disambiguation. LDOW 2014
- Sairam Gurajada, Stephan Seufert, Iris Miliaraki, Martin Theobald: TriAD: A Distributed Shared-Nothing RDF Engine based on Asynchronous Message Passing. SIGMOD 2014
- Maximilian Dylla, Martin Theobald, Iris Miliaraki: Querying and Learning in Probabilistic Databases. Reasoning Web 2014
- Maximilian Dylla, Iris Miliaraki, Martin Theobald: A Temporal-Probabilistic Database Model for Information Extraction. PVLDB 6(14), 2014
- Maximilian Dylla, Iris Miliaraki, Martin Theobald: Top-k Query Processing in Probabilistic Databases with Non-Materialized Views. ICDE 2013
- Ndapandula Nakashole, Mauro Sozio, Fabian Suchanek, Martin Theobald: Query-Time Reasoning in Uncertain RDF Knowledge Bases with Soft and Hard Rules. VLDS 2012
- Mohamed Yahya, Martin Theobald: D2R2: Disk-Oriented Deductive Reasoning in a RISC-Style RDF Engine. RuleML 2011
- Timm Meiser, Maximilian Dylla, Martin Theobald: Interactive Reasoning in Uncertain RDF Knowledge Bases. CIKM 2011
- Ndapandula Nakashole, Martin Theobald, Gerhard Weikum: Scalable Knowledge Harvesting with High Precision and High Recall. WSDM 2011
- Maximilian Dylla, Mauro Sozio, Martin Theobald: Resolving Temporal Conflicts in Inconsistent RDF Knowledge Bases. BTW 2011
- Ndapandula Nakashole, Martin Theobald, Gerhard Weikum: Find your Advisor: Robust Knowledge Gathering from the Web. WebDB 2010
- Anish Das Sarma, Martin Theobald, Jennifer Widom: LIVE: A Lineage-Supported Versioned DBMS. SSDBM 2010
- Anish Das Sarma, Martin Theobald, Jennifer Widom: Exploiting Lineage for Confidence Computation in Uncertain and Probabilistic Databases. ICDE 2008
- Omar Benjelloun, Anish Das Sarma, Alon Y. Halevy, Martin Theobald, Jennifer Widom: Databases with uncertainty and lineage.
 VLDB J. 17(2), 2008