



LarKC

*The Large Knowledge Collider:
a platform for large scale integrated reasoning and Web-search*

FP7 – 215535

D4.6.1 – Strategies & Design for stream reasoning

Coordinator: Emanuele Della Valle (CEFRIEL)

**With contributions from: Davide Francesco Barbieri,
Daniele Braga, Stefano Ceri, and Emanuele Della Valle
(CEFRIEL)**

Quality Assessor: Frank van Harmelen (VUA)

Quality Controller: Zhisheng Huang (VUA)

Document Identifier:	LarKC/2008/D4.6.1/V1.0
Class Deliverable:	LarKC EU-IST-2008-215535
Version:	Version 1.0
Date:	March 31th, 2010
State:	Final
Distribution:	Public



EXECUTIVE SUMMARY

In this first 24 months of LarKC we have been extensively experimenting with processing of RDF streams with C-SPARQL - the extension to SPARQL we proposed in D3.1.

In D2.6.1, we summarise, in the form of a paper - entitled “Continuous Queries and Real-time Analysis of Social Semantic Data with C-SPARQL” accepted at SDOW 2009¹ - the experiments that we conducted starting from the assumption that our C-SPARQL window-based selection over RDF Streams would outperform the standard filter-based selection that one can implement by storing the RDF stream and using SPARQL. The experiments provided evidence that this is the case, and allowed us to show the advantages of using C-SPARQL in the field of Social Data analysis on the Web.

In D3.3, we focus on the design of an execution environment for C-SPARQL queries. The deliverable consist of a paper – entitled “An Execution Environment for C-SPARQL Queries” accepted at the 13th International Conference on Extending Database Technology (EDBT 2010)²- and of an extended description of the internal design of our C-SPARQL Engine. In the paper, we present (a) the features of an execution environment that leverages existing data stream management system technologies; (b) some optimizations in terms of rewriting rules for efficiently exploiting the designed execution environment; and (c) evidence of the effectiveness of our optimizations on a prototype of execution environment.

In this deliverable, we present, in the form of a paper – entitled “Incremental Reasoning on Streams and Rich Background Knowledge” accepted at the 7th Extended Semantic Web Conference (ESWC 2010)³ - a technique for Stream Reasoning, consisting in incremental maintenance of materializations of ontological entailments in the presence of streaming information. Previous work, delivered in the context of deductive databases, describes the use of logic programming for the incremental maintenance of such entailments. Our contribution is a new technique that exploits the nature of streaming data in order to efficiently maintain materializations of ontological entailments. By adding expiration time information to each RDF triple, we show that it is possible to compute a new complete and correct materialization whenever a new window of streaming data arrives, by dropping explicit statements and entailments that are no longer valid, and then computing when the RDF triples within the window will expire. We provide experimental evidence that our approach significantly reduces the time required to compute a new materialization at each window change, and opens up for several further optimizations.

The materialization, which we compute with our approach, is assumed to be limited to a light ontological language. It enable our C-SPARQL engine to work not only in a Basic RDF entailment regime, but also to support RDFS entailment regime. More expressive reasoning, which can be goal driven, are expected to be cascaded after the C-SPARQL engine.

With the “window based selection” of RDF streams (described in D.2.6.1), the optimized C-SPARQL execution environment, which allows to abstract from fine grains data streams to gross grain events, (described in D3.3), and the technique for incremental maintenance of materializations of ontological entailments in the presence of streaming information (described in this deliverable), we completed the plug-in level research described in D3.1. Our future efforts are focused in implementing and deploying running C-SPARQL plug-ins for the LarKC platform.

¹ <http://sdow.semanticweb.org/2009/>

² <http://ldb.epfl.ch/EDBTICDT/>

³ <http://www.eswc2010.org/>



DOCUMENT INFORMATION

IST Project Number	FP7 – 215535	Acronym	LarKC
Full Title	The Large Knowledge Collider: a platform for large scale integrated reasoning and Web-search		
Project URL	http://www.larkc.eu/		
Document URL			
EU Project Officer	Stefano Bertolo		

Deliverable	Number	4.6.1	Title	Strategies & Design for stream reasoning
Work Package	Number	4	Title	Reasoning and Deciding

Date of Delivery	Contractual	M24	Actual	M24
Status	Draft		final ■	
Nature	prototype <input type="checkbox"/> report ■ dissemination <input type="checkbox"/>			
Dissemination level	public ■ consortium <input type="checkbox"/>			

Authors (Partner)	Davide Francesco Barbieri, Daniele Braga, Stefano Ceri, and Emanuele Della Valle (CEFRIEL)			
Responsible Author	Name	Emanuele Della Valle	E-mail	emanuele.dellavalle@polimi.it
	Partner	Cefriel	Phone	+39 (02) 23954-324

Abstract (for dissemination)	This Deliverable presents a technique for Stream Reasoning, consisting in incremental maintenance of materializations of ontological entailments in the presence of streaming information. Previous work, delivered in the context of deductive databases, describes the use of logic programming for the incremental maintenance of such entailments. Our contribution is a new technique that exploits the nature of streaming data in order to efficiently maintain materialized views of RDF triples, which can be used by a reasoner.
Keywords	C-SPARQL, Stream Reasoning, Incremental Maintenance of Materializations

Version Log			
Issue Date	Rev. No.	Author	Change
16.12.2009	0.0	Della Valle	Paper Initialized
17.12.2009	0.1	Barbieri	Some content developed
18.12.2009	0.2	Della Valle	Large contribution to content
19.12.2009	0.3	Barbieri	Fixed Spelling
20.12.2009	0.4	Ceri	Abstract and section 2 refinement
21.12.2009	0.5	Braga	Minor language fixing
22.12.2009	0.6	Ceri	Section 4 refinement
22.12.2009	0.7	Della Valle	New complete version
22.12.2009	0.8	Ceri	References clean up
1.3.2009	0.9	Della Valle	Preparation of camera ready version of the paper
25.3.2009	1.0	Della Valle	Executive summary written



PROJECT CONSORTIUM INFORMATION

Participant's name	Partner	Contact
Semantic Technology Institute Innsbruck, University of Innsbruck	 	Prof. Dr. Dieter Fensel, Semantic Technology Institute (STI), universitaet Innsbruck, Innsbruck, Austria, E-mail: dieter.fensel@sti-innsbruck.at
AstraZeneca AB		Bosse Andersson AstraZeneca Lund, Sweden Email: bo.h.andersson@astrazeneca.com
CEFRIEL - SOCIETA CONSORTILE A RESPONSABILITA LIMITATA		Emanuele Della Valle, CEFRIEL - SOCIETA CONSORTILE A RESPONSABILITA LIMITATA, Milano, Italy, Email: emanuele.dellavalle@cefriel.it
CYCROP, RAZISKOVANJE IN EKSPERIMENTALNI RAZVOJ D.O.O.		Michael Witbrock, CYCROP, RAZISKOVANJE IN EKSPERIMENTALNI RAZVOJ D.O.O., Ljubljana, Slovenia, Email: witbrock@cyc.com
Höchstleistungsrechenzentrum, Universitaet Stuttgart		Georgina Gallizo, Höchstleistungsrechenzentrum, Universitaet Stuttgart, Stuttgart, Germany, Email: gallizo@hlrs.de
MAX-PLANCK GESELLSCHAFT ZUR FOERDERUNG DER WISSENSCHAFTEN E.V.		Dr. Lael Schooler Max-Planck-Institut für Bildungsforschung Berlin, Germany Email: schooler@mpib-berlin.mpg.de
Ontotext Lab, Sirma Group Corp		Atanas Kiryakov, Ontotext Lab, Sofia, Bulgaria Email: atanas.kiryakov@sirma.bg
SALTLUX INC.		Tony Lee, SALTLUX INC, Seoul, Korea, Email: tony@saltlux.com
SIEMENS AKTIENGESELLSCHAFT		Dr. Volker Tresp, SIEMENS AKTIENGESELLSCHAFT, Muenchen, Germany, E-mail: volker.tresp@siemens.com
THE UNIVERSITY OF SHEFFIELD		Prof. Dr. Hamish Cunningham, THE UNIVERSITY OF SHEFFIELD Sheffield, UK, Email: h.cunningham@dcs.shef.ac.uk



<p>VRIJE UNIVERSITEIT AMSTERDAM</p>		<p>Prof. Dr. Frank van Harmelen, VRIJE UNIVERSITEIT AMSTERDAM, Amsterdam, Netherlands, Email: Frank.van.Harmelen@cs.vu.nl</p>
<p>THE INTERNATIONAL WIC INSTITUTE, BEIJING UNIVERSITY OF TECHNOLOGY</p>		<p>Prof. Dr. Ning Zhong, THE INTERNATIONAL WIC INSTITUTE, Mabeshi, Japan, Email: zhong@maebashi-it.ac.jp</p>
<p>INTERNATIONAL AGENCY FOR RESEARCH ON CANCER</p>		<p>Dr. Paul Brennan, INTERNATIONAL AGENCY FOR RESEARCH ON CANCER, Lyon, France, Email: brennan@iarc.fr</p>
<p>INFORMATION RETRIEVAL FACILITY</p>		<p>John Tait, INFORMATION RETRIEVAL FACILITY Vienna, Austria Email : john.tait@ir-facility.org</p>



TABLE OF CONTENTS

1. INTRODUCTION	7
2. BACKGROUND.....	9
2.1. STREAM REASONING.....	9
2.2. EXPRESSING ONTOLOGY LANGUAGES AS RULES	10
2.3. INCREMENTAL MAINTENANCE OF MATERIALIZATIONS	11
3. MAINTAINING MATERIALIZATION OF RDF STREAMS	13
4. IMPLEMENTATION EXPERIENCE.....	17
5. EVALUATION.....	18
6. CONCLUSION AND FUTURE WORK.....	20
7. REFERENCES	20

Incremental Reasoning on Streams and Rich Background Knowledge

Davide Francesco Barbieri, Daniele Braga, Stefano Ceri, Emanuele Della Valle,
and Michael Grossniklaus

Politecnico di Milano – Dipartimento di Elettronica e Informazione
Piazza L. da Vinci, 32 - 20133 Milano – Italy
{dbarbieri, braga, ceri, dellavalle, grossniklaus}@elet.polimi.it

Abstract. This article presents a technique for Stream Reasoning, consisting in incremental maintenance of materializations of ontological entailments in the presence of streaming information. Previous work, delivered in the context of deductive databases, describes the use of logic programming for the incremental maintenance of such entailments. Our contribution is a new technique that exploits the nature of streaming data in order to efficiently maintain a materialization, which can be used for further reasoning in a goal-directed reasoner.

By adding expiration time information to each RDF triple, we show that it is possible to compute a new complete and correct materialization whenever a new window of streaming data arrives, by dropping explicit statements and entailments that are no longer valid, and then computing when the RDF triples within the window will expire. We provide experimental evidence that our approach significantly reduces the time required to compute a new materialization at each window change, and opens up for several further optimizations.

1 Introduction

Streaming data is an important class of information sources. Examples of data streams are Web logs, feeds, click streams, sensor data, stock quotations, locations of mobile users, and so on. Streaming data is received continuously and in real-time, either implicitly ordered by arrival time, or explicitly associated with timestamps. A new class of database systems, called data stream management systems (DSMS), is capable of performing queries over streams [1], but such systems cannot perform complex reasoning tasks. Reasoners, on the other hand, can perform complex reasoning tasks, but they do not provide support to manage *rapidly* changing worlds.

Recently, we have made the first steps into a new research direction: Stream Reasoning [2] is a new multi-disciplinary approach that can provide the abstractions, foundations, methods, and tools required to integrate data streams, the Semantic Web, and reasoning systems. Central to the notion of stream reasoning is a paradigmatic change from persistent knowledge bases and user-invoked reasoning tasks to transient streams and continuous reasoning tasks.

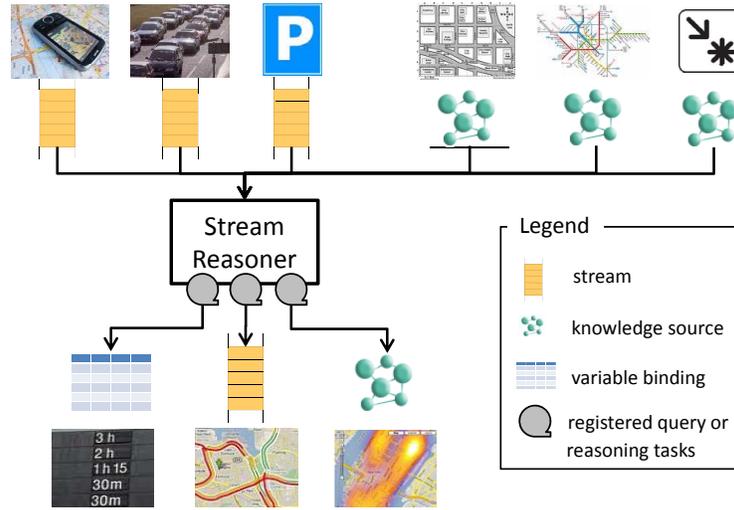


Fig. 1. Mobile Scenario

The first step for enabling Stream Reasoning is the development of languages and systems for querying RDF data also in the form of data streams. Streaming SPARQL [3], Continuous SPARQL (C-SPARQL) [4, 5], and Time-Annotated SPARQL [6] are three recent independent proposals for extending SPARQL to handle both static RDF graphs and transient streams of RDF triples. This paper builds on our previous works on C-SPARQL.

In Fig. 1, we show a Stream Reasoner. It takes several streams of rapidly changing information and several static sources of background knowledge as input. In the context of a mobile scenario, examples of sources of streaming data can be the positions of users, the traffic in the streets, and the availability of parking lots, whereas examples of background knowledge can be the city layout, the public transportation schedules, and the descriptions of points of interest and of events in a given area. Several reasoning tasks, expressed in the form of C-SPARQL queries, are registered into the stream reasoner, and the system continuously generates new answers. These answers can be in the standard SPARQL output form (i.e., variable bindings and graphs) or in the form of streams. In our mobile scenario, for instance, we can register two C-SPARQL queries: one continuously monitors the status of the public transportation system and returns the delays as variable bindings, the other one monitors the sensors for traffic detection and generates a stream of aggregate information for each major road. Current implementations of the proposed SPARQL extensions, however, assume only a simple entailment (see Section 2 of [7]). They do not try to handle reasoning on streaming information, e.g., providing strategical suggestions about how to perform goals.

In existing work on logical reasoning, the knowledge base is always assumed to be static (or slowly evolving). There is work on changing beliefs on the basis of

new observations [8], but the solutions proposed in this area are far too complex to be applicable to gigantic data streams of the kind we image in a mobile context. However, the nature of data streams is different from arbitrary changes, because change occurs in a “regular” way at the points where the streaming data is observed.

In this article, we present a technique for stream reasoning that incrementally maintains a materialization of ontological entailments in the presence of streaming information. We elaborate on previous papers [9, 10] that extend to logic programming results from incremental maintenance of materialized views in deductive databases [11]. Our contribution is a new technique that takes the order in which streaming information arrives at the Stream Reasoner into explicit consideration. By adding expiration time information to each RDF statement, we show that it is possible to compute a new complete and correct materialization by (a) dropping explicit statements and entailments that are no longer valid, and (b) evaluating a maintenance program that propagates insertions of explicit RDF statements as changes to the stored implicit entailments.

The rest of the paper is organized as follows. Section 2 presents a wrap up of the background information needed to understand this paper. In particular, it presents the state of the art in incremental maintenance of materializations of ontologies represented as logic programs. Section 3 presents our major contribution in the form of Datalog rules computing the incremental materialization of ontologies for window-based changes of ontological entailments. In Section 4 we present our implementation experience. Section 5 provides experimental evidence that our approach significantly reduces the time required to compute the new materialization. Finally, we close the paper by sketching future works in Section 6.

2 Background

2.1 Stream Reasoning

A first step toward stream reasoning has been to combine the power of existing data-stream management systems and the Semantic Web [12]. The key idea is to keep streaming data in relational format as long as possible and to bring it to the semantic level as aggregated events [5]. Existing data models, access protocols, and query languages for data-stream management systems and the Semantic Web are not sufficient to do so and, thus, they must be combined.

C-SPARQL [4, 5] introduces the notion of RDF streams as the natural extension of the RDF data model to this scenario, and then extend SPARQL to query RDF streams. An RDF stream is defined as an ordered sequence of pairs, where each pair is constituted by an RDF triple and its timestamp τ .

$$\begin{array}{c} \dots \\ (\langle subj_i, pred_i, obj_i \rangle, \tau_i) \\ (\langle subj_{i+1}, pred_{i+1}, obj_{i+1} \rangle, \tau_{i+1}) \\ \dots \end{array}$$

Fig. 2 shows an example of a C-SPARQL query that continuously queries a RDF stream as well as a static RDF graph. The RDF stream describes the users sitting in trains and trains moving from a station to another one. The RDF graph describes where the stations are located, e.g., a station is in a city, which is in a region.

```

1. REGISTER QUERY WhereCommutersAre COMPUTE EVERY 1sec AS
2. PREFIX ex: <http://example/>
3. SELECT DISTINCT ?user ?type ?x
4. FROM <http://mobileservice.org/meansOfTransportation.rdf>
5. FROM STREAM <http://mobileservice.org/positions.trdf>
6. [RANGE 10sec STEP 1sec]
7. WHERE {
8.   ?user ex:isIn ?x .
9.   ?user a ex:Commuter .
10.  ?x a ?type .
11.  ?user ex:remainingTravelTime ?t .
12.  FILTER (?t >= "PT30M"xsd:duration )
13. }
```

Fig. 2. An example of C-SPARQL query that continuously queries a RDF stream as well as a static RDF graph

At line 1, the `REGISTER` clause instructs the C-SPARQL engine to register a continuous query. The `COMPUTE EVERY` clause states the frequency of every new computation. In line 5, the `FROM STREAM` clause defines the RDF stream of positions used in the query. Next, line 6 defines the window of observation of the RDF stream. Streams, by their very nature, are volatile and consumed on the fly. The C-SPARQL engine, therefore, observes them through a window that contains the stream's most recent elements and that changes over time. In the example, the window comprises RDF triples produced in the last 10 seconds and the window slides every second. The `WHERE` clause is standard SPARQL as it includes a set of matching patterns, which restricts users to be commuters and a `FILTER` clause, which restricts the answers to users whose remaining traveling time is at least 30 minutes. This example shows that, at the time of the presentation in the window, it is possible to compute the time when triples both of the window and of ontological entailments will cease to be valid.

2.2 Expressing Ontology Languages as Rules

Using rules is a best practice (see Section 2.1 of [9]) in implementing the logical entailment supported by ontology languages such as RDF-S [13] and OWL2-RL [14]. For example, Fig. 3 presents the set of rule used by the Jena Generic Rule Engine [15] to compute RDF-S closure. The first rule (`rdfs2`) states that if there is a triple `<?x ?p ?y>` and the domain of the property `?p` is the class

?c (represented by the triple $\langle ?p \text{ rdfs:domain } ?c \rangle$) then the resource ?x is of type ?c (represented by the triple $\langle ?x \text{ rdf:type } ?c \rangle$).

```
[rdfs2: (?x ?p ?y), (?p rdfs:domain ?c) -> (?x rdf:type ?c)]
[rdfs3: (?x ?p ?y), (?p rdfs:range ?c) -> (?y rdf:type ?c)]
[rdfs5a: (?a rdfs:subPropertyOf ?b), (?b rdfs:subPropertyOf ?c)
-> (?a rdfs:subPropertyOf ?c)]
[rdfs5b: (?a rdf:type rdf:Property) -> (?a rdfs:subPropertyOf ?a)]
[rdfs6: (?a ?p ?b), (?p rdfs:subPropertyOf ?q) -> (?a ?q ?b)]
[rdfs7: (?a rdf:type rdfs:Class) -> (?a rdfs:subClassOf ?a)]
[rdfs8: (?a rdfs:subClassOf ?b), (?b rdfs:subClassOf ?c)
-> (?a rdfs:subClassOf ?c)]
[rdfs9: (?x rdfs:subClassOf ?y), (?a rdf:type ?x) -> (?a rdf:type ?y)]
[rdfs10: (?x rdf:type rdfs:ContainerMembershipProperty)
-> (?x rdfs:subPropertyOf rdfs:member)]
[rdf1and4: (?x ?p ?y) -> (?p rdf:type rdf:Property),
(?x rdf:type rdfs:Resource),
(?y rdf:type rdfs:Resource)]
[rdfs7b: (?a rdf:type rdfs:Class) -> (?a rdfs:subClassOf rdfs:Resource)]
```

Fig. 3. Rules Implementing RDF-S in Jena Generic Rule Engine

In the rest of the paper, we adopt logic programming terminology. We refer to a set of rules as a *logic program* (or simply program) and we assume that any RDF graph can be stored in the extension of a single ternary predicate P . Under this assumption, the rule rdfs2 can be represented in Datalog as follows.

$$P(x, rdf : type, c) :- P(p, rdfs : domain, C), P(s, p, y)$$

2.3 Incremental Maintenance of Materializations

Maintenance of a materialization when facts change, i.e., facts are added or removed from the knowledge base, is a well studied problem. The state of the art approach implemented in systems such as KAON¹ is a declarative variant [9] of the delete and re-derive (DRed) algorithm proposed in [16]. DRed incrementally maintains a materialization in three steps.

1. Overestimate the deletions by computing all the direct consequences of a deletion.
2. Prune the overestimated deletions for which the deleted fact can be re-derived from other facts.
3. Insert all derivation which are consequences of added facts.

¹ The Datalog engine is part of the KAON suite, see <http://kaon.semanticweb.org>

More formally, a logic program is composed by a set of rules \mathbf{R} that we can represent as $H :- B_1, \dots, B_n$, where H is the predicate that forms the head of the rule and B_1, \dots, B_n are the predicates that form the body of the rule. If we call the set of predicates in a logic program \mathbf{P} , then we can formally assert that $H, B_i \in \mathbf{P}$. A *maintenance program*, which implements the declarative version of the DRed algorithm, can be automatically derived from the original program with a fixed set of rewriting functions (see Table 2) that uses seven maintenance predicates (see Table 1) [9].

Name	Content of the extension
P	the current materialization
P^{Del}	the deletions
P^{Ins}	the explicit insertion
P^{Red}	the triples marked for deletion which have alternative derivations
P^{New}	the materialization after the execution of the maintenance program
P^+	the net insertions required to maintain the materialization
P^-	the net deletions required to maintain the materialization

Table 1. The maintenance predicates (derived from [9])

Given a materialized predicate P and the set of extensional insertions P^{Ins} to and deletions P^{Dels} from P , the goal of the rewriting functions is the definition of two maintenance predicates P^+ and P^- , such that the extensions of P^+ and P^- contain the net insertions and deletions, respectively, that are needed to incrementally maintain the materialization of P .

Predicate		
Name	Generator Parameter	Rewriting Result
δ_1^{New}	$P \in \mathbf{P}$	$P^{New} :- P, not P^{Del}$
δ_2^{New}	$P \in \mathbf{P}$	$P^{New} :- P^{Red}$
δ_3^{New}	$P \in \mathbf{P}$	$P^{New} :- P^{Ins}$
δ^+	$P \in \mathbf{P}$	$P^+ :- P^{Ins}, not P$
δ^-	$P \in \mathbf{P}$	$P^- :- P^{Del}, not P^{Ins}, not P^{Red}$
Rule		
Name	Generator Parameter	Rewriting Result
δ^{Red}	$H :- B_1, \dots, B_n$	$H^{Red} :- H^{Del}, B_1^{New}, \dots, B_n^{New}$
δ^{Del}	$H :- B_1, \dots, B_n$	$\{H^{Del} :- B_1, \dots, B_{i-1}, B_i^{Del}, B_{i+1}, \dots, B_n\}$
δ^{Ins}	$H :- B_1, \dots, B_n$	$\{H^{Ins} :- B_1^{New}, \dots, B_{i-1}^{New}, B_i^{Ins}, B_{i+1}^{New}, \dots, B_n^{New}\}$

Table 2. Rewriting functions (derived from [9])

We can divide the rewriting functions shown in Table 2 in two groups. One group of functions apply to predicates, while the other group of functions apply to rules. The former functions use the predicates defined in Table 1 to introduce the rules that will store the materialization after the execution of the maintenance program in the extension of the predicate P^{New} . The latter functions introduce the rules that populate the extensions of the predicates P^{Del} , P^{Red} , and P^{Ins} . These three rewriting functions are executed for each rule that has the predicate P as head. While the function δ^{Red} rewrites each rule in exactly one maintenance rule, the two functions δ^{Del} and δ^{Ins} rewrite each rule with n bodies B_i into n maintenance rules.

To exemplify how these rewriting functions work in practice, let us return to the scenario exemplified in Sect. 2.1. To describe that scenario, we introduced the predicate $isIn$ that captures the respective position of moving objects (e.g., somebody is in a train, the train is in station, somebody else is in a car, the car is in a parking lot, etc.). A simple ontology for a mobility scenario could express transitivity and be represented using the following Datalog rule.

$$(R) \text{ isIn}(x, z) :- \text{isIn}(x, y), \text{isIn}(y, z)$$

By applying the rewriting functions presented in Table 2 to the rule (R) and the predicate $isIn$, we obtain the maintenance program shown in Table 3. Each row of the table contains the applied rewriting function and the rewritten maintenance rule.

Rule	Rewriting Function
$isIn^{New}(x, y) :- isIn(x, y), not\ isIn^{Del}(x, y)$	$\delta_1^{New}(isIn)$
$isIn^{New}(x, y) :- isIn^{Red}(x, y)$	$\delta_2^{New}(isIn)$
$isIn^{New}(x, y) :- isIn^{Ins}(x, y)$	$\delta_3^{New}(isIn)$
$isIn^+(x, y) :- isIn^{Ins}(x, y), not\ isIn(x, y)$	$\delta^+(isIn)$
$isIn^-(x, y) :- isIn^{Del}(x, y), not\ isIn^{Ins}(x, y), not\ isIn^{Red}(x, y)$	$\delta^-(isIn)$
$isIn^{Red}(x, z) :- isIn^{Del}(x, z), isIn^{New}(x, y), isIn^{New}(y, z)$	$\delta^{Red}(R)$
$isIn^{Del}(x, z) :- isIn^{Del}(x, y), isIn(y, z)$	$\delta^{Del}(R)$
$isIn^{Del}(x, z) :- isIn(x, y), isIn^{Del}(y, z)$	$\delta^{Del}(R)$
$isIn^{Ins}(x, z) :- isIn^{Ins}(x, y), isIn^{new}(y, z)$	$\delta^{Ins}(R)$
$isIn^{Ins}(x, z) :- isIn^{new}(x, y), isIn^{Ins}(y, z)$	$\delta^{Ins}(R)$

Table 3. The maintenance program automatically derived from a program containing only the rule R by applying the rewriting functions show in Table 2

3 Maintaining Materialization of RDF Streams

As we explained earlier in this paper, incremental maintenance of materializations of ontological entailments after knowledge changes is a well studied problem. However, additions or removals of facts from the knowledge base induced

TS	Triples in the Window	Entailments in the Window
1	A $\xrightarrow{[11]}$ B	
2	A $\xrightarrow{[11]}$ B $\xrightarrow{[12]}$ C	A $\xrightarrow{[11]}$ C
3	A $\xrightarrow{[11]}$ B $\xrightarrow{[12]}$ C $\xrightarrow{[13]}$ D	A $\xrightarrow{[11]}$ B $\xrightarrow{[12]}$ C $\xrightarrow{[13]}$ D A $\xrightarrow{[11]}$ C $\xrightarrow{[12]}$ D
4	A $\xrightarrow{[11]}$ B $\xrightarrow{[12]}$ C $\xrightarrow{[13]}$ D A $\xrightarrow{[14]}$ E $\xrightarrow{[14]}$ D	A $\xrightarrow{[11]}$ B $\xrightarrow{[12]}$ C $\xrightarrow{[13]}$ D A $\xrightarrow{[14]}$ E $\xrightarrow{[14]}$ D A $\xrightarrow{[11]}$ C $\xrightarrow{[12]}$ D
...
12	A $\xrightarrow{[11]}$ B $\xrightarrow{[12]}$ C $\xrightarrow{[13]}$ D A $\xrightarrow{[14]}$ E $\xrightarrow{[14]}$ D	A $\xrightarrow{[14]}$ E $\xrightarrow{[14]}$ D B $\xrightarrow{[12]}$ C $\xrightarrow{[13]}$ D
13	A $\xrightarrow{[14]}$ E $\xrightarrow{[14]}$ D C $\xrightarrow{[13]}$ D	A $\xrightarrow{[14]}$ D

Fig. 4. Our approach to incrementally maintain the materialization at work.

by data streams are governed by windows, which have a known expiration time. The intuition behind our approach is straightforward. If we tag each RDF triple (both explicitly inserted and entailed) with a *expiration time* that represents the last moment in which it will be in the window, we can compute a new complete and correct materialization by dropping RDF triples that are no longer in the window and then evaluate a maintenance program that

1. computes the entailments derived by the inserts,
2. annotates each entailed triple with a expiration time, and
3. eliminates from the current state all copies of derived triples except the one with the highest timestamp.

Note that this approach supports the immediate deletions of both window facts and entailed triples which are dropped by inspection to their expiration times. Instead it requires some extra work for managing insertions as new timestamps need to be computed. This approach is more effective than overestimating the deletions and then computing re-derivations, as we will demonstrate in Section 5.

Figure 4 illustrates our approach. Let us assume that we have a stream of triples in which all the triples use the same predicate *isIn* introduced in Section 2.3. Let us also assume that we register a simple C-SPARQL query that observes an RDF stream through a sliding window of 10 seconds and computes the transitive closure of the *isIn* property.

In the 1st second of execution, the triple $\langle A \text{ isIn } B \rangle$ enters the window. We tag the triple with the expiration time 11 (i.e., it will be valid until the 11th second) and no derivation occurs. The transitive closure only contains that triple. In the 2nd second the triple $\langle B \text{ isIn } C \rangle$ enters the window. We can tag it with the expiration time 12 and we can materialize the entailed triple $\langle A \text{ isIn } C \rangle$. As the triple $\langle A \text{ isIn } B \rangle$ expires in the 11th second, the entailed triple $\langle A \text{ isIn } C \rangle$ also expires then and, thus, we tag it with the expiration time 11 (i.e., Step 2 of our approach). As the 11th second passes, we will have to just drop the triples tagged with 11 and the materialization will be up to date (i.e., Step 1 of our approach).

Let us then assume that in the 3rd second, the triple $\langle C \text{ isIn } D \rangle$ enters the window. We tag it with the expiration time 13 and compute two entailments: the triple $\langle B \text{ isIn } D \rangle$ with expiration time 12 and the triple $\langle A \text{ isIn } D \rangle$ with expiration time 11. In the 4th second, the two triples $\langle A \text{ isIn } E \rangle$ and $\langle E \text{ isIn } D \rangle$ enter the window. Both triples are tagged with the expiration time 14. We also derive the entailed triple $\langle A \text{ isIn } D \rangle$ with time expiration 14. The triple $\langle A \text{ isIn } D \rangle$ was previously derived, but its expiration time was 11 and, therefore, that triple is dropped. The rest of Fig. 4 shows how triples are deleted when they expire.

More formally, our logic program is composed of a set of rules \mathbf{R} that we can represent as $H[T] :- B_1[T_1], \dots, B_n[T_n]$, where H is the predicate that forms the head of the rule and it is valid until T . $B_1[T_1], \dots, B_n[T_n]$ are the n predicates that form the body of the rule with their respective n expiration times $T_1 \dots T_n$. As in the case illustrated in Section 2.3, we can formally assert that $H, B_i \in \mathbf{P}$ where \mathbf{P} denotes the set of predicates in a logic program.

A maintenance program, which implements our approach in a declarative way, can be automatically be derived from the original program with a fixed set of rewriting functions (see Table 5) that uses five maintenance predicates (see Table 4) inspired by the approach of Volz et al. [9].

Name	Content of the extension
P	the current materialization
P^{Ins}	the triples that enter the window
P^{New}	the triples which are progressively added to the materialization
P^{Old}	the triples for which re-derivations with a longer expiration time were materialized
P^+	the net insertions required to maintain the materialization
P^-	the net deletions required to maintain the materialization

Table 4. The maintenance predicates of our approach

Given a materialized predicate P and set of extensional insertions P^{Ins} determined by the new triple entering the window, the goal of the rewriting functions is the definition of the maintenance predicates P^+ and P^- whose extension contains the net insertions and the net deletions needed to incrementally main-

tain the materialization of P . The extension of the maintenance predicate P^- contains the extensions of predicate P that expires as well as the extension of predicate P^{Old} . In Table 5 we formally defines our rewriting functions. Note that P^{++} is only an auxiliary predicate with not special meaning.

Predicate		
Name	Generator Parameter	Rewriting Result
Δ_1^{New}	$P \in \mathbf{P}$	$P^{New}[T] :- P[T], not P[T_1], T_1 = (now - 1)$
Δ_2^{New}	$P \in \mathbf{P}$	$P^{New}[T] :- P^{Ins}[T], not P^{Old}[T]$
Δ_1^{Old}	$P \in \mathbf{P}$	$P^{Old}[T] :- P^{Ins}[T_1], P[T], T_1 > T$
Δ_2^{Old}	$P \in \mathbf{P}$	$P^{Old}[T] :- P^{Ins}[T_1], P^{Ins}[T], T_1 > T$
Δ_1^-	$P \in \mathbf{P}$	$P^- [T] :- P[T_1], T_1 = (now - 1), not P^{Ins}[T_1]$
Δ_2^-	$P \in \mathbf{P}$	$P^- [T] :- P^{Old}[T]$
Δ^{++}	$P \in \mathbf{P}$	$P^{++}[T] :- P^{New}[T], not P[T_1]$
Δ^+	$P \in \mathbf{P}$	$P^+[T] :- P^{++}[T], not P^{Old}[T_1]$
Rule		
Name	Generator Parameter	Rewriting Result
Δ^{Ins}	$H :- B_1, \dots, B_n$	$\{H^{Ins}[T] :- B_1^{New}[T_1], \dots, B_{i-1}^{New}[T_{i-1}], B_i^{Ins}[T_i], B_{i+1}^{New}[T_{i+1}], \dots, B_n^{New}[T_n], T = \min(T_1, \dots, T_n)\}$

Table 5. The rewriting functions of our approach

By applying the rewriting functions presented in Table 5 to the rule (R) and the predicate $isIn$ defined in Section 2.3, we obtain the maintenance program shown in Table 6.

Rule	Function
$isIn^{New}(x, y)[T] :- isIn(x, y)[T], not isIn(x, y)[T_1], T_1 = (now - 1)$	$\Delta_1^{New}(isIn)$
$isIn^{New}(x, y)[T] :- isIn^{Ins}(x, y)[T], not isIn^{Old}(x, y)[T]$	$\Delta_2^{New}(isIn)$
$isIn^{Old}(x, y)[T] :- isIn^{Ins}(x, y)[T_1], isIn(x, y)[T], T_1 > T$	$\Delta_1^{Old}(isIn)$
$isIn^{Old}(x, y)[T] :- isIn^{Ins}(x, y)[T_1], isIn^{Ins}(x, y)[T], T_1 > T$	$\Delta_2^{Old}(isIn)$
$isIn^-(x, y)[T] :- isIn(x, y)[T_1], T_1 = (now - 1), not isIn^{Ins}(x, y)[T_1]$	$\Delta_1^-(isIn)$
$isIn^-(x, y)[T] :- isIn^{Old}(x, y)[T]$	$\Delta_2^-(isIn)$
$isIn^{++}(x, y)[T] :- isIn^{New}(x, y)[T], not isIn(x, y)[T_1]$	$\Delta^{++}(isIn)$
$isIn^+(x, y)[T] :- isIn^{++}(x, y)[T], not isIn^{Old}(x, y)[T_1]$	$\Delta^+(isIn)$
$isIn^{Ins}(x, z)[T] :- isIn^{Ins}(x, y)[T_1], isIn^{New}(y, z)[T_2], T = \min(T_1, T_2)$	$\Delta^{Ins}(R)$
$isIn^{Ins}(x, z)[T] :- isIn^{New}(x, y)[T_1], isIn^{Ins}(y, z)[T_2], T = \min(T_1, T_2)$	$\Delta^{Ins}(R)$

Table 6. The maintenance program automatically derived from a program containing only the rule R by applying the rewriting functions show in Table 5

4 Implementation Experience

Figure 5 illustrates the architecture of our current prototype, implemented by using the Jena Generic Rule Engine. The *Incremental Maintainer* component orchestrates the maintenance process. It keeps the current materialization in the *Permanent Space* and uses the *Working Space* to compute the net inserts and deletes. Both spaces consist of an RDF store for the triples and a hashtable which caters for efficient management of the expiration time associated with each triple.

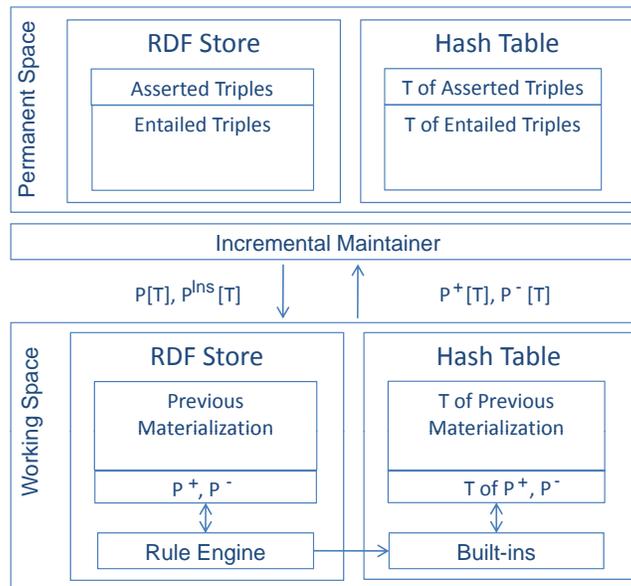


Fig. 5. Overview of the prototype implementation

The maintenance program (see Fig. 6) is loaded into the rule engine that operates over the RDF store in the working space. The management of expiration times is performed by using four custom built-ins, *GetVT*, *GetDiffVT*, *SetVT* and *DelVT*, that are triggered by the maintenance program². *GetVT* retrieves the expiration time of a triple from the hashtable; *GetDiffVT* gets possible other expiration times of a given triple and is used to efficiently implement the rules generated by Δ_2^{Old} ; *SetVT* sets the expiration time of a triple in the hashtable; *DelVT* deletes the expiration time of a triple from the hashtable.

The maintenance process is carried out as follows. When the system is started up, the background knowledge is loaded into the permanent space. Then, the

² For more information on how to write built-ins for Jena Generic Rule Engine see [15]

```

[New1: (?A isIn ?B), GetVT(?A isIn ?B, ?T), noValue(?A isInExp ?B)
      -> (?A isInNew ?B), SetVT(?A isInNew ?B, ?T)]
[New2: (?A isInIns ?B), GetVT(?A isInIns ?B, ?T), noValue(?A isInOld ?B)
      -> (?A isInNew ?B), SetVT(?A isInNew ?B, ?T)]
[Old1: (?A isInIns ?B), GetVT(?A isInIns ?B, ?T1),
      (?A isIn ?B), GetVT(?A isIn ?B, ?T), lessThan(?T, ?T1)
      -> (?A isInOld ?B), DelVT(?A isInIns ?B, ?T)]
[Old2: (?A isInIns ?B), GetVT(?A isInIns ?B, ?T1),
      (?A isInIns ?B), GetDiffVT(?A isIn ?B, ?T1, ?T), lessThan(?T, ?T1)
      -> (?A isInOld ?B), DelVT(?A isInIns ?B, ?T)]
[Rem1: (?A isInExp ?B), GetVT(?A isInExp ?B, ?T), noValue(?A isInIns ?B)
      -> (?A isInRem ?B), DelVT(?A isInExp ?B, ?T) ]
[Rem2: (?A isInOld ?B) -> (?A isInRem ?B) ]
[Add2: (?A isInNew ?B), GetVT(?A isInNew ?B, ?T), noValue(?A isIn ?B)
      -> (?A isInAdd2 ?B), SetVT(?A isInAdd2 ?B, ?T) ]
[Add1: (?A isInAdd2 ?B), GetVT(?A isInAdd2 ?B, ?T), noValue(?A isInOld ?B)
      -> (?A isInAdd ?B), SetVT(?A isInAdd ?B, ?T) ]
[Ins1: (?A isInIns ?B), GetVT(?A isInIns ?B, ?T1),
      (?B isInNew ?C), GetVT(?B isInNew ?C, ?T2), min(?T1, ?T2, ?T)
      -> (?A isInIns ?C), SetVT(?A isInIns ?C, ?T)]
[Ins2: (?A isInNew ?B), GetVT(?A isInNew ?B, ?T1),
      (?B isInIns ?C), GetVT(?B isInIns ?C, ?T2), min(?T1, ?T2, ?T)
      -> (?A isInIns ?C), SetVT(?A isInIns ?C, ?T)]

```

Fig. 6. The maintenance program shown in Table 6 implemented in Jena Generic Rule Engine

maintenance program is evaluated on the background knowledge and the extension of all predicates P is stored in the RDF store. The expiration time of all triples is set to a default value which indicates that they cannot expire. As the window slides over the stream(s), the incremental maintainer:

- (a) puts all triples entering the window in the extension of P^{Ins} ,
- (b) loads the current materialization and P^{Ins} in the working space,
- (c) copies the expiration times from the permanent space into the working space,
- (d) evaluates the maintenance program,
- (e) updates the RDF store in the permanent space by adding the extension of P^+ and removing the extension of P^- ,
- (f) updates the hash tables by changing the expiration time of the triples in the extension of P^+ and removing from the table the triples of P^- , and
- (g) clears the working space for a new evaluation.

5 Evaluation

This section reports on the evaluation we carried out using various synthetically generated data sets that use the transitive property *isIn*. Although we limit our experiments to the transitive property, the test is significant because widely used

vocabularies in Web ontological languages are transitive (e.g., `rdfs:subClassOf`, `rdfs:subPropertyOf`, `owl:sameAs`, `owl:equivalentProperty`, `owl:equivalentClass` and all properties of type `owl:TransitiveProperty`). Moreover, transitive properties are quite generative in terms of entailments and, thus, stress the system.

Our synthetic data generator generates trees of triples all using *isIn* as property. We can control the depth of the tree and the number of trees generated. All generated triples are stored in a pool. An experiment consists of measuring the time needed to compute a new materialization based on the given the background knowledge, the triples in the window as well as the triples that enter and exit the window at each step. When we start an experiment, we first extract a subset of triples from the pool to form the background knowledge. Then, we stream the rest of the triples from the pool. We control both the dimension of the window over the stream of triples and the number of triples entering and exiting the window at each step.

In our experiments we compare three approaches: (a) the *naive* approach of recomputing the entire materialization at each step, (b) the maintenance program shown in Table 3 implementing [9], denoted as *incremental-volz*, and (c) our maintenance program shown in Tables 6 and in Fig. 6, denoted as *incremental-stream*. Intuitively, the naive approach is dominated with a small number of streaming triples, and dominates when streaming triples are a large fraction of the materialization.

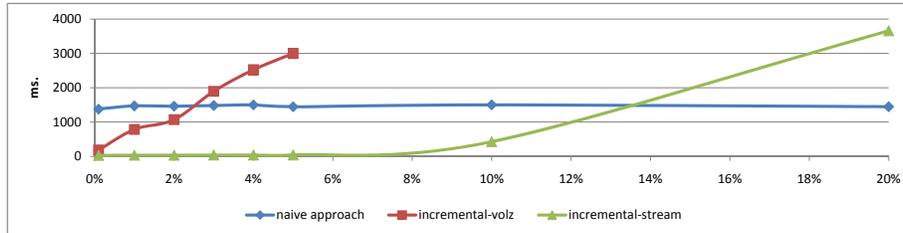


Fig. 7. Evaluation results: the time (ms) required to maintain the materialization as a function of the percentage of the background knowledge subject to change

We run multiple experiments³ using different settings of the experimental environment, by changing the size of the background knowledge, the size of the window, and the number of triples entering and exiting the window at each step. In Fig. 7, we plot the results of one of these experiments (which qualitatively are very similar). We compare the materialization maintenance time as a function of the percentage of the background knowledge subject to change. As one can read from the graph, the incremental-volz [9] approach is faster than the naive approach only if the changes induced by the streaming triples encompass less than 2.5% of the background knowledge. Our incremental-stream approach is

³ We run all experiment on an Intel® Core™ Duo 2.20 GHz with 2 GB of RAM

an order of magnitude faster than incremental-volz for up to 0.1% of changes and continues to be two orders of magnitude faster up to 2.5% of changes. It no longer pays off with respect to the naive approach when the percentage of change is above 13%.

6 Conclusion and Future Work

In this paper, we have shown how previous work from the field of deductive databases can be applied to the maintenance of ontological entailments with data streams. Our approach is an extension of the algorithm developed by Volz et al. [9], that uses logic programming to maintain materializations incrementally. Data streams use the notion of windows to extract snapshots from streams, that are then processed by the query evaluator; we leverage this fact to define the triples that are inserted into and deleted from the materialization. We have also presented an implementation as an extension of the Jena Generic Rule Engine; our implementation uses hash tables to manage triple expiration time. We have shown that our approach outperforms existing approaches when the window size is a fraction (below 10%) of the knowledge base: this assumption holds for all known data stream applications.

We foresee several extensions to this work. With our approach, at insertion time we explicitly remove old triples which have multiple derivations, but we are considering the option of keeping all derivations and simply let them expire when they expire, thus simplifying also insertions. Of course, this requires programs (e.g., our C-SPARQL engine) to be aware of the existence of multiple instances of the same triple, with different expiration times, and ignore all but one of such instances. Another open problem is the application of our approach to several queries over the same streams, with several windows that move at different intervals. A possible solution to this problem is to build the notion of “maximal common sub-window” and then apply the proposed algorithm to them. This is an original instance of multi-query optimization, that is indeed possible when queries are preregistered (as with stream databases and C-SPARQL). Finally, we intend to explore a “lazy” approach to materialization, in which only entailments that are needed to answer registered queries are computed. In our future work, we plan to address these issues.

Acknowledgements

The work described in this paper has been partially supported by the European project LarKC (FP7-215535). Michael Grossniklaus’s work is carried out under SNF grant number PBEZ2-121230.

References

1. Garofalakis, M., Gehrke, J., Rastogi, R.: Data Stream Management: Processing High-Speed Data Streams (Data-Centric Systems and Applications). Springer-Verlag New York, Inc. (2007)

2. Della Valle, E., Ceri, S., van Harmelen, F., Fensel, D.: It's a Streaming World! Reasoning upon Rapidly Changing Information. *IEEE Intelligent Systems* **24**(6) (2009) 83–89
3. Bolles, A., Grawunder, M., Jacobi, J.: Streaming SPARQL – Extending SPARQL to Process Data Streams. In: *Proc. Europ. Semantic Web Conf. (ESWC)*. (2008) 448–462
4. Barbieri, D.F., Braga, D., Ceri, S., Della Valle, E., Grossniklaus, M.: C-SPARQL: SPARQL for Continuous Querying. In: *Proc. Intl. Conf. on World Wide Web (WWW)*. (2009) 1061–1062
5. Barbieri, D.F., Braga, D., Ceri, S., Grossniklaus, M.: An Execution Environment for C-SPARQL Queries. In: *Proc. Intl. Conf. on Extending Database Technology (EDBT)*. (2010)
6. Rodriguez, A., McGrath, R., Liu, Y., Myers, J.: Semantic Management of Streaming Data. In: *Proc. Intl. Workshop on Semantic Sensor Networks (SSN)*. (2009)
7. McBride, B., Hayes, P.: RDF Semantics. W3C Recommendation (2004) <http://www.w3.org/TR/rdf-mt/>.
8. Gaerdenfors, P., ed.: *Belief Revision*. Cambridge University Press (2003)
9. Volz, R., Staab, S., Motik, B.: Incrementally maintaining materializations of ontologies stored in logic databases. *J. Data Semantics* **2** (2005) 1–34
10. Staudt, M., Jarke, M.: Incremental maintenance of externally materialized views. In Vijayaraman, T.M., Buchmann, A.P., Mohan, C., Sarda, N.L., eds.: *VLDB, Morgan Kaufmann* (1996) 75–86
11. Ceri, S., Widom, J.: Deriving production rules for incremental view maintenance. In Lohman, G.M., Sernadas, A., Camps, R., eds.: *VLDB, Morgan Kaufmann* (1991) 577–589
12. Della Valle, E., Ceri, S., Barbieri, D.F., Braga, D., Campi, A.: A First Step Towards Stream Reasoning. In: *Proc. Future Internet Symposium (FIS)*. (2008) 72–81
13. Brickley, D., Guha, R.: RDF Vocabulary Description Language 1.0: RDF Schema. W3C Recommendation (2004) <http://www.w3.org/TR/rdf-schema/>.
14. Motik, B., Grau, B.C., Horrocks, I., Wu, Z., Fokoue, A., Lutz, C.: Owl 2 web ontology language: Profiles. W3C Recommendation (2009) <http://www.w3.org/TR/owl2-profiles/>.
15. Reynolds, D.: Jena 2 inference support (2009) <http://jena.sourceforge.net/inference/>.
16. Gupta, A., Mumick, I.S., Subrahmanian, V.S.: Maintaining views incrementally. In Buneman, P., Jajodia, S., eds.: *SIGMOD Conference, ACM Press* (1993) 157–166