

NATIVE XML SUPPORT FOR SEMISTRUCTURED PROBABILISTIC
DATA MANAGEMENT

A Thesis
Presented to
the Faculty of California Polytechnic State University
San Luis Obispo

In Partial Fulfillment
of the Requirements for the Degree
Master of Science in Computer Science

by
Evan Pierce Rosson

AUTHORIZATION FOR REPRODUCTION OF MASTER'S THESIS

I reserve the reproduction rights of this thesis for a period of seven years from the date of submission. I waive reproduction rights after the time span has expired.

Signature

Date

APPROVAL PAGE

TITLE: Native XML Support for Semistructured Probabilistic Data Management

AUTHOR: Evan Pierce Rosson

DATE SUBMITTED: June 2007

Dr. Alexander Dekhtyar
Advisor or Committee Chair _____
Signature

Dr. John Clements
Committee Member _____
Signature

Dr. Clark Turner
Committee Member _____
Signature

Abstract

Native XML Support for Semistructured Probabilistic Data Management

Evan Pierce Rosson

Many applications require storing and querying probabilistic information; for example, the risk analysis software used by insurance companies. Probabilistic databases are designed to store such data and support queries using operations based on probability theory. Semistructured databases, often based on XML, allow storage of data that may not strictly conform to a schema, which permits data imported from sources with many different schema. The Semistructured Probabilistic Database Management System (SPDBMS) combines these two approaches. It stores Semistructured Probabilistic Objects (SPOs), probability distributions of variables with discrete and finite domains expressed as XML documents. The SQL-like Semistructured Probabilistic Object Query Language (SPOQL) is used to query and manipulate SPOs using operations based on probability theory.

We present a native XML implementation of SPDBMS, better suited to the semistructured nature of this data than the original relational backend. This is implemented using XQuery, a functional query language for processing XML; and ExistDB, an open-source XML database. The performance of this new implementation is compared with the existing relational implementation. We also implement a new SPDBMS query operation, MIX, and a distinction between LEFT and RIGHT JOIN queries.

Contents

List of Figures	viii
1 Introduction	1
1.1 Probabilistic Databases	1
1.2 Our Semistructured Probabilistic Database	1
1.3 Motivation for an XML backend	2
1.4 Outline	3
2 Background	4
2.1 Data Model	4
2.2 SPOQL and SP-Algebra Operations	7
2.2.1 Simple Selection	8
2.2.2 Projection	10
2.2.3 Conditionalization	13
2.2.4 Cartesian Product	13
2.2.5 Join	15
2.2.6 Mix	15
2.3 Data Modification	16
2.3.1 Operations	16
2.4 Related Work	17
3 SPDBMS Design	19
3.1 Existing System Design	19

<i>CONTENTS</i>	vi
3.2 XML Database Adapter Design	21
3.3 XQuery Interface Design	23
4 XQuery API Implementation	25
4.1 Simple Selection	25
4.1.1 OR	25
4.1.2 AND	26
4.2 Projection on Variables	27
4.3 Conditionalization	29
4.4 Cartesian Product	29
4.5 Join	32
4.6 Mix	34
5 Experiments	36
5.1 Design	36
5.1.1 Experimental Variables	36
5.1.2 Construction and Execution	37
5.2 Results and Analysis	40
5.2.1 Module Compilation Speed	40
5.2.2 Results for Projection on Variables	41
5.2.3 Memory Usage	42
6 Conclusion and Future Work	43
7 Bibliography	44

<i>CONTENTS</i>	vii
A XQuery	46
B Experimental Results	68
B.1 Baseline experiments	68
B.2 2 variables, 100 SPOs	73
B.2.1 Simple Selection	73
B.2.2 Select on Context	76
B.2.3 Select on Conditional	80
B.2.4 Select on Variables	84
B.2.5 Select on Table	88
B.2.6 Project on Context	92
B.2.7 Project on Conditional	96
B.2.8 Project on Variables	100
B.2.9 Conditional	104
B.2.10 Conditional	108
B.2.11 Select on Probability	112
B.2.12 Complex Select Conditions	116
B.2.13 Complex Project Conditions	121
C SPO Schema Document	126

List of Figures

1	Example: SPO table format	4
2	Example: SPO XML format	5
3	SP-Algebra and SPOQL operations summary	8
4	Example: projection on context	10
5	Example: projection on conditional	10
6	Example: projection on variables	11
7	Example: conditionalization	12
8	Example: product	14
9	Original SPDBMS architecture	20
10	Where-OR implementation	26
11	Projection on variables implementation	28
12	Conditionalization implementation	30
13	Product implementation	31
14	Left join implementation	33
15	Right join implementation	34
16	Left mix implementation	34
17	Right mix implementation	35
18	Baseline experiment results 1	38
19	Baseline experiment results 2	39
20	Projection on variables experiment results	41

1 Introduction

1.1 Probabilistic Databases

Many applications require storing and querying probabilistic information: stock market prediction, risk analysis software used by insurance companies, and image recognition, to name a few [10] [11]. Such probabilistic data is poorly suited to a traditional relational database - operations common to probabilistic data are not provided by relational databases. A *probabilistic database* allows one to apply transformations to the database's content, based on the laws of probability theory, and perform queries based on these probabilistic transformations. Probabilistic databases have been the subject of much research, and a number of different models have been proposed ([11], [7] and section 2.4).

1.2 Our Semistructured Probabilistic Database

The existing structured probabilistic systems cited above lack the flexibility to work well with real-world data taken from multiple sources, or probabilistic data that otherwise fails to match a common schema [11]. A solution is presented in [9] in the form of a Semistructured Probabilistic Database (SPDB). SPDB is a system for storing probabilistic data as Semistructured Probabilistic Objects, SPOs. An SPO specifies a probability distribution for any number of random variables of finite domain, along with non-probabilistic data; the structure described in detail in section 2.1. Semistructured probabilistic

relations store an arbitrary collection of SPOs; as semistructured data, no common schema is needed for the storage of SPOs in a relation.

SPOs are manipulated via a semistructured probabilistic algebra, also described in [9], through which the probability distribution of a set of SPOs can be queried. SP-algebra is implemented via a Semistructured Probabilistic Object Query Language (SPOQL), an SQL-like language. SPOQL and SP-algebra support - among other operations - conditionalizing SPOs to recalculate probability based on known values, and creation of joint probability distributions from multiple SPOs via join and Cartesian product.

1.3 Motivation for an XML backend

XML is particularly well-suited to the semistructured nature of SPO data. In addition, because SPDB currently presents SPOs to the user in XML form, storing SPO data in a relational database is unnatural and the conversion adds a great deal of overhead. We propose an XML backend to SPDB, implemented with an XML-based database and XQuery, to alleviate performance problems and simplify the system as a whole. Our solution is implemented using the open-source XML database ExistDB.

Our primary contribution is the design, implementation, testing, and performance evaluation of the native XML backend described above. We also implement a distinction between LEFT and RIGHT JOIN operations and implement the MIX operation; these operations are specified in earlier works, but were not implemented until now. We also provide a suite of automated

tests for verifying the correctness of SPDB.

1.4 Outline

Section 2 provides background information on the project, including the semantics of each SPOQL operation. Section 3 describes the design of the SPDBMS server and its integration with the native XML backend, and section 4 goes into further detail on the design and implementation of the XQuery SPDB library. Section 5 discusses the design and results of experiments which test the performance of the new storage.

ω : SPO-1		
university : Cal Poly		
CSC560	CSC305	P
A	A	0.34476
A	B	0.32192
B	A	0.05507
CSC101 = A		

Figure 1: An example of an SPO in table format. “University: Cal Poly” is context information. CSC560 and CSC305 are random variables. The rows of grades for each random variable and the probability for each form the probability table. The precondition of an A in CSC101 is conditional information.

2 Background

2.1 Data Model

SPOs - *Semistructured Probabilistic Objects* - are the structure SPDB uses to represent probabilistic data. An SPO is a tuple consists of four parts: Context, Variables, the Probability Table, and Conditionals [9].

Examples of SPOs in tabular and XML formats are shown in figures 1 and 2, respectively.

More formally, consider the set \mathcal{R} of all possible relational attributes; and \mathcal{V} , the set of all possible random variables and their domains. An SPO is a tuple $S = \langle T, V, P, C \rangle$, where: [11, section 3] [9]

- T is a relational tuple of a *semistructured schema* over \mathcal{R} , all possible relational attributes.
- $V = \{v_1, \dots, v_q\} \subseteq \mathcal{V}$ is the set of *participating random variables* in S ,

```

<spo>
  <context>
    <elem><name>university</name><val>Cal Poly</val></elem>
  </context>
  <table>
    <variable><name>CSC560</name> <name>CSC305</name></variable>
    <row> <val>A</val> <val>A</val> <P>0.34476</P> </row>
    <row> <val>A</val> <val>B</val> <P>0.32192</P> </row>
    <row> <val>B</val> <val>A</val> <P>0.05507</P> </row>
  </table>
  <conditional>
    <elem><name>CSC101</name> <val>B</val></elem>
  </conditional>
</spo>

```

Figure 2: An example of an SPO in XML format. “University: Cal Poly” is context information. CSC560 and CSC305 are random variables. The rows of grades for each random variable and the probability for each form the probability table. The precondition of an A in CSC101 is conditional information. This SPO is equivalent to the example in table format from figure 1.

where \mathcal{V} is the set of all possible random variables, and $V \neq \emptyset$.

- $P : \text{dom}(V) \longrightarrow \mathcal{P}[0, 1]$ is the *probability table* of s . P need not be complete: total probability may be less than one.
- $C = \{(u_1, X_1), \dots, (u_s, X_s)\}$, where $\{u_1, \dots, u_s\} = U \subseteq \mathcal{V}$ and $X_i \subseteq \text{dom}(u_i)$, $1 \leq i \leq n$ such that $V \cap U = \emptyset$, is the *conditional* of S .

To elaborate on this definition, less formally:

- *Participating random variables* are the names of each column seen in figures 1 and 2. *Participating* variables are part of the probability table, and distinct from the variables found within an SPO's *conditionals* but not in the table. Each variable has a finite *domain*, represented as the range of values for each variable in the probability table and conditionals. In the above example, CSC560 and CSC305 are all participating random variables.
- The *probability table* specifies the probability distribution for each set of events. Each *column* in the table is a participating random variable. Each *row* in the table contains a set of values for each variable within the domain of that variable, and the probability that this combination of values will occur.

A probability table is *complete* if it contains a probability for every possible instance - that is, every possible combination of values within the domain of each variable. Incomplete tables may be stored and

queried as well. The table in our example above is not complete - observe that a row for (CSC560=B, CSC305=B) is missing, and other possible values for each random variable (in our example, grades of C, D, and F) may account for other missing rows as well.

- *Conditional* information represents random variables with a known value or restricted to a subset of its domain. When analyzing a probability distribution, we often have prior information on the value of some variables. Our definition represents this restriction as (u, X) , where u is a random variable and X is a subset of the domain of u . Unlike context, conditionals are subject to change: *conditionalization* queries may add new conditions to the SPO.
- *Context* is any supporting information for this probability distribution. Context may contain any set of name-value pairs whose values are certain, and these values has no effect on the probability calculations performed elsewhere. In the example of figures 1 and 2, *university* is part of the context - it is known in advance, not a random variable, and plays no part in probability calculations.

2.2 SPOQL and SP-Algebra Operations

SPOQL - Semistructured Probability Object Query Language - is an SQL-like language used to extract information from the database. Like SQL, SPOQL queries have *no side effects* - the operations discussed below do not

Operation	Section	SP-Algebra	SPOQL Syntax
Selection	2.2.1	$\sigma_c(S)$	SELECT * FROM S WHERE c
Projection	2.2.2	$\pi_f(S)$	SELECT f FROM S
Conditionalization	2.2.3	$\mu_d(S)$	SELECT * FROM S CONDITIONAL d
Cartesian Product	2.2.4	$S_1 \times S_2$	SELECT * FROM S_1 TIMES S_2
Left Join	2.2.5	$S_1 \times S_2$	SELECT * FROM S_1 [LEFT] JOIN S_2
Right Join	2.2.5	$S_1 \times S_2$	SELECT * FROM S_1 RIGHTJOIN S_2
Left Mix	2.2.6	$S_1 \otimes_L S_2$	SELECT * FROM S_1 [LEFT]MIX S_2
Right Mix	2.2.6	$S_1 \otimes_R S_2$	SELECT * FROM S_1 RIGHTMIX S_2

Figure 3: A summary of conceptual operations on SPOs, the relevant section in this paper, the SP-algebra symbols, and the corresponding SPOQL syntax. For left join/left mix, SPOQL in [brackets] is optional: left is the default direction if none is specified.

modify the contents of the database.

SP-Algebra, specified in [11], is the theoretical basis for SPOQL; similar to relational algebra as a basis for SQL. SPDBMS input is in the form of SPOQL. In the following sections, we'll outline each SP-algebra operation.

The table in figure 3 summarizes each operation, including the SP-algebra symbol, SPOQL syntax, and section in which it is discussed.

2.2.1 Simple Selection

The familiar SELECT statement is used to view the SPOs stored in a given relation.

WHERE statements, also familiar from SQL, can be used to restrict the set of SPOs viewed. WHERE conditions may be based on equality or comparisons with the following fields:

- Selection on Variables: Show only SPOs that have the given variable somewhere in their probability table, with any value.

```
SELECT * FROM relation WHERE var.CSC101 IN V
```

- Selection on Table: Show only SPOs where the given variable exists with the given value.

```
SELECT * FROM relation WHERE tbl.CSC101 = A
```

- Selection on Probability: Show only rows in the probability table where probability is within the given bounds.

```
SELECT * FROM relation WHERE tbl.prob > 0.10
```

- Selection on Conditionals: Show only SPOs where the given conditional exists with the given value.

```
SELECT * FROM relation WHERE cnd.CSC101 = A
```

- Selection on Context: Show only SPOs containing the given context element.

```
SELECT * FROM relation WHERE cnt.year = 1999
```

Multiple WHERE conditions may be combined using AND and OR operators:

```
SELECT * FROM relation WHERE cnt.year = 1999 AND cnd.CSC101 = A
```

$\omega: S$		
university : Cal Poly		
department : CSC		
CSC560	CSC305	P
A	A	0.34476
A	B	0.32192
B	A	0.05507
CSC101 = A		
CSC102 = A		

$\omega: \pi_{cnt.university}(S)$		
university : Cal Poly		
CSC560	CSC305	P
A	A	0.34476
A	B	0.32192
B	A	0.05507
CSC101 = A		
CSC102 = A		

Figure 4: Example of the projection operation on context information. The context information 'department' is removed and 'university' is selected.

$\omega: S$		
university : Cal Poly		
department : CSC		
CSC560	CSC305	P
A	A	0.34476
A	B	0.32192
B	A	0.05507
CSC101 = A		
CSC102 = A		

$\omega: \pi_{cnd.CSC101}(S)$		
university : Cal Poly		
department : CSC		
CSC560	CSC305	P
A	A	0.34476
A	B	0.32192
B	A	0.05507
CSC101 = A		

Figure 5: Example of the projection operation on a conditional. The conditional 'CSC101' is selected and 'CSC102' is removed.

2.2.2 Projection

Projection simplifies an SPO by returning the SPO with only a subset of its original variables, conditionals, or context values.

Projection on Conditional and Projection on Context Conditional and context projection are easy to understand: given an initial SPO, these select only a subset of that SPO's conditional or context tuples, respectively. See the examples in figures 4 and 5.

$\omega: S$		
university : Cal Poly		
CSC560	CSC305	P
A	A	0.34476
A	B	0.32192
B	A	0.05507
CSC101 = A		

(Step 1)	
$\omega: \pi_{var.CSC305}(S)$	
university : Cal Poly	
CSC305	P
A	0.34476
B	0.32192
A	0.05507
CSC101 = A	

(Step 2)	
$\omega: \pi_{var.CSC305}(S)$	
university : Cal Poly	
CSC305	P
A	0.39983
B	0.32192
CSC101 = A	

Figure 6: Example of the projection operation on a variable, shown in two steps for clarification. Step 1 removes all variables not selected by the projection; step 2 merges all duplicate rows. In step 1, observe the two rows whose events are identical: both have only a single variable with value 'A'. These rows are merged in step 2 by summing their probabilities.

Projection on Variables Much like relational projection, SPO projection on variables modifies the SPO by selecting only a subset of its columns. Projection on variables may not result in duplicate rows - if all values for two rows are identical after a projection, the rows are merged and their probabilities are summed. The example in figure 6 illustrates this process.

Unlike the conditionalization operation (2.2.3), an SPO's conditionals are not modified by projection on a variable.

$\omega: S$		
university : Cal Poly		
CSC560	CSC305	P
A	A	0.300
A	B	0.100
B	A	0.200
CSC101 = A		

(Step 1)		
$\omega: \mu_{var.CSC560=A}(S)$		
university : Cal Poly		
CSC560	CSC305	P
A	A	$0.300 * (0.600/0.400)$
A	B	$0.100 * (0.600/0.400)$
B	A	—
CSC101 = A		
CSC560 = A		

(Step 2)		
$\omega: \mu_{var.CSC560=A}(S)$		
university : Cal Poly		
CSC560	CSC305	P
A	A	0.450
A	B	0.150
CSC101 = A		
CSC560 = A		

(Final step)	
$\omega: \mu_{var.CSC560=A}(S)$	
university : Cal Poly	
CSC305	P
A	0.450
B	0.150
CSC101 = A	
CSC560 = A	

Figure 7: Example of the conditionalization operation, shown in four steps for clarity. In the first step, we add the new conditional value 'CSC560 = A' and show our calculations for new probabilities for each existing row. Probabilities are scaled such that the total probability after the conditionalization is equal to the total before the conditionalization: in this example, the pre-conditionalized total is 0.6 and the intermediate total is 0.400 (0.300 + 0.100), so the remaining rows' probability is scaled up by $(0.6 / 0.4) = 1.5$. Step 2 removes rows that are incompatible with the new conditional (here, rows where $CSC560 \neq A$). In the final step, the conditionalized column is finally removed from the probability table.

2.2.3 Conditionalization

Conditionalization fixes the value of a given variable. This has the effect of removing columns from the probability table, like projection, but the variable is not removed from the simulation - it is instead moved to the *conditional* section with the specified value. Rows without the specified value are removed from the probability table - after fixing the value for this variable, these rows need not be analyzed - and the probability of all remaining rows is normalized. An example of this process is shown in figure 7.

If a `CONDITIONAL` operation on variable V is applied to an SPO which does not reference V , that SPO's `conditional` data is updated, but the SPO's probability table is unchanged. This makes sense: requiring a condition to an irrelevant variable should have no effect on the probability of each result.

The sum of probabilities for a complete table is unchanged after conditionalization. For complete tables, the sum of probabilities is always 1. Incomplete tables may be conditionalized as well.

2.2.4 Cartesian Product

A Cartesian product constructs a joint probability distribution from two *product-compatible* SPOs. Two SPOs are product-compatible if they have matching conditionals, and have no random variables in common [11]. The result looks similar to a Cartesian product in a relational database - each pair of rows in the probability table is returned. One important difference from relational databases is, of course, the probability of each row. Probabilities

$\omega: S_1$ <hr/> <table style="width: 100%; border-collapse: collapse;"> <tr> <td style="width: 50%; padding: 2px;">university :</td> <td style="padding: 2px;">Cal Poly</td> </tr> <tr> <td style="border-top: 1px solid black; padding: 2px;">CSC560</td> <td style="border-top: 1px solid black; padding: 2px;">P</td> </tr> <tr> <td style="padding: 2px;">A</td> <td style="padding: 2px;">0.4</td> </tr> <tr> <td style="padding: 2px;">B</td> <td style="padding: 2px;">0.5</td> </tr> <tr> <td colspan="2" style="border-top: 1px solid black; padding: 2px;">CSC101 = A</td> </tr> </table>	university :	Cal Poly	CSC560	P	A	0.4	B	0.5	CSC101 = A		$\omega: S_2$ <hr/> <table style="width: 100%; border-collapse: collapse;"> <tr> <td style="width: 50%; padding: 2px;">university :</td> <td style="padding: 2px;">Cal Poly</td> </tr> <tr> <td style="border-top: 1px solid black; padding: 2px;">CSC305</td> <td style="border-top: 1px solid black; padding: 2px;">P</td> </tr> <tr> <td style="padding: 2px;">A</td> <td style="padding: 2px;">0.1</td> </tr> <tr> <td style="padding: 2px;">B</td> <td style="padding: 2px;">0.8</td> </tr> <tr> <td colspan="2" style="border-top: 1px solid black; padding: 2px;">CSC101 = A</td> </tr> </table>	university :	Cal Poly	CSC305	P	A	0.1	B	0.8	CSC101 = A					
university :	Cal Poly																								
CSC560	P																								
A	0.4																								
B	0.5																								
CSC101 = A																									
university :	Cal Poly																								
CSC305	P																								
A	0.1																								
B	0.8																								
CSC101 = A																									
<table border="1" style="margin: auto; border-collapse: collapse;"> <tr> <td colspan="3" style="padding: 5px;">$\omega: S_1 \times S_2$</td> </tr> <tr> <td colspan="3" style="border-top: 1px solid black; padding: 5px;">university : Cal Poly</td> </tr> <tr> <td style="border-top: 1px solid black; padding: 2px;">CSC560</td> <td style="border-top: 1px solid black; padding: 2px;">CSC305</td> <td style="border-top: 1px solid black; padding: 2px;">P</td> </tr> <tr> <td style="padding: 2px;">A</td> <td style="padding: 2px;">A</td> <td style="padding: 2px;">0.04</td> </tr> <tr> <td style="padding: 2px;">A</td> <td style="padding: 2px;">B</td> <td style="padding: 2px;">0.32</td> </tr> <tr> <td style="padding: 2px;">B</td> <td style="padding: 2px;">A</td> <td style="padding: 2px;">0.05</td> </tr> <tr> <td style="padding: 2px;">B</td> <td style="padding: 2px;">B</td> <td style="padding: 2px;">0.40</td> </tr> <tr> <td colspan="3" style="border-top: 1px solid black; padding: 2px;">CSC101 = A</td> </tr> </table>		$\omega: S_1 \times S_2$			university : Cal Poly			CSC560	CSC305	P	A	A	0.04	A	B	0.32	B	A	0.05	B	B	0.40	CSC101 = A		
$\omega: S_1 \times S_2$																									
university : Cal Poly																									
CSC560	CSC305	P																							
A	A	0.04																							
A	B	0.32																							
B	A	0.05																							
B	B	0.40																							
CSC101 = A																									

Figure 8: Example of the Cartesian product operation. We begin with two SPOs for this operation: S_1 and S_2 . These are product compatible: *conditionals are identical* and they have *no common participating random variables*. The result contains all participating random variables from both SPOs. Rows are a Cartesian product of rows from the original SPOs, much like a Cartesian product in a relational database; probability for each row is the product of the probabilities from each row that created it.

for each pair of rows are multiplied to find the probability used in the product row. This operation is illustrated in figure 8.

2.2.5 Join

A join is similar to a Cartesian product - it constructs a joint probability distribution from two *join-compatible* SPOs. Two SPOs are join-compatible if their conditionals match, as with product-compatibility; and they must have *at least one* common random variable [11].

To perform a join, one of the two SPOs is *conditionalized* (Conditionalization is explained in section 2.2.3.) Results differ depending on which side is conditionalized. We may specify which side is conditionalized by using LEFTJOIN or RIGHTJOIN.

2.2.6 Mix

Cartesian product and join are mutually exclusive. Mix chooses the appropriate operation for a pair of SPOs: if the pair is join-compatible, mix executes a join; if the pair is product-compatible,

While a pair of SPOs cannot be *both* product- and join-compatible, note that a pair can be *neither* product- or join-compatible if their conditionals are not equal.

As discussed in the previous section, JOIN may have a direction: LEFTJOIN or RIGHTJOIN. Since MIX may execute joins, we specify the direction of joins by using two types of mixes as well: LEFTMIX and RIGHTMIX. The direction

of a mix is irrelevant for a pair of product-compatible SPOs; product does not have a direction.

2.3 Data Modification

All SPOs in the database are stored in a parent *relation*. Each relation is a named grouping of any set of SPOs. Any two SPOs may be part of the same relation - because SPOs are *semistructured* data, no schema must be enforced.

TIMES, JOIN, and MIX operations (see sections 2.2.4, 2.2.5, and 2.2.6 respectively) can operate on any two relations, provided the SPOs each relation contains are product-compatible or join-compatible.

2.3.1 Operations

- CREATE Creates a new relation in the database. A name for the relation must be provided. Relations are empty when created - use *INSERT* to populate relations with SPO data.
- INSERT Adds new SPO data to a given relation. SPO data is provided via an XML file.
- DELETE Remove all SPOs from a relation matching the given WHERE conditions (see section 2.2.1).
- DROP Removes an existing relation, given the name of the relation to drop. All SPOs in the relation are destroyed.

2.4 Related Work

SPDBMS has been the subject of much past work. Zhao and Dekhtyar defined SP-algebra and implemented the original system, including its relational backend [11]. Mathias specified the SPOQL language and implemented it on top of Zhao and Dekhtyar’s work [9].

Numerous models for probabilistic relational databases have been developed[11]. A data model proposed by Barbará, Garcia-Molina and Porter [5] extends the relational model to support uncertain data. Each relation must have a deterministic key, and all other attributes may be either deterministic (certain) or stochastic (uncertain). Stochastic attributes have an associated probability. Attributes may also be specified as independent or interdependent: the probability of an independent attribute is not tied to that of any other attribute; the probability of an interdependent attribute is tied to another attribute.

Cavallo and Pittarelli [8] incorporate probability into relational databases as a tuple (V, Δ, dom, p) : V is a set of unique attributes, δ is the set of attribute domains, dom provides the mapping $V \rightarrow \Delta$, and p is a probability distribution over V . Their model required a sum of 1 for the probabilities of all tuples in a relation. They define projection and join operators.

MystiQ [6] is a relatively recent probabilistic database system, implemented as middle ware on a standard relational database. MystiQ manipulates probabilistic data stored in a standard relational database; however, MystiQ has a different concept of probabilistic data than our work. MystiQ issues queries on a standard relational database using ‘fuzzy’ queries, where

query restrictions are not completely defined. *MystiQ* will return standard relational rows as results, each row paired with a probability that it was what the user wanted.

These probabilistic databases based on the relational model have inherent weaknesses. Many works describe a single probabilistic object as a single database row with attributes whose values are uncertain. We instead represent a probabilistic object as a set of random variables that form a probability distribution, where the probability of multiple combinations of values can be specified. Another strength of our work over related work is its semistructured nature - SPOs from a single relation, representing the same type of object, need not have identical schema. This allows us to import and work with data from multiple sources with varying schema.

3 SPDBMS Design

3.1 Existing System Design

The original SPDB server design, as described in [9, section 2.2], is shown in figure 9. The SPDB server stores all SPO data in a relational database. Clients connect to the server via TCP/IP, and manipulate stored data or issue SPOQL queries. The SPO Request Dispatcher parses SPOQL queries into SP-algebra strings, further parses SP-algebra strings into an SP-algebra based abstract tree, and passes the final structure to the database adapter for execution. The database adapter constructs DBMS-specific SQL based on this structure to either update the SPOs/relations as requested by an `INSERT`, `DELETE`, `CREATE`, or `DROP` operation; or the SQL required to construct the SPO XML returned by the application for a `SELECT` operation [2]. In the original design, these operations may require the execution of any number of SQL statements against the underlying database.

Decoupling the database adapter from the rest of the application in this design allows new database backends to be constructed with minimum effort - most application logic is not tied to the DBMS backend.

The implementation of this database adapter design required a JDBC interface to the database. All other details of database access are left to the database adapter implementer.

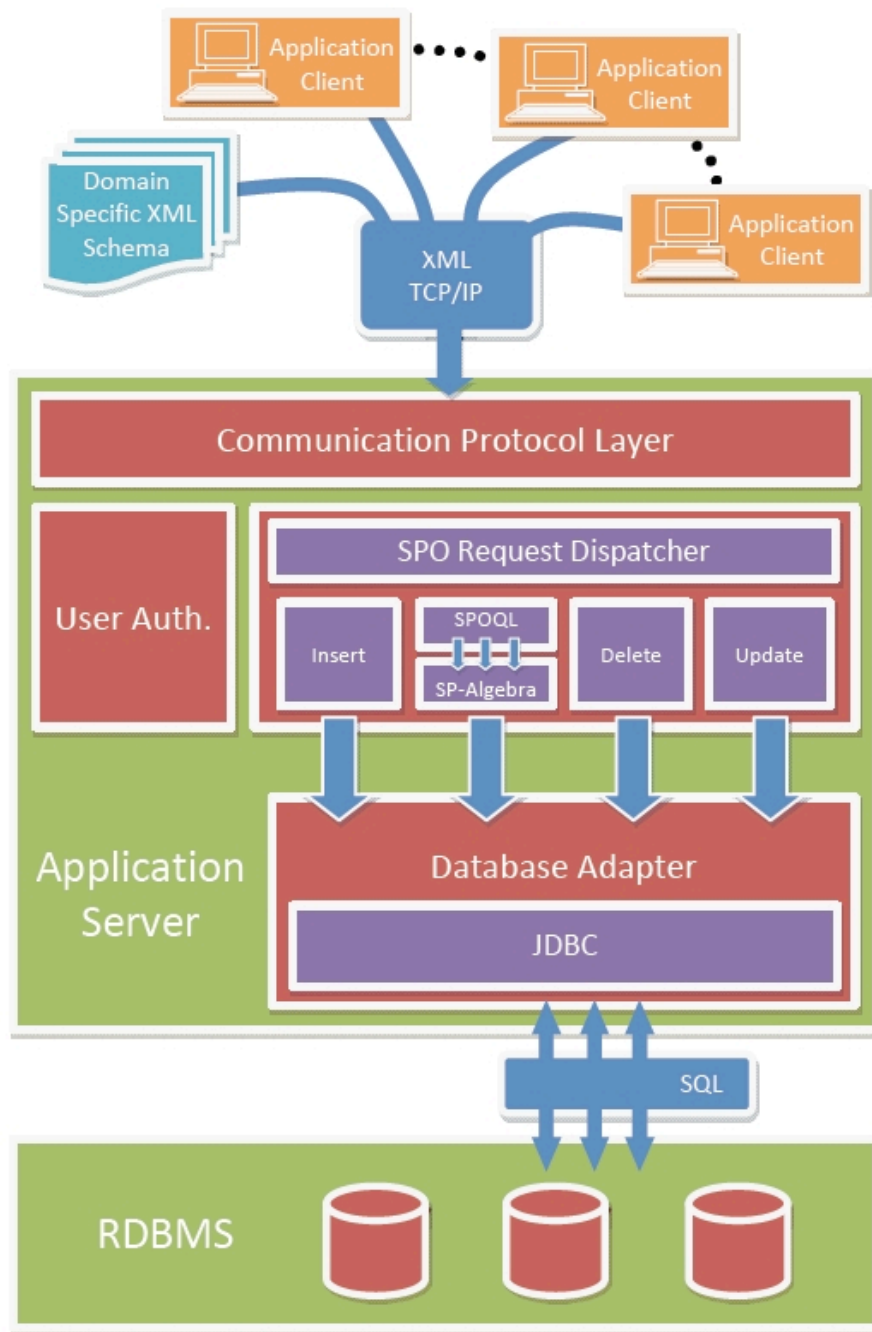


Figure 9: The original SPDBMS system architecture, as described in [9, section 2.2]. (Image source: [2, figure 2])

3.2 XML Database Adapter Design

The SPDB server design (figure 9) specifies a relational database beneath the database adapter, but the implementation does not enforce this; we generalize this design to allow XML databases.

The design for our Exist-DB adapter for SPDB consists of two stages:

- Query planning: given the abstract query tree sent to the database adapter from the query parser, construct a single XQuery statement.
- Query execution: execute the planned XQuery statement.

This is modeled after the query planning and execution stages employed by many relational databases. The advantages of this separation have been shown many times for relational databases, and these advantages (discussed below) are applicable here as well. ¹

One advantage of this approach is that *only one* XQuery statement is constructed per SPOQL query. This minimizes the overhead required by database round-trips. This is a dramatic improvement over the current Oracle implementation's ad-hoc approach to querying the underlying database, where many round-trips are required per SPOQL query. An advantage of a native XML backend is that no post-processing of data returned from XQuery is necessary - SPDB returns the same XML data that XQuery excels at processing.

¹In particular, we appreciated having query execution decoupled from query planning when creating automated tests for query planning.

All SPO-related logic is processed by the XQuery backend as XQuery functions, rather than Java-based SPO server. This has advantages similar to stored procedures in a relational environment: among other things, the caller need not be aware of the XQuery implementation. Removing this logic from the caller also allows us to take full advantage of XQuery's ability to manipulate semistructured XML data.

All XQuery commands sent to the server begin with the following, in order to access this XQuery function library

```
import module spdb="http://www.csc.calpoly.edu/~erosson/spdb"
as "resource:spdb.xqm";
```

`spdb.xqm` is the name of our XQuery function library, implemented as an XQuery module. `http://www.csc.calpoly.edu/~erosson/spdb` is the XML namespace required for the module. According to [1] this module is only compiled once, and cached for all queries thereafter. (See our experimental results, section 5.2.1 for further discussion of this.)

In addition, the strategy described above greatly simplifies query planning. For example, consider the following SPOQL statement:

```
SELECT * FROM relation1 LEFTJOIN relation2
```

The SPOQL query planner translates this to reasonably concise and readable XQuery:

```
import module spdb="http://www.csc.calpoly.edu/~erosson/spdb"
```

```
as "resource:spdb.xqm";
spdb:leftjoin(spdb:relation("relation1"), spdb:relation("relation2"))
```

The complex logic behind JOIN is abstracted behind an interface of XQuery functions.

3.3 XQuery Interface Design

Each SP-Algebra operation specified in section 2 is implemented as an XQuery function. These XQuery functions provide all of the system's probability analysis features - the SPDB server merely translates SPOQL queries into the matching XQuery calls.

All functions in the XQuery API - that is, any function referenced from the SPDB server corresponding to an SP-Algebra operation - will always accept input and require output of the same form:

```
<spos>
  <spo>
    ...data...
  </spo>
  <spo>
    ...data...
  </spo>
  ...more SPOs...
</spos>
```

This can be verified by examining the method signature for any XQuery function called from SPDB:

```
declare function spdb:sp-some-function($spos as element(spos), ...)
as element(spos);
```

This signature allows the output of any one function to be passed to another. Arbitrary complex SPOQL queries can be planned before any results are computed.

4 XQuery API Implementation

See appendix A for the source of the XQuery module. All functions are written by Evan Pierce Rosson unless another author is noted. To summarize, Dustin Anderson implemented all simple selection WHERE-clauses (section 2.2.1), and projection on context and conditionals (2.2.2). Evan Pierce Rosson implemented all other operations: the current projection on variables (2.2.2); AND and OR conditions used in WHERE-clauses; CONDITIONAL, JOIN, TIMES, and MIX operations (2.2.3, 2.2.4, 2.2.5, and 2.2.6 respectively).

Note that there are a number of auxiliary functions present. Functions called by the SPDB server are differentiated from helper functions only by function name - 'public' functions have names beginning with 'sp-'. All other functions are exposed only to allow for thorough unit testing.

The following sections will discuss the implementation of each function implemented by Evan Pierce Rosson, including relevant code.

4.1 Simple Selection

4.1.1 OR

The implementation of OR is straightforward. Each parameter of the OR operation is evaluated, and the union of the results is returned. To avoid returning duplicates - SPOs included in both arguments to OR - the built-in function `distinct-deep` performs a distinct union - SPOS present in both

```

1 declare function spdb:sp-where-or($spos1 as element(
    spos), $spos2 as element(spos)) as element(spos)
    {
2   element spos {$spos1/@*, functx:distinct-deep((
    $spos1/spo, $spos2/spo)) }
3 };

```

Figure 10: WHERE condition: OR implementation. “WHERE X OR Y” becomes “sp-where-or(X, Y)”

arguments are included only once.

4.1.2 AND

AND could have been implemented in our XQuery library as an intersection. Instead, we chose a solution that required no XQuery function implementation, and a simple implementation in the SPDBMS server. The SPOQL query `SELECT * FROM Relation WHERE X AND Y` is translated to XQuery similar to `spdb:X(spdb:Y(spdb:relation("Relation")))`. The requirements of WHERE filter X are applied only to the results of filter Y, not the entire relation. This leads to a more efficient implementation and more concise XQuery.

In the above example, the XQuery statements `spdb:Y(spdb:X(spdb:relation("Relation")))` and `spdb:X(spdb:Y(spdb:relation("Relation")))` are equivalent. One of these may be faster than the other: we would prefer to execute the faster, more selective (fewer results) condition first. This minimizes the number of re-

sults that need to be processed by the slower, less selective outer function. Currently we do not attempt this optimization: order of functions in the generated XQuery is based solely on the order of `WHERE` conditions in the SPOQL input. Estimating the cost of each function for this sort of optimization is a good opportunity for future work (section 6).

4.2 Projection on Variables

The general algorithm used in our implementation of projection on variables is:

- Remove all projected columns from the result.
- Merge all *duplicate rows* by summing their probabilities.

Recall from section 2.2.2 that projection on variables must not create duplicate rows. Rows that are identical (except for probability) after removing a column are to be merged by summing their probabilities into a single row.

Identifying duplicate rows proved to be an expensive operation in XQuery. Given n rows, every possible *pair of rows* must be compared, for a total of $\frac{n^2}{2}$, comparisons.

Oracle uses a hash table, with hashed rows as keys, to achieve linear complexity. XQuery, however, lacks the concept of a hash table. We hoped ExistDB's XQuery processor would be able to optimize this case. This, and further optimization possibilities, are further discussed in our experimental results, section 5.2.2.

```

1 declare function spdb:rows-project-var($rawrows as
    element(row)*,
2 $visiblenums as xs:integer*) as element(row)*
3 {
4   let $rowsviscols := for $row in $rawrows
5     return element row {$row/@*,
6       $row/val[position() = $visiblenums],
7       $row/P
8     }
9   return for $row1 at $i in $rowsviscols
10     (: where no previous rows have identical values
11       :)
11     where fn:empty($rowsviscols[position() lt $i][
12       not(val != $row1/val)])
12     return element row {$row1/@*,
13       $row1/val,
14       element P {$row1/P/@*,
15         (: sum all rows with identical values. skip
16           previous rows: we checked those already
17           :)
16         sum(for $p in $rowsviscols[position() ge $i
17           ][not(val != $row1/val)]/P
17           return xs:decimal($p))
18       }
19     }
20 };

```

Figure 11: Implementation of projection on variables.

4.3 Conditionalization

Our XQuery conditionalization implementation is shown in figure 12.

Notice that it's possible for a conditionalization to change an SPO's `<conditional>` data without modifying its probability table, as described in 2.2.3. `conditionalize-table` will return an unmodified probability table if the conditionalized variable is not present in the table. However, its caller `conditionalize` will always modify the SPO's conditional information, even if the probability table is unchanged.

At first, conditionalization appears very similar to projection on variables (sections 2.2.2, 4.2) - both remove a column and multiple rows from the table. Conditionalization, however, has no need to identify and merge duplicate rows. Projection on variables must compare against other rows to determine duplicates; conditionalization simply chooses which rows to remove based on the new condition, and has no need to compare with other rows in the same table. Thus, conditionalization does not have the same complexity problems as projection on variables, and is expected to run in linear time.

4.4 Cartesian Product

The implementation of Cartesian product ("TIMES"), specified in section 2.2.4, is shown in figure 13.

Note that `TIMES` has two double `FOR` loops: one iterating each possible pair of SPOs; another iterating each possible pair of rows in each SPO's table. If

```

1 declare function spdb:conditionalize-table($table as
    element(table), $name as xs:string, $value as xs
    :string) as element(table) {
2   let $num := spdb:table-column-num($table/variable/
    name, $name)
3   (: If var isn't in table, return the original
    table.
4   This isn't silent failure, but an irrelevant
    conditionalization. :)
5   return if (fn:empty($num)) then $table else
6   let $newrows := $table/row[exists(val[$num][text()
    eq $value])]
7   let $oldprob := sum(for $i in $table/row/P return
    xs:decimal($i))
8   let $newprob := sum(for $i in $newrows/P return xs
    :decimal($i))
9   let $probmult := if ($newprob gt 0) then $oldprob
    div $newprob else 1
10  return element table {
11    element variable {$table/variable/name[position
    () ne $num]},
12    for $row in $newrows
13      let $p := xs:decimal($row/P) * $probmult
14      return element row {$row/val[position() ne
    $num], element P{$p}}
15  }
16 };
17
18 declare function spdb:conditionalize($spo as element
    (spo), $name as xs:string, $value as xs:string)
    as element(spo) {
19  element spo {$spo/@*,
20    $spo/context,
21    spdb:conditionalize-table($spo/table, $name,
    $value),
22    element conditional {
23      $spo/conditional/*,
24      element elem {element name{$name}, element val
    {$value}}
25    }
26  }
27 };

```

Figure 12: Conditionalization implementation.

```

1 declare function spdb:table-product($table1 as
    element(table), $table2 as element(table)) as
    element(table) {
2   if (not(spdb:common-vars-empty($table1, $table2)))
        then () else
3   element table {
4     element variable {$table1/variable/name, $table2
        /variable/name},
5     spdb:rows-product($table1/row, $table2/row)
6   }
7 };
8
9 declare function spdb:rows-product($rows1 as element
    (row)*, $rows2 as element(row)*) as element(row)*
    {
10  for $row1 in $rows1
11    for $row2 in $rows2
12      let $prob := element P { $row1/P/number() *
        $row2/P/number() }
13      return element row {$row1/val, $row2/val,
        $prob}
14 };
15
16 declare function spdb:sp-product($spos1 as element(
    spos), $spos2 as element(spos)) as element(spos)
    {
17  element spos {
18    for $spo1 in $spos1/spo
19      for $spo2 in $spos2/spo
20        where (spdb:product-compatible($spo1, $spo2)
        )
21        return spdb:product($spo1, $spo2)
22  }
23 };

```

Figure 13: Partial Cartesian product ("TIMES") implementation.

all SPOs are product-compatible and have many rows, this operation can take a very long time. These nested loops are difficult to avoid: the specification of `TIMES` (section 2.2.4) requires that it return $n*m$ SPOs given two relations of size n and m ; and $r*s$ rows per SPO given a pair of SPOs with row counts r and s . These nested loops perform exactly this many iterations.

4.5 Join

Our `JOIN` implementation is shown in figure 14. `JOIN` has complexity problems similar to those of `TIMES`. Two double `FOR` loops are required, for each pair of SPOs and each pair of rows.

As with `TIMES`, all SPOs are join-compatible and have many rows, this operation can take a very long time. These nested loops are difficult to avoid: the specification of `JOIN` (section 2.2.5) requires that it return $n*m$ SPOs given two relations of size n and m ; and $r*s$ rows per SPO given a pair of SPOs with row counts r and s . These nested loops perform exactly this many iterations.

Our implementation removes all common columns from the `JOINED` SPOs by conditionalizing on variables common to both SPOs (`spdb:leftjoin-conditionalize()`, figure 14). Without any common columns, the two SPOs are product-compatible and we perform a Cartesian product.

The implementation of `RIGHTJOIN` is shown in figure 17: the parameters are simply switched to form a left join. This means variables from the second parameter to a `RIGHTJOIN` are displayed before the variables from the first


```

1 declare function spdb:leftjoin-conditionalize-table(
    $table as element(table), $vars as element(name)
    *, $vals as element(val)*) as element(table) {
2   if (fn:empty($vars)) then $table else
3   let $var := $vars[1]
4   let $val := $vals[1]
5   return spdb:leftjoin-conditionalize-table(element
        table {$table/@*,
6         element variable {$table/variable/name[. ne $var
            ]}],
7         spdb:conditionalize-table($table, $var, $val)/
            row
8   }, $vars[position() ne 1], $vals[position() ne 1])
9 };
10
11 declare function spdb:table-leftjoin($table1 as
    element(table), $table2 as element(table)) as
    element(table)? {
12   let $common := spdb:common-vars($table1/variable,
        $table2/variable)/name
13   return if (fn:empty($common)) then () else
14   element table {$table1/@*,
15         element variable {$table1/variable/name, $table2
            /variable/name[not(. = $common)]}],
16         for $row1 in $table1/row
17         let $values := spdb:row-column-vals($table1/
            variable/name, $row1, $common/text())
18         let $condtable2 := spdb:leftjoin-
            conditionalize-table($table2, $common,
                $values)
19         return spdb:rows-product($row1, $condtable2/
            row)
20   }
21 };

```

Figure 14: Partial JOIN implementation. Some parts are very similar to TIMES and have been omitted from this figure.

```

1 declare function spdb:sp-rightjoin($spos1 as element
    (spos), $spos2 as element(spos)) as element(spos)
    {
2   spdb:sp-leftjoin($spos2, $spos1)
3 };

```

Figure 15: Right-join implementation, based on left-join.

```

1 declare function spdb:table-leftmix($table1 as
    element(table), $table2 as element(table)) as
    element(table) {
2   let $common := spdb:common-vars($table1/variable,
    $table2/variable)/name
3   return if (fn:empty($common))
4   then spdb:table-product($table1, $table2)
5   else spdb:table-le      $table1, $table2)
6 };

```

Figure 16: Mix implementation.

parameter in the probability table, which may be unintuitive. This is acceptable, however, as the order in which variables are displayed in the results is not significant.

The names `LEFTJOIN` and `RIGHTJOIN` are used in SPOQL queries, instead of `LEFT JOIN` and `RIGHT JOIN`, to avoid making large changes to the existing SPOQL parser.

4.6 Mix

As described in section 2.2.6, `MIX` simply chooses between `JOIN` and `TIMES` for each pair of SPOs, depending on their compatibility. Note that, while

```
1 declare function spdb:sp-rightmix($spos1 as element(  
    spos), $spos2 as element(spos)) as element(spos)  
    {  
2   spdb:sp-leftmix($spos2, $spos1)  
3  };
```

Figure 17: Right-mix implementation, based on left-mix.

JOIN and TIMES are mutually exclusive, there exist SPOs which are neither join- nor product-compatible that will not be included in MIX's results.

The implementation of RIGHTMIX is similar to that of RIGHTJOIN (4.5), as shown in figure 17.

The names LEFTMIX and RIGHTMIX are used in SPOQL queries, instead of LEFT MIX and RIGHT MIX, for the same reason as the similar care for JOIN operations (section 4.5): to avoid making large changes to the existing SPOQL parser.

5 Experiments

5.1 Design

A major motivation for our native XML SPDBMS backend was improving system performance: the existing Oracle implementation was reported to have trouble with both memory and processing time. We need to justify performance improvements with appropriate experiments. Our experiments focus on measuring the time required to process queries.

5.1.1 Experimental Variables

Factors we planned to vary in experiments include:

- Relation size: The number of SPOs included in each relation
- Number of random variables: The average number of variables in each SPO
- Domain size: the typical number of possible values for each random variable
- Join-compatibility and Product-Compatibility: the number of results expected for JOIN, TIMES, and MIX operations. (Recall the explanations of join- and product-compatibility from sections 2.2.5 and 2.2.4.)

Later, we discovered that execution speed was also dependent on our *XQuery module size* as it influenced the time required to compile the module

for every query, and we set up further experiments to show this. See section 5.2.1 for details.

5.1.2 Construction and Execution

A custom XML generator was written to provide test data for experiments varying the above parameters. An existing set of experimental data and queries, provided with the existing source code from [9], was included and adapted to our needs.

To collect timing information, the SPDB server was instrumented to log two pieces of data - time spent performing database queries, and time spent processing the data outside of the database. All times recorded are in milliseconds. Graphs are generated using Gnuplot.

Experiments are automated by a collection of custom shell scripts. Given a set of test data, a list of files containing experimental queries, and a properties file specifying which database backend to use (ExistDB or Oracle): we start the SPDB server, clean out the database, populate the database with test data required throughout the experiment, and run all queries specified for the experiment. After experiments are complete, graphs of the results are generated using custom Gnuplot scripts.

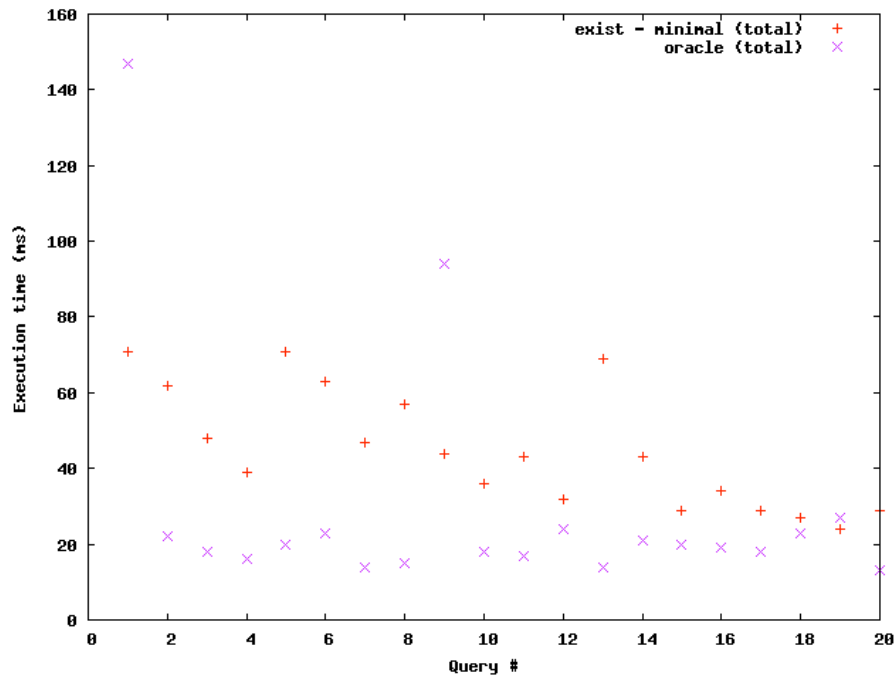


Figure 18: Baseline experiments: running the simplest possible query (SELECT * FROM Relation) on an empty database. Exist tends to lag behind Oracle a bit, but the difference is reasonable.

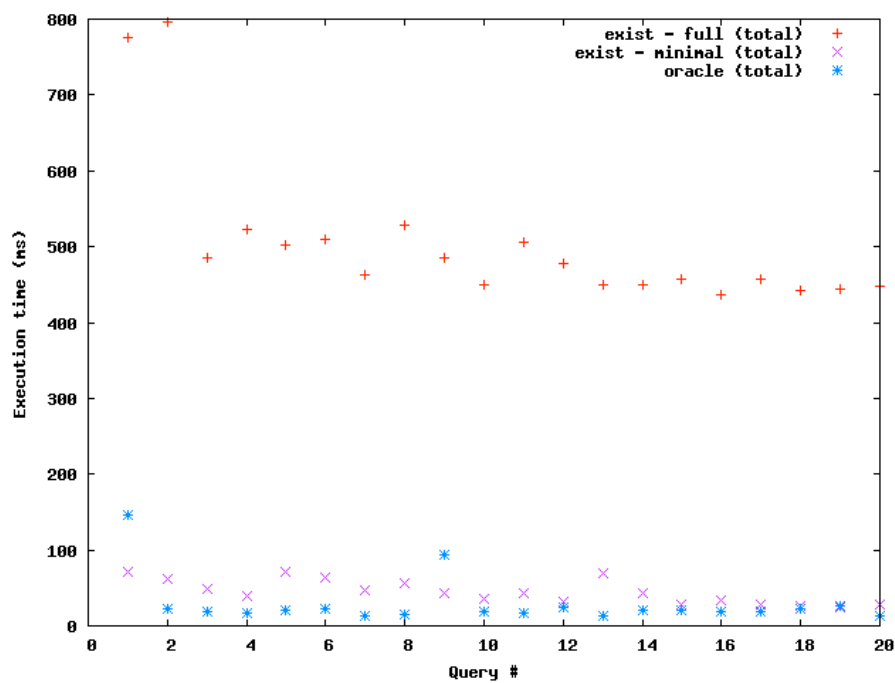


Figure 19: Baseline experiments: running the simplest possible query (SELECT * FROM Relation) on an empty database. "Exist(minimal)" uses the smallest possible XQuery module to run this set of experiments, whereas "Exist(full)" uses the full module supporting all SPDBMS operations. Using the full XQuery module dramatically slows Exist.

5.2 Results and Analysis

5.2.1 Module Compilation Speed

Early experiments showed our ExistDB backend to be dramatically slower than Oracle, even for a baseline measurement - the simplest possible operation, `SELECT * FROM Relation` performed on an empty database. Further experiments, shown in figures 18 and 19, showed that reducing our XQuery module to only a minimum set of functions needed to support this operation dramatically improved our results. ExistDB appeared to be recompiling our XQuery module once for *every XQuery execution*, or once per SPOQL statement!

This module does not change between runs and should only need to be compiled once to be used for all queries. ExistDB's built-in function libraries suggest that this is possible, but we were unable to determine how to implement this ourselves. ExistDB's documentation suggests that this should be done automatically[1]:

XQuery modules executed via the REST interface, the XQueryServlet or XQueryGenerator are automatically cached: the compiled expression will be added to an internal pool of prepared queries. The next time a query or module is loaded from the same location, it will not be compiled again. Instead, the already compiled code is reused.

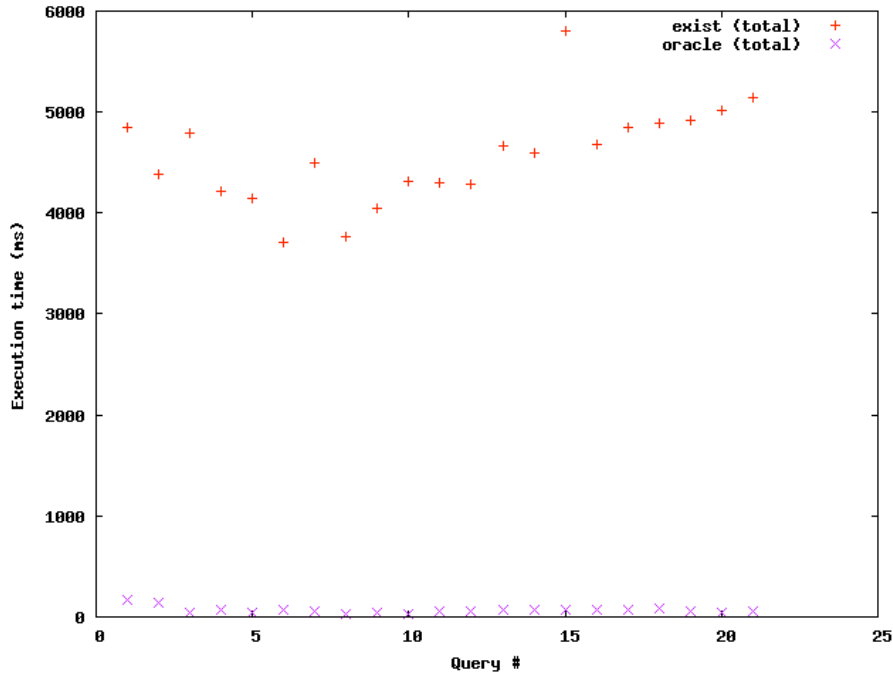


Figure 20: Projection on variables experiments. Performance compared to the original relational implementation was disappointing.

5.2.2 Results for Projection on Variables

Figure 20 shows our XQuery backend’s results for the projection on variables operation (section 2.2.2). These results were particularly disappointing. Section 4.2 discusses the relevant code.

When projected columns are removed from the result, some rows will have the same values for all remaining columns - probabilities for these rows must be merged. Identifying and merging these duplicate rows using XQuery is expensive: each *pair of rows* must be checked for equality, leading to a complexity of $O(n^2)$ for n rows. It’s certainly possible to optimize this

case in XQuery. We could manually construct indexes on discrete values for each column in such a way that each row has a unique integer index based on the values of its random variables. When one column is removed and indexes recalculated, rows with identical values would also have identical indexes. Exist is much more likely to be able to optimize the comparison of a single integer, improving performance and complexity. This has not yet been implemented.

5.2.3 Memory Usage

While experiments focused on execution time rather than memory usage, it was apparent that our XQuery implementation handled memory more efficiently than the existing Oracle implementation. Oracle ran out of memory on some experiments with large joins and products; ExistDB was able to complete the same experiments. This problem with the Oracle implementation existed in previous implementations, and our focus is on the new ExistDB implementation, so we made no attempt to diagnose or repair the Oracle backend's memory problems.

6 Conclusion and Future Work

We have presented a native XML database backend to an XML-centric semistructured probabilistic database. Our original goal of improved execution time over the previous Oracle backend was not met. We analyzed the problems behind these disappointing results and developed several ideas for relevant future work:

- Improve the execution time of all queries by recompiling and reloading our XQuery module, or by compiling only the functions required for a particular operation.
- Port our XQuery module to other XML database backends. Compare their performance.
- Improve the execution time of projection on variables, perhaps using the indexing scheme described in section .
- Implement the AND optimization discussed in section 4.1.2.
- Improve the execution time of join and product operations.

7 Bibliography

References

- [1] Existdb documentation: Xquery documentation: Xquery caching.
- [2] ANDERSON, D., REED, A., ROSSON, E., AND SIDEROPOULOS, A. Native xml support for semistructured probabalistic data management. Tech. rep., Computer Science Department, California Polytechnic State University, Dec 2007.
- [3] ANDERSON, D., REED, A., ROSSON, E., AND SIDEROPOULOS, A. Native xml support for semistructured probabalistic data management. Tech. rep., Computer Science Department, California Polytechnic State University, Oct 2007.
- [4] ANDERSON, D., REED, A., ROSSON, E., AND SIDEROPOULOS, A. Native xml support for semistructured probabalistic data management: A progress report. Tech. rep., Computer Science Department, California Polytechnic State University, Nov 2007.
- [5] BARBARÁ, D., GARCIA-MOLINA, H., AND PORTER, D. The management of probabilistic data. *IEEE Trans. on Knowledge and Data Engineering* 4 (1992), 487–502.

- [6] BOULOS, J., DALVI, N., MANDHANI, B., MATHUR, S., RE, C., AND SUCIU, D. Mystiq: A system for finding more answers by using probabilities.
- [7] CAVALLO, R., AND PITTARELLI, M. The theory of probabilistic databases. pp. 71–81.
- [8] CAVALLO, R., AND PITTARELLI, M. The theory of probabilistic databases. In *Proc. VLDB'87* (1987), pp. 71–81.
- [9] DEKHTYAR, A., GUTTI, P., AND MATHIAS, K. Structured queries for semistructured probabilistic data. *ACM TDM* (2006).
- [10] PEARL, J. *Probabilistic Reasoning in Intelligent Systems*. Morgan Kaufmann, 1988.
- [11] ZHAO, W., DEKHTYAR, A., AND GOLDSMITH, J. A framework for management of semistructured probabilistic data. *Journal of Intelligent Information Systems* 25, 3 (2004), 293–332.

A XQuery

```

1 module namespace spdb="http://www.csc.calpoly.edu/~
   erosson/spdb";
2
3 declare function spdb:sp-select-var($spos as element
   (spos), $varname as xs:string)
4 as element(spos)
5 {
6     element spos {
7         $spos/spo[table/variable/name = $varname]
8     }
9 };
10
11 (: @author Dustin Anderson :)
12 declare function spdb:sp-select-table($sp as element
   (spos),
13 $varname as xs:string, $value as xs:string, $comp as
   xs:string) as element(spos)
14 {
15 let $ret := (
16 for $spo in $sp//variable/name[text() eq $varname]/
   ancestor::spo
17 let $n := ( for $var in $spo//variable
18     for $e at $i in $var/name
19     return if ( $e eq $varname) then $i else ()
20 )
21 let $rows := (
22     if ($comp eq "=") then
23         $spo//row/val[position()=$n and (text() eq
   $value)]/parent::row
24     else if ($comp eq "!=") then
25         $spo//row/val[position()=$n and (text() ne
   $value)]/parent::row
26     else if ($comp eq "<") then
27         $spo//row/val[position()=$n and (text() lt
   $value)]/parent::row

```

```

28     else if ($comp eq "<=") then
29         $spo//row/val[position()=$n and (text() le
           $value)]/parent::row
30     else if ($comp eq ">") then
31         $spo//row/val[position()=$n and (text() gt
           $value)]/parent::row
32     else if ($comp eq ">=") then
33         $spo//row/val[position()=$n and (text() ge
           $value)]/parent::row
34     else ()
35     )
36     let $variables := $rows/parent::*/*variable
37     let $context := $rows/ancestor::*/*context
38     let $table := element table {$variables, $rows}
39     let $conditional := $rows/ancestor::*/*conditional
40     let $spo := <spo path="{ $spo/@path }">{$context,
           $table, $conditional}</spo>
41     return $spo
42     )
43     return <spos>{$ret}</spos>
44 };
45
46 (: @author Dustin Anderson :)
47 declare function spdb:sp-select-prob($sp as element(
           spos), $value as xs:string,
48 $comp as xs:string) as element(spos)
49 {
50     let $spos := (
51     for $node in $sp//spo
52         let $context := $node/descendant::context
53         let $conditional := $node/descendant::conditional
54         let $variables := $node/descendant::variable
55         let $rows := (
56             if ($comp eq "=") then
57                 $node/descendant::P[text() eq $value]/parent
                   ::row
58             else if ($comp eq "!=") then

```

```

59     $node/descendant::P[text() ne $value]/parent
        ::row
60     else if ($comp eq "<") then
61         $node/descendant::P[text() lt $value]/parent
            ::row
62     else if ($comp eq "<=") then
63         $node/descendant::P[text() le $value]/parent
            ::row
64     else if ($comp eq ">") then
65         $node/descendant::P[text() gt $value]/parent
            ::row
66     else if ($comp eq ">=") then
67         $node/descendant::P[text() ge $value]/parent
            ::row
68     else ()
69     )
70     let $table := element table {$variables, $rows}
71     let $spo := <spo path="{ $node/@path }">{ $context,
        $table, $conditional }</spo>
72     return $spo
73     )
74     return <spos> { $spos } </spos>
75 };
76
77 (: @author Dustin Anderson :)
78 declare function spdb:sp-select-conditional($sp as
        element(spos),
79 $varname as xs:string, $value as xs:string, $comp as
        xs:string) as element(spos)
80 {
81     let $thespos := (
82         for $node in $sp//spo
83         let $elem := (
84             if ($comp eq "=") then
85                 $node/descendant::conditional/child::elem[(
                    descendant::name/text() eq $varname) and (
                    descendant::val eq $value)]
86             else if ($comp eq "!=") then

```



```

87     $node/descendant::conditional/child::elem[(
      descendant::name/text() eq $varname) and (
      descendant::val ne $value)]
88   else if ($comp eq "<") then
89     $node/descendant::conditional/child::elem[(
      descendant::name/text() eq $varname) and (
      descendant::val lt $value)]
90   else if ($comp eq "<=") then
91     $node/descendant::conditional/child::elem[(
      descendant::name/text() eq $varname) and (
      descendant::val le $value)]
92   else if ($comp eq ">") then
93     $node/descendant::conditional/child::elem[(
      descendant::name/text() eq $varname) and (
      descendant::val gt $value)]
94   else if ($comp eq ">=") then
95     $node/descendant::conditional/child::elem[(
      descendant::name/text() eq $varname) and (
      descendant::val ge $value)]
96   else()
97   )
98   let $spo := $elem/ancestor::spo
99   return $spo
100  )
101  return <spos> {$thespos} </spos>
102 };
103
104 (: @author Dustin Anderson :)
105 declare function spdb:sp-select-context($sp as
      element(spos),
106 $varname as xs:string, $value as xs:string, $comp as
      xs:string) as element(spos)
107 {
108   let $thespos := (
109     for $node in $sp//spo
110     let $elem := (
111       if ($comp eq "=") then

```

```

112     $node/descendant::context/child::elem[descendant
        ::name/text() eq $varname and descendant::val
        eq $value]
113     else if ($comp eq "!=") then
114     $node/descendant::context/child::elem[descendant
        ::name/text() eq $varname and descendant::val
        ne $value]
115     else if ($comp eq "<") then
116     $node/descendant::context/child::elem[descendant
        ::name/text() eq $varname and descendant::val
        lt $value]
117     else if ($comp eq "<=") then
118     $node/descendant::context/child::elem[descendant
        ::name/text() eq $varname and descendant::val
        le $value]
119     else if ($comp eq ">") then
120     $node/descendant::context/child::elem[descendant
        ::name/text() eq $varname and descendant::val
        gt $value]
121     else if ($comp eq ">=") then
122     $node/descendant::context/child::elem[descendant
        ::name/text() eq $varname and descendant::val
        ge $value]
123     else ()
124     )
125     let $spo := $elem/ancestor::spo
126     return $spo
127     )
128
129 return <spos>{$thespos}</spos>
130 };
131
132 (: @author Dustin Anderson :)
133 declare function spdb:sp-project-context($sp as
        element(spos),
134 $varnames as xs:string*) as element(spos)
135 {
136 let $thespos := (

```

```

137   for $varname in $varnames
138     for $node in $sp//spo
139       let $conditional := $node/descendant::conditional
140       let $variables := $node/descendant::variable
141       let $rows := $node/descendant::row
142       let $table := element table {$variables, $rows}
143       let $elem := $node/descendant::context/child::
           elem[descendant::name/text() eq $varname]
144       let $context := <context> {$elem} </context>
145       let $spo := <spo path="{ $node/@path }">{$context,
           $table, $conditional}</spo>
146       return $spo
147     )
148
149   return <spos>{$thespos}</spos>
150 };
151
152 (: @author Dustin Anderson :)
153 declare function spdb:sp-project-conditional($sp as
           element(spos),
154 $varnames as xs:string*) as element(spos)
155 {
156   let $thespos := (
157     for $varname in $varnames
158       for $node in $sp//spo
159         let $context := $node/descendant::context
160         let $variables := $node/descendant::variable
161         let $rows := $node/descendant::row
162         let $table := element table {$variables, $rows}
163         let $elem := $node/descendant::conditional/child
           ::elem[descendant::name/text() eq $varname]
164         let $conditional := <conditional> {$elem} </
           conditional>
165         let $spo := <spo path="{ $node/@path }">{$context,
           $table, $conditional}</spo>
166         return $spo
167     )
168   return <spos>{$thespos}</spos>

```

```
169 };  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179 (: http://www.xqueryfunctions.com/xq/funcctx_value-  
    intersect.html :)  
180 declare namespace funcctx = "http://www.funcctx.com";  
181 declare function funcctx:value-intersect  
182   ( $arg1 as xs:anyAtomicType* ,  
183     $arg2 as xs:anyAtomicType* ) as xs:  
    anyAtomicType* {  
184  
185   distinct-values($arg1[.=$arg2])  
186 } ;  
187  
188 declare function funcctx:is-node-in-sequence-deep-  
    equal  
189   ( $node as node()? ,  
190     $seq as node()* ) as xs:boolean {  
191  
192   some $nodeInSeq in $seq satisfies deep-equal(  
     $nodeInSeq,$node)  
193 } ;  
194  
195 declare function funcctx:distinct-deep  
196   ( $nodes as node()* ) as node()* {  
197  
198   for $seq in (1 to count($nodes))  
199   return $nodes[$seq][not(funcctx:is-node-in-  
     sequence-deep-equal(  
200     .,$nodes[position() < $seq  
        ])))]
```

```
201 } ;
202
203
204 declare function spdb:QName($code as xs:string) as
      xs:QName {
205   fn:QName("http://www.csc.calpoly.edu/~erosson/spdb
      /errors", $code)
206 };
207 declare function spdb:error($code as xs:string, $msg
      as xs:string) as node() {
208   fn:error(spdb:QName($code), fn:concat("SPO error:
      ", $msg))
209 };
210 declare function spdb:error($msg as xs:string) as
      node() {
211   spdb:error("???", $msg)
212 };
213
214 declare function spdb:relation($name as xs:string)
      as element(spos) {
215   let $ret := document($name)/spos
216   return if (fn:empty($ret)) then spdb:error("
      relation_dne",fn:concat("Relation does not
      exist: ", $name)) else
217   $ret
218 };
219
220 declare function spdb:cond-normalize($cond as
      element(conditional)) as element(conditional) {
221   element conditional {
222     for $elem in $cond/element
223     order by $elem/name
224     return $elem
225   }
226 };
227
228 (:
229 Match a SPO's <cond> blocks. Order doesn't matter.
```

```
230 :)
231 declare function spdb:cond-match($cond1 as element(
      conditional), $cond2 as element(conditional)) as
      xs:boolean {
232   deep-equal(spdb:cond-normalize($cond1), spdb:cond-
      normalize($cond2))
233 };
234
235 declare function spdb:common-vars($vars1 as element(
      variable), $vars2 as element(variable)) as
      element(variable) {
236   element variable {
237     for $name in functx:value-intersect($vars1/name/
      text(), $vars2/name/text())
238     return element name {$name}
239   }
240 };
241
242 declare function spdb:common-vars-empty($elem1 as
      element(), $elem2 as element()) as xs:boolean {
243   fn:empty(spdb:common-vars($elem1//variable, $elem2
      //variable)/name)
244 };
245
246 declare function spdb:product-compatible($spo1 as
      element(spo), $spo2 as element(spo)) as xs:
      boolean {
247   if (spdb:cond-match($spo1/conditional, $spo2/
      conditional))
248   then spdb:common-vars-empty($spo1, $spo2)
249   else false()
250 };
251
252 declare function spdb:join-compatible($spo1 as
      element(spo), $spo2 as element(spo)) as xs:
      boolean {
253   if (spdb:cond-match($spo1/conditional, $spo2/
      conditional))
```

```
254     then not(spdb:common-vars-empty($spo1, $spo2))
255     else false()
256 };
257
258 declare function spdb:table-product($table1 as
      element(table), $table2 as element(table)) as
      element(table) {
259     if (not(spdb:common-vars-empty($table1, $table2)))
          then () else
260     element table {
261         element variable {$table1/variable/name, $table2
          /variable/name},
262         spdb:rows-product($table1/row, $table2/row)
263     }
264 };
265
266 declare function spdb:rows-product($rows1 as element
      (row)*, $rows2 as element(row)*) as element(row)*
      {
267     for $row1 in $rows1
268     for $row2 in $rows2
269     let $prob := element P { $row1/P/number() *
          $row2/P/number() }
270     return element row {$row1/val, $row2/val,
          $prob}
271 };
272
273 declare function spdb:context-merge($con1 as element
      (context), $con2 as element(context)) as element(
      context) {
274     element context {
275         functx:distinct-deep(($con1/elem, $con2/elem))
276     }
277 };
278
279 declare function spdb:binary-spo-wrap($spo1 as
      element(spo), $spo2 as element(spo), $table as
      element(table)) as element(spo) {
```

```
280     element spo {
281         spdb:context-merge($spo1/context, $spo2/context)
282         ,
283         $table,
284     }
285 };
286
287 declare function spdb:product($spo1 as element(spo),
288     $spo2 as element(spo)) as element(spo)? {
289     if (not(spdb:product-compatible($spo1, $spo2)))
290         then () else
291         spdb:binary-spo-wrap($spo1, $spo2, spdb:table-
292             product($spo1/table, $spo2/table))
293 };
294
295 declare function spdb:sp-product($spos1 as element(
296     spos), $spos2 as element(spos)) as element(spos)
297     {
298     element spos {
299         for $spo1 in $spos1/spo
300         for $spo2 in $spos2/spo
301         return spdb:product($spo1, $spo2)
302     }
303 };
304
305 declare function spdb:table-column-num($vars as
306     element(name)*, $name as xs:string) as xs:integer
307     ? {
308     spdb:table-column-nums($vars, $name)
309 };
310
311 declare function spdb:table-column-nums($vars as
312     element(name)*, $names as xs:string*) as xs:
313     integer* {
314     for $var at $index in $vars
315     where $var/text() = $names
316     return $index
```



```
308 };
309
310 declare function spdb:table-column-vals-by-nums(
    $rows as element(row)*, $num as xs:integer) as
    element(val)* {
311     $rows/val[position() eq $num]
312 };
313
314 declare function spdb:row-column-vals-by-nums($row
    as element(row), $nums as xs:integer*) as element
    (val)* {
315     $row/val[position() = $nums]
316 };
317
318 declare function spdb:table-column-vals($table as
    element(table), $name as xs:string) as element(
    val)* {
319     let $num := spdb:table-column-num($table/variable/
        name, $name)
320     return if (fn:empty($num)) then $table/row/val
        else
321         spdb:table-column-vals-by-nums($table/row, $num)
322 };
323
324 declare function spdb:row-column-vals($vars as
    element(name)*, $row as element(row), $names as
    xs:string*) as element(val)* {
325     spdb:row-column-vals-by-nums($row, spdb:table-
        column-nums($vars, $names))
326 };
327
328 declare function spdb:conditionalize-table($table as
    element(table), $name as xs:string, $value as xs
    :string) as element(table) {
329     let $num := spdb:table-column-num($table/variable/
        name, $name)
330     (: If var isn't in table, return the original
        table.
```

```

331   This isn't silent failure, but an irrelevant
      conditionalization. :)
332   return if (fn:empty($num)) then $table else
333   let $newrows := $table/row[exists(val[$num][text()
      eq $value])]
334   let $oldprob := sum(for $i in $table/row/P return
      xs:decimal($i))
335   let $newprob := sum(for $i in $newrows/P return xs
      :decimal($i))
336   let $probmult := if ($newprob gt 0) then $oldprob
      div $newprob else 1
337   return element table {
338     element variable {$table/variable/name[position
      () ne $num]},
339     for $row in $newrows
340       let $p := xs:decimal($row/P) * $probmult
341       return element row {$row/val[position() ne
      $num], element P{$p}}
342   }
343 };
344
345 declare function spdb:conditionalize($spo as element
      (spo), $name as xs:string, $value as xs:string)
      as element(spo) {
346   element spo {$spo/@*,
347     $spo/context,
348     spdb:conditionalize-table($spo/table, $name,
      $value),
349     element conditional {
350       $spo/conditional/*,
351       element elem {element name{$name}, element val
      {$value}}
352     }
353   }
354 };
355
356 declare function spdb:sp-conditionalize($spos as
      element(spos), $name as xs:string, $value as xs:

```

```

        string) as element(spos) {
357   element spos {$spos/@*,
358     for $spo in $spos/spo
359       return spdb:conditionalize($spo, $name, $value
        )
360   }
361 };
362 declare function spdb:leftjoin-conditionalize-table(
    $table as element(table), $vars as element(name)
    *, $vals as element(val)*) as element(table) {
363   if (fn:empty($vars)) then $table else
364   let $var := $vars[1]
365   let $val := $vals[1]
366   return spdb:leftjoin-conditionalize-table(element
    table {$table/@*,
367     element variable {$table/variable/name[. ne $var
    ]}),
368     spdb:conditionalize-table($table, $var, $val)/
    row
369   }, $vars[position() ne 1], $vals[position() ne 1])
370 };
371
372 declare function spdb:table-leftjoin($table1 as
    element(table), $table2 as element(table)) as
    element(table)? {
373   let $common := spdb:common-vars($table1/variable,
    $table2/variable)/name
374   return if (fn:empty($common)) then () else
375   element table {$table1/@*,
376     element variable {$table1/variable/name, $table2
    /variable/name[not(. = $common)]},
377     for $row1 in $table1/row
378       let $values := spdb:row-column-vals($table1/
    variable/name, $row1, $common/text())
379     let $condtable2 := spdb:leftjoin-
    conditionalize-table($table2, $common,
    $values)

```

```
380         return spdb:rows-product($row1, $condtable2/
           row)
381     }
382 };
383
384 declare function spdb:leftjoin($spo1 as element(spo)
           , $spo2 as element(spo)) as element(spo)? {
385     if (not(spdb:join-compatible($spo1, $spo2))) then
           () else
386     spdb:binary-spo-wrap($spo1, $spo2, spdb:table-
           leftjoin($spo1/table, $spo2/table))
387 };
388
389 declare function spdb:sp-leftjoin($spos1 as element(
           spos), $spos2 as element(spos)) as element(spos)
           {
390     element spos {
391         for $spo1 in $spos1/spo
392         for $spo2 in $spos2/spo
393         return spdb:leftjoin($spo1, $spo2)
394     }
395 };
396
397 declare function spdb:sp-rightjoin($spos1 as element
           (spos), $spos2 as element(spos)) as element(spos)
           {
398     spdb:sp-leftjoin($spos2, $spos1)
399 };
400
401 declare function spdb:table-leftmix($table1 as
           element(table), $table2 as element(table)) as
           element(table) {
402     let $common := spdb:common-vars($table1/variable,
           $table2/variable)/name
403     return if (fn:empty($common)) then spdb:table-
           product($table1, $table2) else spdb:table-
           1         $table1, $table2)
404 };
```

```
405
406 declare function spdb:leftmix($spo1 as element(spo),
    $spo2 as element(spo)) as element(spo)? {
407   if (not(spdb:cond-match($spo1/conditional, $spo2/
    conditional))) then () else
408   spdb:binary-spo-wrap($spo1, $spo2, spdb:table-
    leftmix($spo1/table, $spo2/table))
409 };
410
411 declare function spdb:sp-leftmix($spos1 as element(
    spos), $spos2 as element(spos)) as element(spos)
    {
412   element spos {
413     for $spo1 in $spos1/spo
414       for $spo2 in $spos2/spo
415         return spdb:leftmix($spo1, $spo2)
416   }
417 };
418
419 declare function spdb:sp-rightmix($spos1 as element(
    spos), $spos2 as element(spos)) as element(spos)
    {
420   spdb:sp-leftmix($spos2, $spos1)
421 };
422
423 declare function spdb:sp-where-or($spos1 as element(
    spos), $spos2 as element(spos)) as element(spos)
    {
424   element spos {$spos1/@*, functx:distinct-deep((
    $spos1/spo, $spos2/spo)) }
425 };
426
427 (: Existdb-specific: http://exist.sourceforge.net/
    update\_ext.html :)
428 declare function spdb:sp-update-insert($dest as
    element(spos), $data as element(spos)) as node()*
    {
429   update insert $data/spo into $dest
```

```
430 };
431
432 declare function spdb:sp-update-delete($dest as xs:
      string, $matches as element(spos)) as node()* {
433   for $match in $matches/spo
434     for $spo in spdb:relation($dest)/spo
435       where deep-equal($match, $spo)
436         return update delete $spo
437 };
438
439 declare function spdb:rows-project-var($rawrows as
      element(row)*, $visiblenums as xs:integer*) as
      element(row)* {
440   let $rowsviscols := for $row in $rawrows
441     return element row {$row/@*,
442       $row/val[position() = $visiblenums],
443       $row/P
444   }
445   return for $row1 at $i in $rowsviscols
446     (: where no previous rows have identical values
447       :)
447     where fn:empty($rowsviscols[position() lt $i][
448       not(val != $row1/val)])
448     return element row {$row1/@*,
449       $row1/val,
450       element P {$row1/P/@*,
451         (: sum all rows with identical values. skip
452           previous rows: we checked those already
453           :)
452         sum(for $p in $rowsviscols[position() ge $i
453           ][not(val != $row1/val)]/P
454           return xs:decimal($p))
454     }
455   }
456 };
457
458 declare function spdb:table-project-var($table as
      element(table), $visiblevars as xs:string*) as
```

```

    element(table) {
459   let $visiblenums := spdb:table-column-nums($table/
        variable/name, $visiblevars)
460   return element table {$table/@*,
461     element variable {$table/variable/@*,
462       $table/variable/name[position() = $visiblenums
        ]
463     },
464     spdb:rows-project-var($table/row, $visiblenums)
465   }
466 };
467
468 declare function spdb:project-var($spo as element(
    spo), $visiblevars as xs:string*) as element(spo)
    {
469   element spo {$spo/@*,
470     $spo/context,
471     spdb:table-project-var($spo/table, $visiblevars)
472     ,
473     $spo/conditional
474   };
475
476 declare function spdb:sp-project-var($spos as
    element(spos), $visiblevars as xs:string*) as
    element(spos) {
477   element spos {$spos/@*,
478     for $spo in $spos/spo
479       return spdb:project-var($spo, $visiblevars)
480   }
481 };
482
483 declare function spdb:rows-project-var-2($rawrows as
    element(row)*, $visiblenums as xs:integer*) as
    element(row)* {
484   let $rowsviscols := for $row in $rawrows
485     return element row {$row/@*,
486       $row/val[position() = $visiblenums],

```

```

487     $row/P
488   }
489   return for $row1 at $i in $rowsviscols
490     (: where no previous rows have identical values
491       :)
491     where fn:empty($rowsviscols[position() lt $i][
492       deep-equal(val, $row1/val)])
492     return element row {$row1/@*,
493       $row1/val,
494       element P {$row1/P/@*,
495         (: sum all rows with identical values. skip
496           previous rows: we checked those already
497           :)
496         sum(for $p in $rowsviscols[position() ge $i
497           ][deep-equal(val, $row1/val)]/P
498           return xs:decimal($p))
498       }
499     }
500 };
501
502 declare function spdb:table-project-var-2($table as
503   element(table), $visiblevars as xs:string*) as
504   element(table) {
503   let $visiblenums := spdb:table-column-nums($table/
504     variable/name, $visiblevars)
504   return element table {$table/@*,
505     element variable {$table/variable/@*,
506       $table/variable/name[position() = $visiblenums
507       ]
507     },
508     spdb:rows-project-var-2($table/row, $visiblenums
509     )
509   }
510 };
511
512 declare function spdb:project-var-2($spo as element(
513   spo), $visiblevars as xs:string*) as element(spo)
514   {

```



```

513   element spo {$spo/@*,
514     $spo/context,
515     spdb:table-project-var-2($spo/table,
516       $visiblevars),
517   }
518 };
519
520 declare function spdb:sp-project-var-2($spos as
521   element(spos), $visiblevars as xs:string*) as
522   element(spos) {
523     element spos {$spos/@*,
524       for $spo in $spos/spo
525         return spdb:project-var-2($spo, $visiblevars)
526     }
527 };
528
529 (: @author Dustin Anderson :)
530 declare function spdb:sp-project-var-3($sp as
531   element(spos), $visiblevars as xs:string*) as
532   element(spos)
533 {
534   let $hiddenvars := $sp//variable/name[. !=
535     $visiblevars]
536   let $thespos := (
537     (:Only iterate over spos that have CS113
538       in them as <name> elements:)
539     for $varname in $hiddenvars
540     for $node in $sp/descendant::variable/name[text() eq
541       $varname]/ancestor::spo
542     (: Get the position of 'CS113' in the <
543       names> :)
544     let $n := ( for $w in $node//variable
545       for $e at $i in $w/child::name
546       return
547         if ( $e eq $varname) then $i
548         else ()

```

```

541         )
542
543         (:Remove $varnames from the <variable>
           element:)
544     let $variable := (
545         for $v in $node//table/child::variable
546         let $name := $v/child::name[position() !=
           $n]
547         return element variable{$name}
548     )
549
550     (:Remove elements in table that match
       with position of $varnames:)
551     let $rowswithP := (
552         for $r in $node//table/row
553         let $vals := $r/child::val[position() !=
           $n]
554         return element row{$vals, $r/child::P
           }
555     )
556
557     (:Remove elements in table that match
       with position of $varnames:)
558     let $rowswithoutP := (
559         for $r in $node//table/row
560         let $vals := $r/child::val[position() !=
           $n]
561         return element row{$vals}
562     )
563
564     let $table := (
565     for $xx at $ii in $rowswithoutP
566     for $yy at $jj in $rowswithoutP
567     return if (($xx eq $yy) and ($jj > $ii))
           then
568     element row {$xx/val, <P>{($rowswithP[
           position()=$ii]/P + $rowswithP[
           position()=$jj]/P)}</P>}

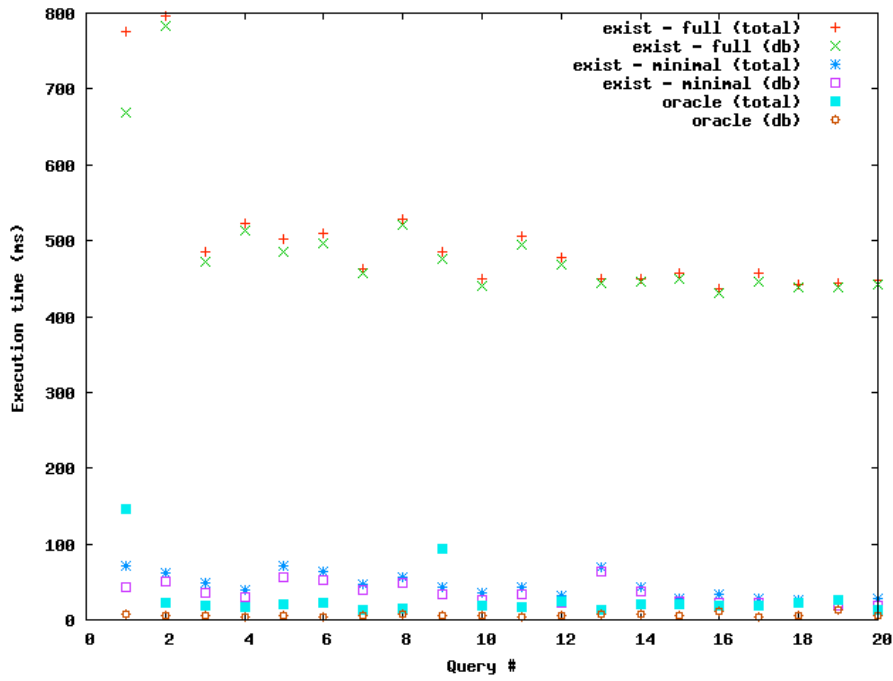
```

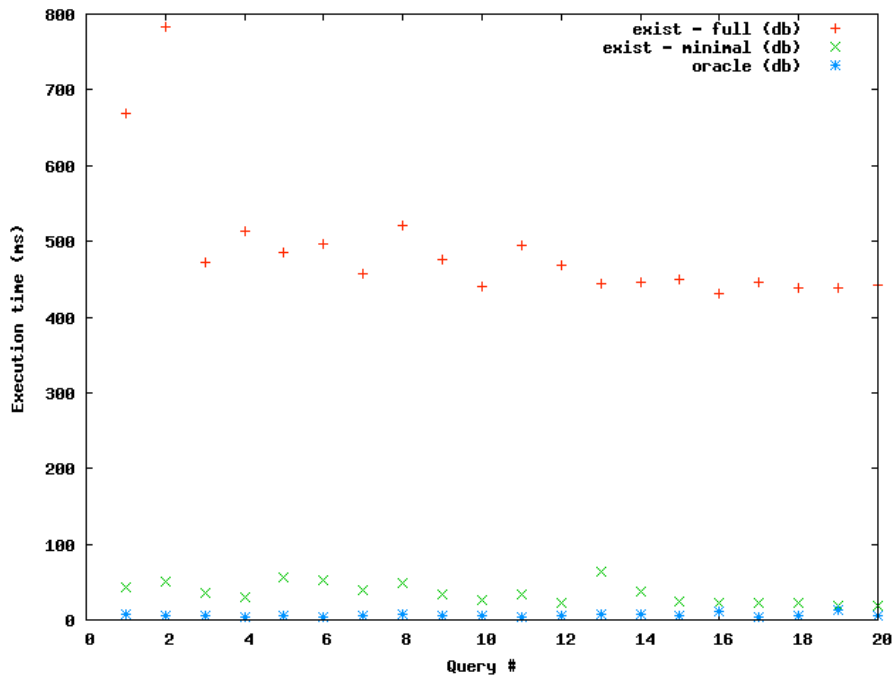
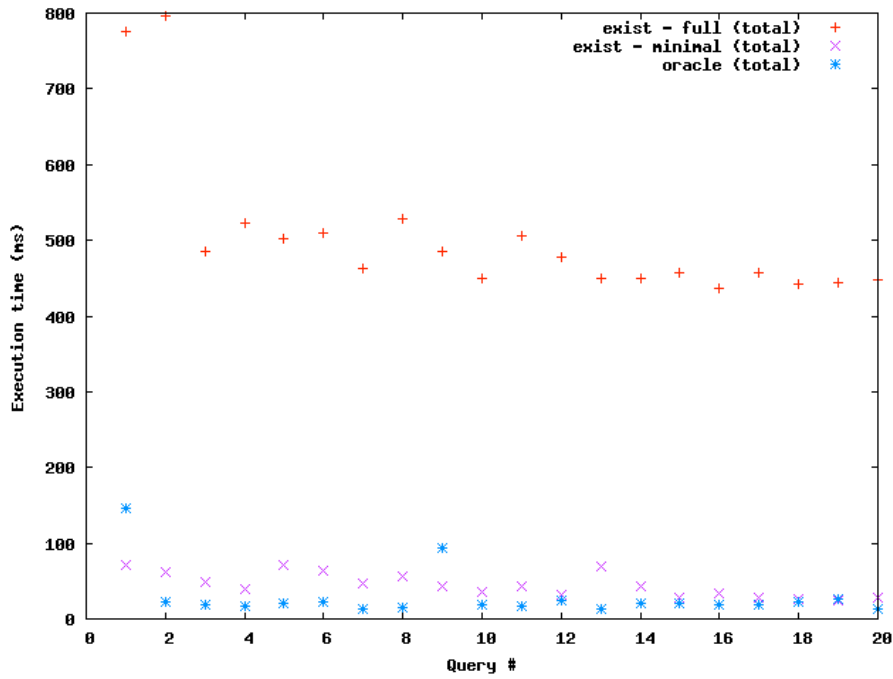
```
569         else
570             ()
571     )
572
573     let $context := $node//context
574     let $conditional := $node//conditional
575
576     return <spo path="{ $node/@path }">{ $context, <
        table>{ $variable, $table}</table>,
        $conditional}</spo>
577 )
578
579 return <spos>{ $thespos}</spos>
580 };
```

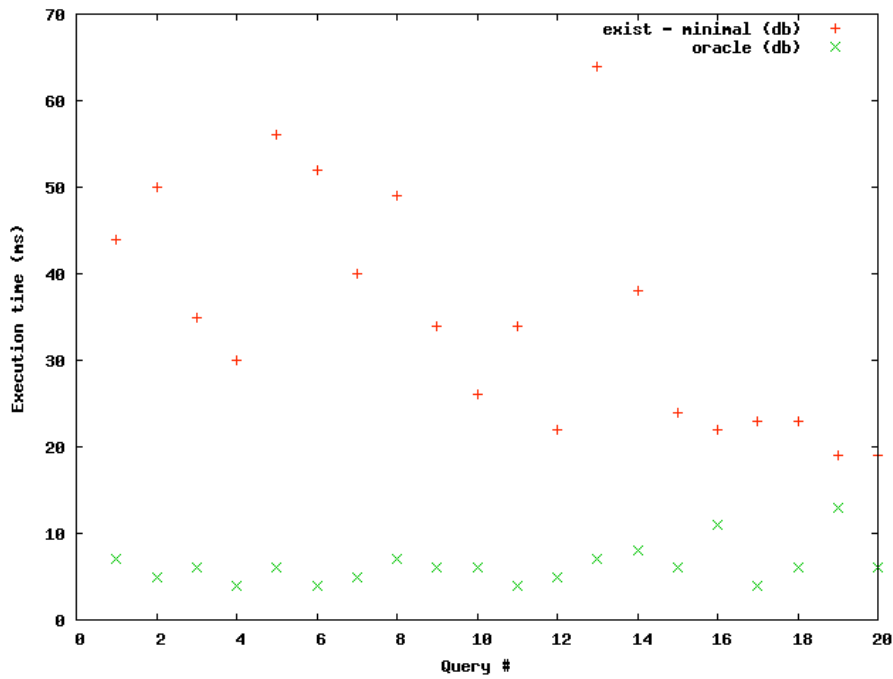
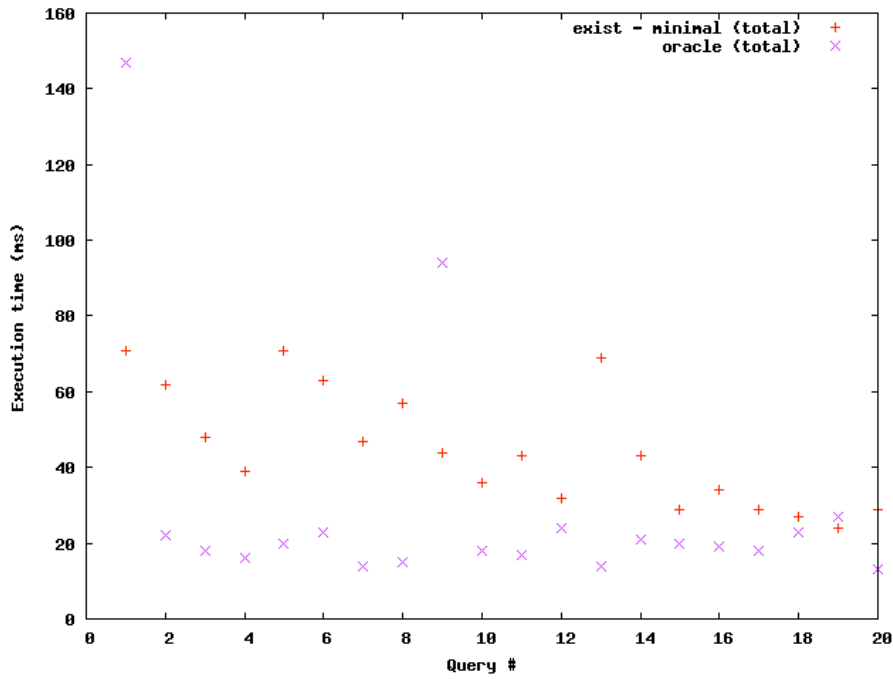
B Experimental Results

B.1 Baseline experiments

Comparison







Oracle Data

```
1 QUERY XML select * from empty: (time: 147 db: 7)
2 QUERY XML select * from empty: (time: 22 db: 5)
3 QUERY XML select * from empty: (time: 18 db: 6)
4 QUERY XML select * from empty: (time: 16 db: 4)
5 QUERY XML select * from empty: (time: 20 db: 6)
6 QUERY XML select * from empty: (time: 23 db: 4)
7 QUERY XML select * from empty: (time: 14 db: 5)
8 QUERY XML select * from empty: (time: 15 db: 7)
9 QUERY XML select * from empty: (time: 94 db: 6)
10 QUERY XML select * from empty: (time: 18 db: 6)
11 QUERY XML select * from empty: (time: 17 db: 4)
12 QUERY XML select * from empty: (time: 24 db: 5)
13 QUERY XML select * from empty: (time: 14 db: 7)
14 QUERY XML select * from empty: (time: 21 db: 8)
15 QUERY XML select * from empty: (time: 20 db: 6)
16 QUERY XML select * from empty: (time: 19 db: 11)
17 QUERY XML select * from empty: (time: 18 db: 4)
18 QUERY XML select * from empty: (time: 23 db: 6)
19 QUERY XML select * from empty: (time: 27 db: 13)
20 QUERY XML select * from empty: (time: 13 db: 6)
```

Exist (minimal) Data

```
1 QUERY XML select * from empty: (time: 71 db: 44)
2 QUERY XML select * from empty: (time: 62 db: 50)
3 QUERY XML select * from empty: (time: 48 db: 35)
4 QUERY XML select * from empty: (time: 39 db: 30)
5 QUERY XML select * from empty: (time: 71 db: 56)
6 QUERY XML select * from empty: (time: 63 db: 52)
7 QUERY XML select * from empty: (time: 47 db: 40)
8 QUERY XML select * from empty: (time: 57 db: 49)
9 QUERY XML select * from empty: (time: 44 db: 34)
10 QUERY XML select * from empty: (time: 36 db: 26)
11 QUERY XML select * from empty: (time: 43 db: 34)
12 QUERY XML select * from empty: (time: 32 db: 22)
13 QUERY XML select * from empty: (time: 69 db: 64)
14 QUERY XML select * from empty: (time: 43 db: 38)
15 QUERY XML select * from empty: (time: 29 db: 24)
16 QUERY XML select * from empty: (time: 34 db: 22)
```

17 QUERY XML select * from empty: (time: 29 db: 23)
18 QUERY XML select * from empty: (time: 27 db: 23)
19 QUERY XML select * from empty: (time: 24 db: 19)
20 QUERY XML select * from empty: (time: 29 db: 19)

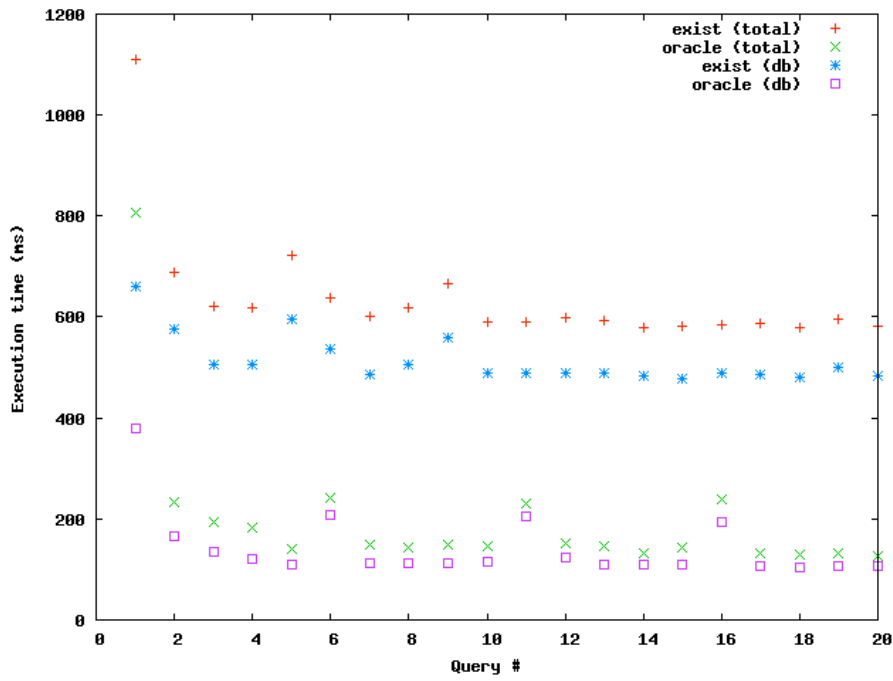
Exist (full) Data

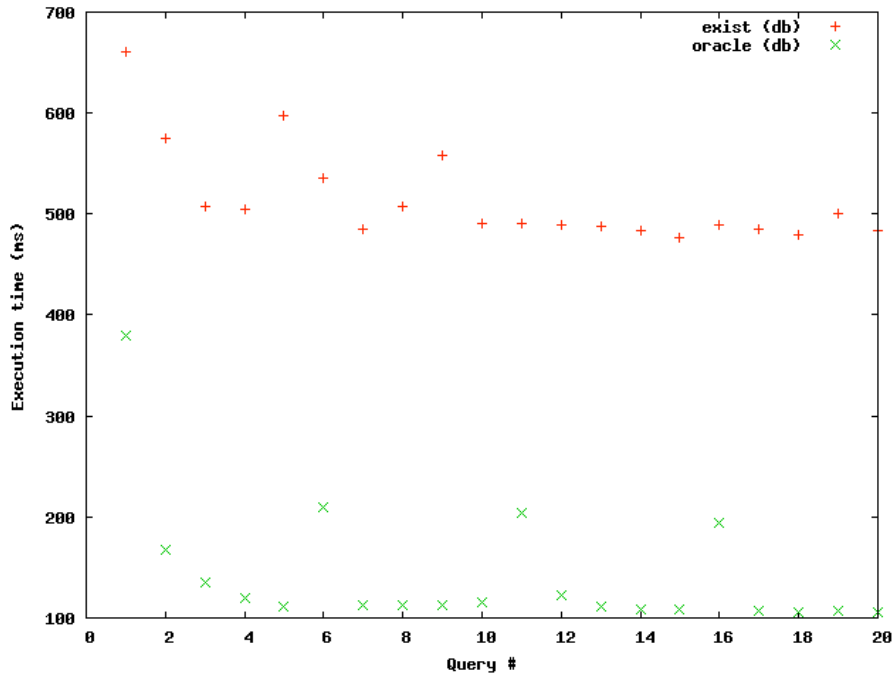
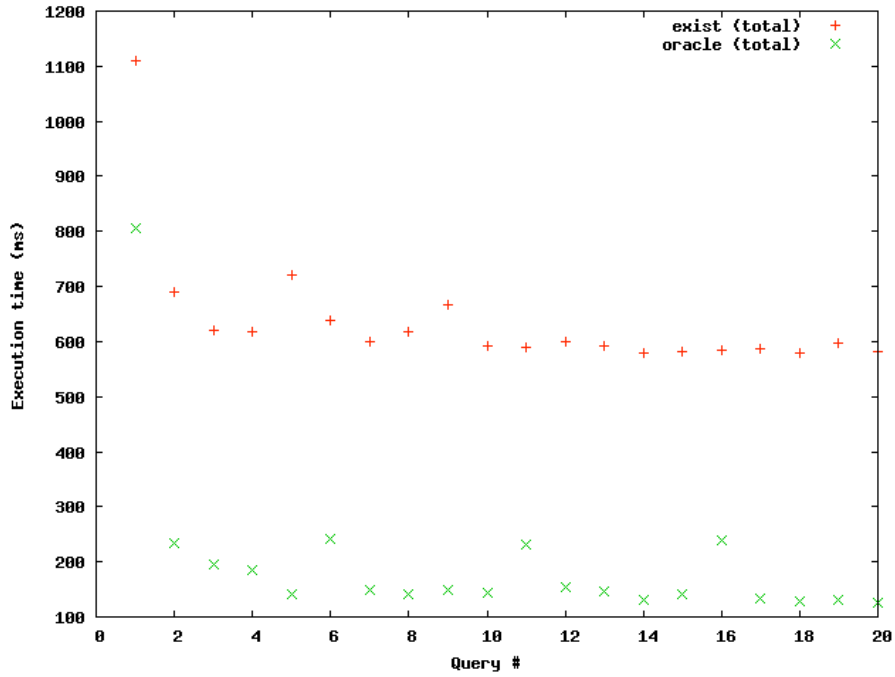
1 QUERY XML select * from empty: (time: 776 db: 668)
2 QUERY XML select * from empty: (time: 796 db: 784)
3 QUERY XML select * from empty: (time: 486 db: 472)
4 QUERY XML select * from empty: (time: 523 db: 513)
5 QUERY XML select * from empty: (time: 502 db: 486)
6 QUERY XML select * from empty: (time: 509 db: 497)
7 QUERY XML select * from empty: (time: 463 db: 458)
8 QUERY XML select * from empty: (time: 528 db: 520)
9 QUERY XML select * from empty: (time: 486 db: 476)
10 QUERY XML select * from empty: (time: 450 db: 440)
11 QUERY XML select * from empty: (time: 506 db: 494)
12 QUERY XML select * from empty: (time: 477 db: 468)
13 QUERY XML select * from empty: (time: 450 db: 444)
14 QUERY XML select * from empty: (time: 449 db: 445)
15 QUERY XML select * from empty: (time: 457 db: 450)
16 QUERY XML select * from empty: (time: 436 db: 430)
17 QUERY XML select * from empty: (time: 457 db: 445)
18 QUERY XML select * from empty: (time: 443 db: 438)
19 QUERY XML select * from empty: (time: 444 db: 439)
20 QUERY XML select * from empty: (time: 447 db: 442)

B.2 2 variables, 100 SPOs

B.2.1 Simple Selection

Comparison





Oracle Data

```
1 QUERY XML select * from First: (time: 806 db: 379)
2 QUERY XML select * from First: (time: 234 db: 167)
3 QUERY XML select * from First: (time: 195 db: 135)
4 QUERY XML select * from First: (time: 184 db: 120)
5 QUERY XML select * from First: (time: 141 db: 111)
6 QUERY XML select * from Second: (time: 241 db: 209)
7 QUERY XML select * from Second: (time: 148 db: 113)
8 QUERY XML select * from Second: (time: 142 db: 113)
9 QUERY XML select * from Second: (time: 148 db: 113)
10 QUERY XML select * from Second: (time: 145 db: 116)
11 QUERY XML select * from Third: (time: 231 db: 204)
12 QUERY XML select * from Third: (time: 153 db: 123)
13 QUERY XML select * from Third: (time: 146 db: 111)
14 QUERY XML select * from Third: (time: 131 db: 109)
15 QUERY XML select * from Third: (time: 142 db: 109)
16 QUERY XML select * from Fourth: (time: 239 db: 194)
17 QUERY XML select * from Fourth: (time: 133 db: 107)
18 QUERY XML select * from Fourth: (time: 129 db: 105)
19 QUERY XML select * from Fourth: (time: 131 db: 107)
20 QUERY XML select * from Fourth: (time: 127 db: 106)
```

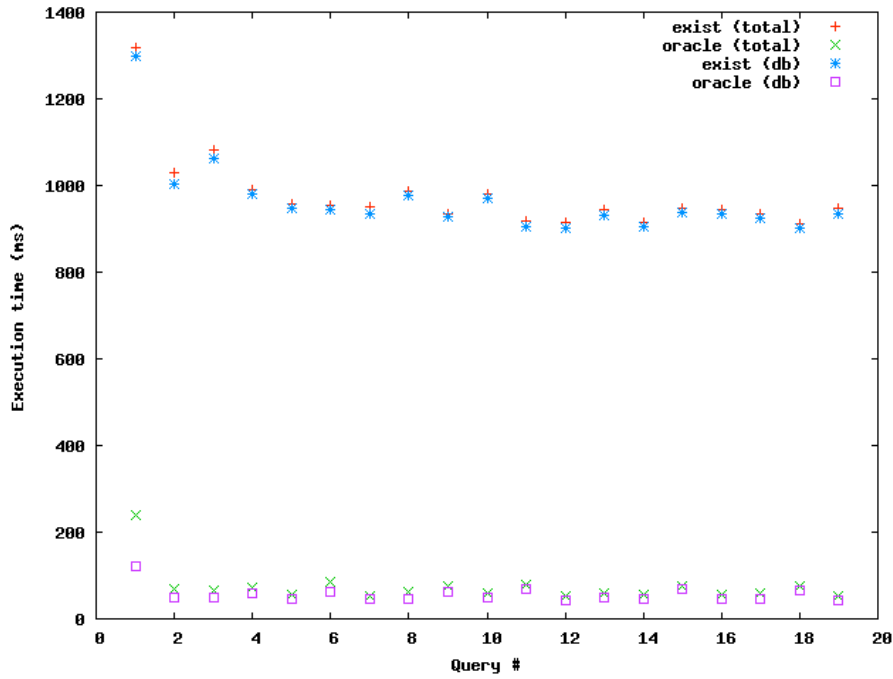
Exist Data

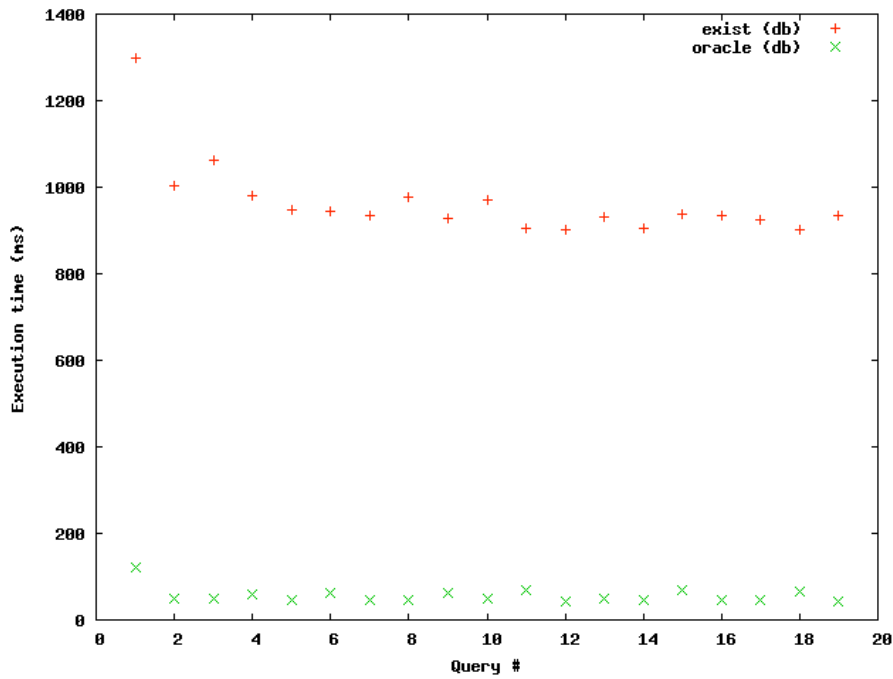
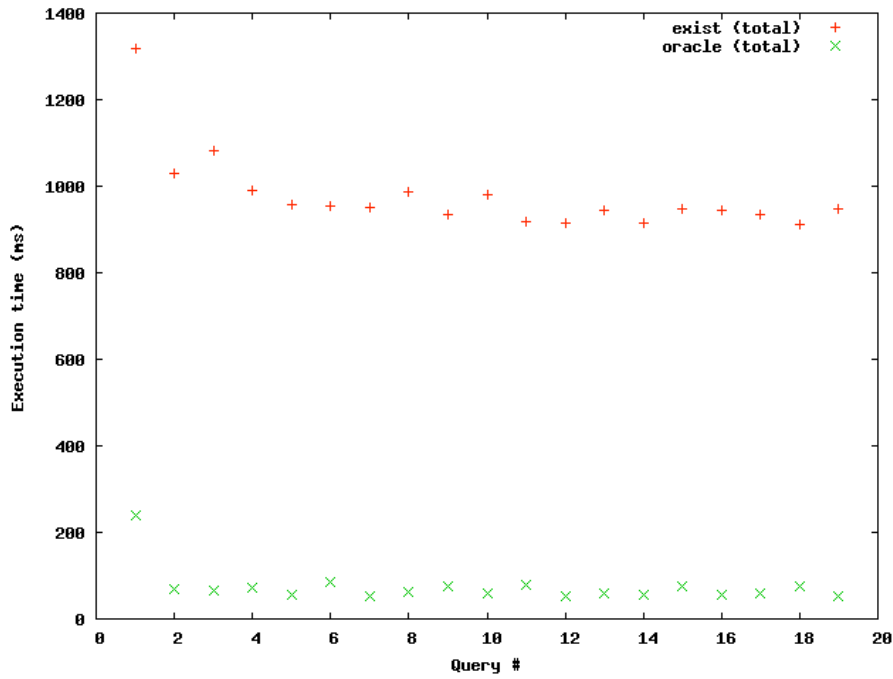
```
1 QUERY XML select * from First: (time: 1109 db: 660)
2 QUERY XML select * from First: (time: 689 db: 575)
3 QUERY XML select * from First: (time: 620 db: 507)
4 QUERY XML select * from First: (time: 617 db: 505)
5 QUERY XML select * from First: (time: 722 db: 597)
6 QUERY XML select * from Second: (time: 639 db: 536)
7 QUERY XML select * from Second: (time: 601 db: 485)
8 QUERY XML select * from Second: (time: 618 db: 507)
9 QUERY XML select * from Second: (time: 667 db: 558)
10 QUERY XML select * from Second: (time: 591 db: 490)
11 QUERY XML select * from Third: (time: 589 db: 490)
12 QUERY XML select * from Third: (time: 599 db: 489)
13 QUERY XML select * from Third: (time: 593 db: 488)
14 QUERY XML select * from Third: (time: 580 db: 483)
15 QUERY XML select * from Third: (time: 581 db: 477)
16 QUERY XML select * from Fourth: (time: 584 db: 489)
```

- 17 QUERY XML select * from Fourth: (time: 588 db: 485)
- 18 QUERY XML select * from Fourth: (time: 579 db: 480)
- 19 QUERY XML select * from Fourth: (time: 597 db: 500)
- 20 QUERY XML select * from Fourth: (time: 581 db: 483)

B.2.2 Select on Context

Comparison





Oracle Data

```
1 QUERY XML select * from First where cnt.college =XX:
    (time: 238 db: 120)
2 QUERY XML select * from First where cnt.comments=00:
    (time: 69 db: 50)
3 QUERY XML select * from First where cnt.year=1976: (
    time: 64 db: 49)
4 QUERY XML select * from First where cnt.year=1995: (
    time: 71 db: 60)
5 QUERY XML select * from First where cnt.comments=GG:
    (time: 55 db: 47)
6 QUERY XML select * from Second where cnt.semester =
    PP: (time: 86 db: 63)
7 QUERY XML select * from Second where cnt.major=XX: (
    time: 54 db: 45)
8 QUERY XML select * from Second where cnt.major=II: (
    time: 62 db: 46)
9 QUERY XML select * from Second where cnt.year=1999:
    (time: 74 db: 61)
10 QUERY XML select * from Second where cnt.semester=FF
    : (time: 58 db: 48)
11 QUERY XML select * from Third where cnt.major=MM: (
    time: 80 db: 70)
12 QUERY XML select * from Third where cnt.instructor=
    HH: (time: 53 db: 42)
13 QUERY XML select * from Third where cnt.semester=TT:
    (time: 60 db: 50)
14 QUERY XML select * from Third where cnt.comments=PP:
    (time: 56 db: 45)
15 QUERY XML select * from Fourth where cnt.comments=WW
    : (time: 77 db: 68)
16 QUERY XML select * from Fourth where cnt.year=1981:
    (time: 56 db: 45)
17 QUERY XML select * from Fourth where cnt.major=DD: (
    time: 58 db: 45)
18 QUERY XML select * from Fourth where cnt.major=WW: (
    time: 76 db: 64)
19 QUERY XML select * from Fourth where cnt.semester=GG
    : (time: 54 db: 44)
```

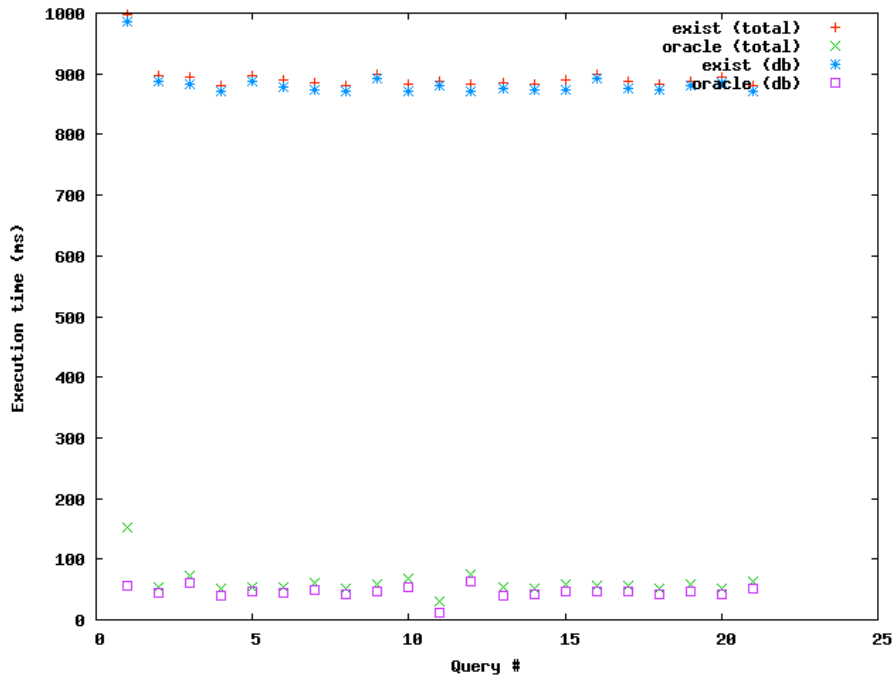
Exist Data

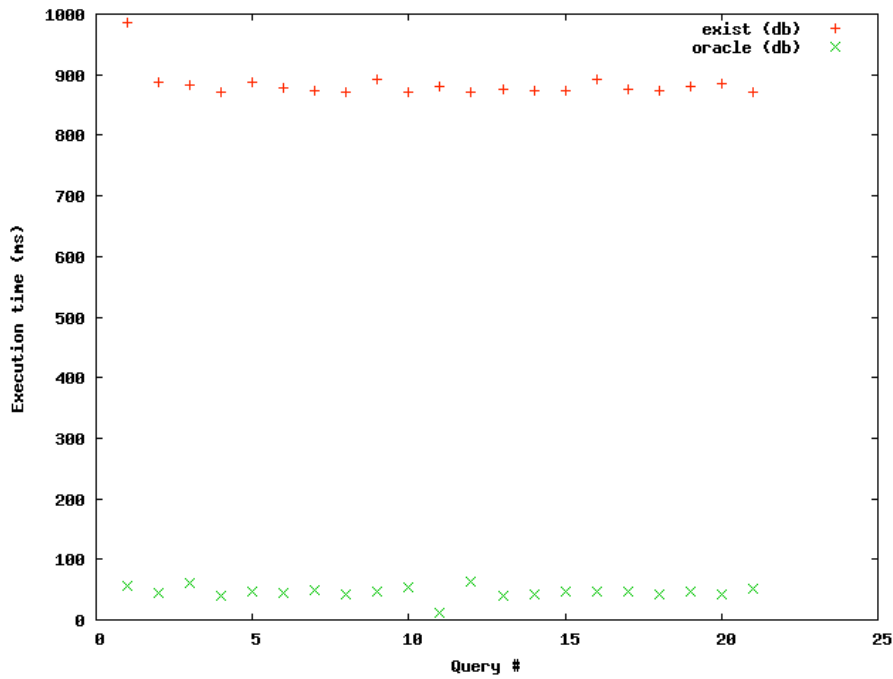
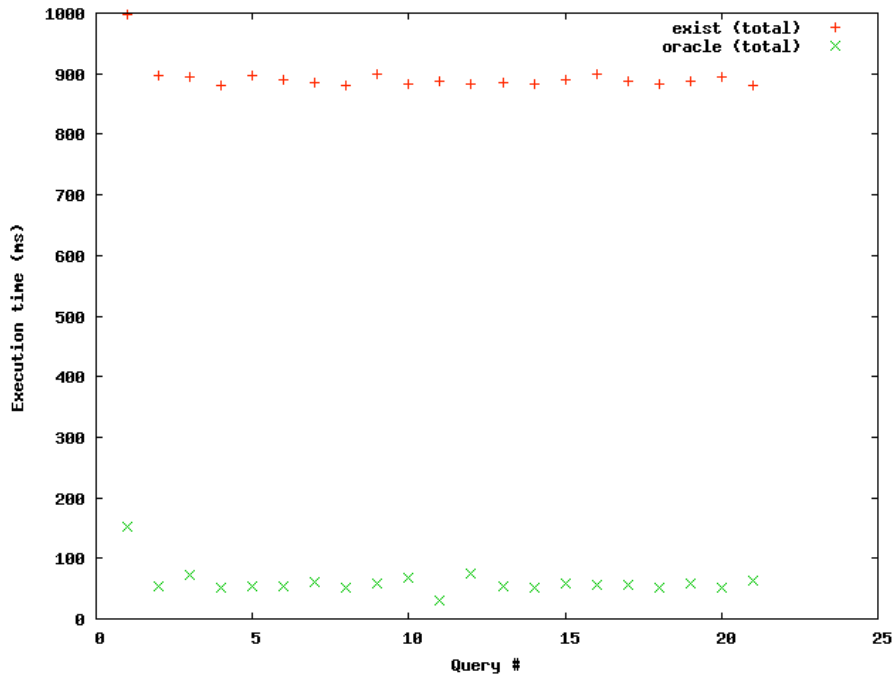
- 1 QUERY XML select * from First where cnt.college =XX:
(time: 1318 db: 1298)
- 2 QUERY XML select * from First where cnt.comments=00:
(time: 1029 db: 1004)
- 3 QUERY XML select * from First where cnt.year=1976: (
time: 1083 db: 1061)
- 4 QUERY XML select * from First where cnt.year=1995: (
time: 989 db: 980)
- 5 QUERY XML select * from First where cnt.comments=GG:
(time: 958 db: 948)
- 6 QUERY XML select * from Second where cnt.semester =
PP: (time: 955 db: 943)
- 7 QUERY XML select * from Second where cnt.major=XX: (
time: 950 db: 936)
- 8 QUERY XML select * from Second where cnt.major=II: (
time: 986 db: 977)
- 9 QUERY XML select * from Second where cnt.year=1999:
(time: 936 db: 928)
- 10 QUERY XML select * from Second where cnt.semester=FF
: (time: 979 db: 971)
- 11 QUERY XML select * from Third where cnt.major=MM: (
time: 918 db: 906)
- 12 QUERY XML select * from Third where cnt.instructor=
HH: (time: 916 db: 902)
- 13 QUERY XML select * from Third where cnt.semester=TT:
(time: 944 db: 931)
- 14 QUERY XML select * from Third where cnt.comments=PP:
(time: 914 db: 904)
- 15 QUERY XML select * from Fourth where cnt.comments=WW
: (time: 947 db: 937)
- 16 QUERY XML select * from Fourth where cnt.year=1981:
(time: 943 db: 934)
- 17 QUERY XML select * from Fourth where cnt.major=DD: (
time: 934 db: 926)
- 18 QUERY XML select * from Fourth where cnt.major=WW: (
time: 913 db: 902)

19 QUERY XML select * from Fourth where cnt.semester=GG
: (time: 947 db: 935)

B.2.3 Select on Conditional

Comparison





Oracle Data

```
1 QUERY XML select * from First where cnd.CS120=B: (
    time: 153 db: 57)
2 QUERY XML select * from First where cnd.CS101=B: (
    time: 53 db: 44)
3 QUERY XML select * from First where cnd.CS106=A: (
    time: 73 db: 60)
4 QUERY XML select * from First where cnd.CS143=B: (
    time: 52 db: 40)
5 QUERY XML select * from First where cnd.CS103=B: (
    time: 55 db: 47)
6 QUERY XML select * from Second where cnd.CS143=B: (
    time: 53 db: 44)
7 QUERY XML select * from Second where cnd.CS146=A: (
    time: 62 db: 50)
8 QUERY XML select * from Second where cnd.CS135=A: (
    time: 51 db: 43)
9 QUERY XML select * from Second where cnd.CS102=A: (
    time: 59 db: 47)
10 QUERY XML select * from Second where cnd.CS127=B: (
    time: 67 db: 55)
11 QUERY XML select * from Second where cnd.CS144=A: (
    time: 30 db: 12)
12 QUERY XML select * from Third where cnd.CS144=A: (
    time: 75 db: 63)
13 QUERY XML select * from Third where cnd.CS138=A: (
    time: 54 db: 40)
14 QUERY XML select * from Third where cnd.CS112=B: (
    time: 51 db: 42)
15 QUERY XML select * from Third where cnd.CS140=A: (
    time: 59 db: 48)
16 QUERY XML select * from Third where cnd.CS122=A: (
    time: 57 db: 46)
17 QUERY XML select * from Fourth where cnd.CS117=B: (
    time: 57 db: 46)
18 QUERY XML select * from Fourth where cnd.CS124=A: (
    time: 51 db: 41)
19 QUERY XML select * from Fourth where cnd.CS107=A: (
    time: 58 db: 46)
```

- 20 QUERY XML select * from Fourth where cnd.CS128=A: (
time: 51 db: 42)
- 21 QUERY XML select * from Fourth where cnd.CS130=B: (
time: 63 db: 51)

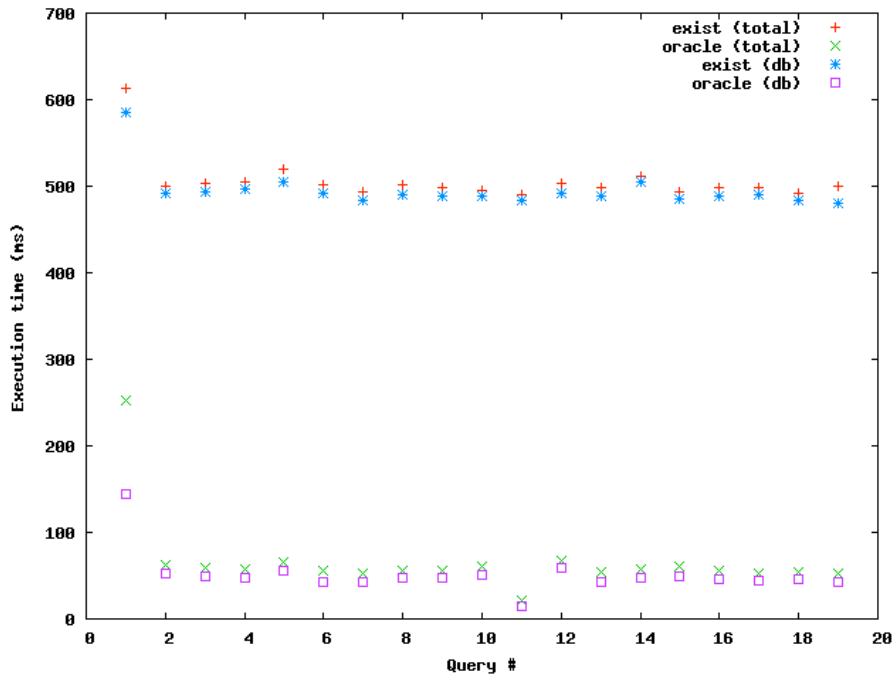
Exist Data

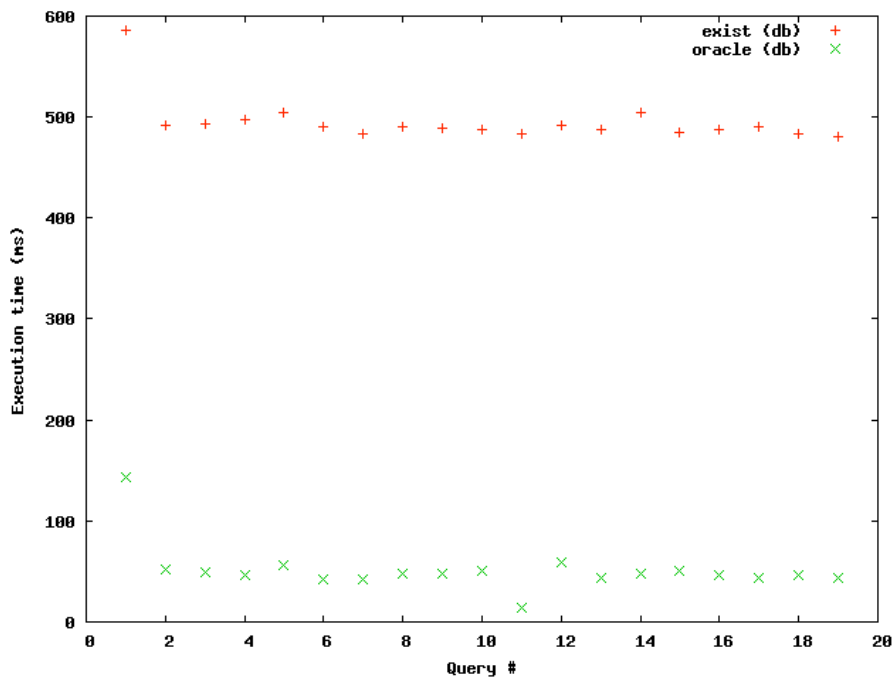
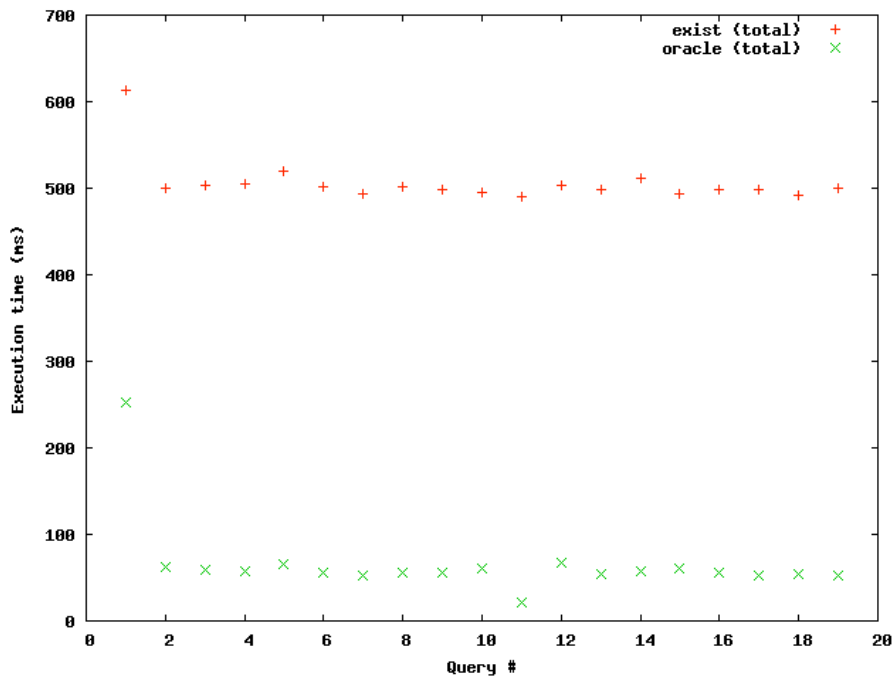
- 1 QUERY XML select * from First where cnd.CS120=B: (
time: 998 db: 987)
- 2 QUERY XML select * from First where cnd.CS101=B: (
time: 896 db: 887)
- 3 QUERY XML select * from First where cnd.CS106=A: (
time: 894 db: 883)
- 4 QUERY XML select * from First where cnd.CS143=B: (
time: 880 db: 872)
- 5 QUERY XML select * from First where cnd.CS103=B: (
time: 896 db: 888)
- 6 QUERY XML select * from Second where cnd.CS143=B: (
time: 890 db: 878)
- 7 QUERY XML select * from Second where cnd.CS146=A: (
time: 885 db: 874)
- 8 QUERY XML select * from Second where cnd.CS135=A: (
time: 880 db: 872)
- 9 QUERY XML select * from Second where cnd.CS102=A: (
time: 899 db: 892)
- 10 QUERY XML select * from Second where cnd.CS127=B: (
time: 883 db: 871)
- 11 QUERY XML select * from Second where cnd.CS144=A: (
time: 887 db: 880)
- 12 QUERY XML select * from Third where cnd.CS144=A: (
time: 883 db: 871)
- 13 QUERY XML select * from Third where cnd.CS138=A: (
time: 885 db: 877)
- 14 QUERY XML select * from Third where cnd.CS112=B: (
time: 882 db: 873)
- 15 QUERY XML select * from Third where cnd.CS140=A: (
time: 890 db: 873)
- 16 QUERY XML select * from Third where cnd.CS122=A: (
time: 900 db: 892)

- 17 QUERY XML select * from Fourth where cnd.CS117=B: (time: 888 db: 877)
- 18 QUERY XML select * from Fourth where cnd.CS124=A: (time: 882 db: 874)
- 19 QUERY XML select * from Fourth where cnd.CS107=A: (time: 888 db: 880)
- 20 QUERY XML select * from Fourth where cnd.CS128=A: (time: 894 db: 885)
- 21 QUERY XML select * from Fourth where cnd.CS130=B: (time: 881 db: 872)

B.2.4 Select on Variables

Comparison





Oracle Data

```
1 QUERY XML select * from First where var.CS299 in V:
  (time: 252 db: 144)
2 QUERY XML select * from First where var.CS205 in V:
  (time: 63 db: 52)
3 QUERY XML select * from First where var.CS286 in V:
  (time: 59 db: 49)
4 QUERY XML select * from First where var.CS268 in V:
  (time: 57 db: 47)
5 QUERY XML select * from Second where var.CS211 in V:
  (time: 65 db: 56)
6 QUERY XML select * from Second where var.CS141 in V:
  (time: 56 db: 42)
7 QUERY XML select * from Second where var.CS279 in V:
  (time: 52 db: 42)
8 QUERY XML select * from Second where var.CS269 in V:
  (time: 56 db: 48)
9 QUERY XML select * from Second where var.CS251 in V:
  (time: 55 db: 48)
10 QUERY XML select * from Third where var.CS127 in V:
  (time: 61 db: 51)
11 QUERY XML select * from Third where var.CS200 in V:
  (time: 21 db: 14)
12 QUERY XML select * from Third where var.CS112 in V:
  (time: 67 db: 59)
13 QUERY XML select * from Third where var.CS224 in V:
  (time: 54 db: 43)
14 QUERY XML select * from Third where var.CS239 in V:
  (time: 57 db: 48)
15 QUERY XML select * from Fourth where var.CS124 in V:
  (time: 60 db: 50)
16 QUERY XML select * from Fourth where var.CS157 in V:
  (time: 56 db: 46)
17 QUERY XML select * from Fourth where var.CS270 in V:
  (time: 52 db: 44)
18 QUERY XML select * from Fourth where var.CS129 in V:
  (time: 54 db: 46)
19 QUERY XML select * from Fourth where var.CS166 in V:
  (time: 52 db: 43)
```

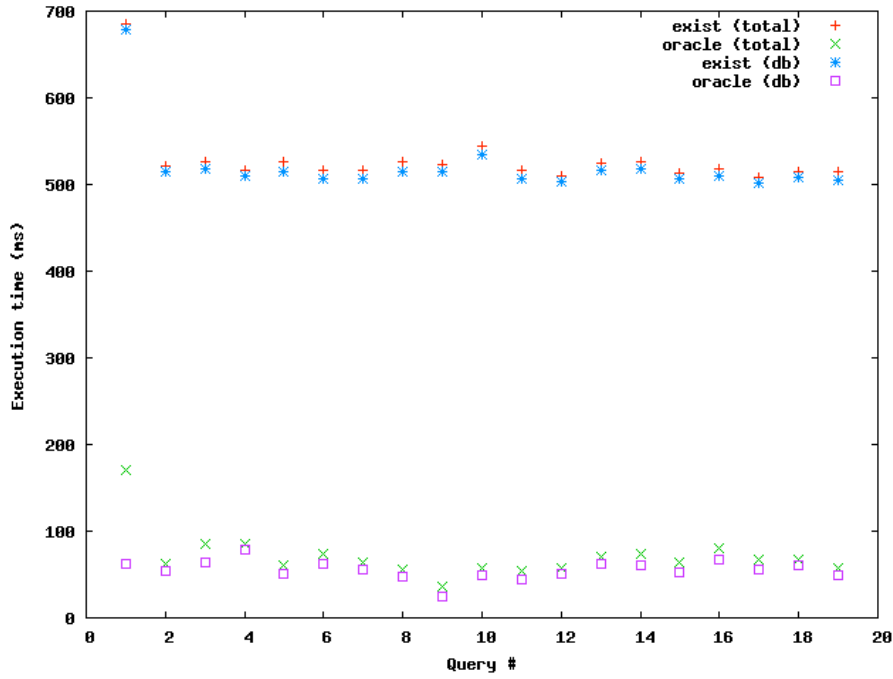
Exist Data

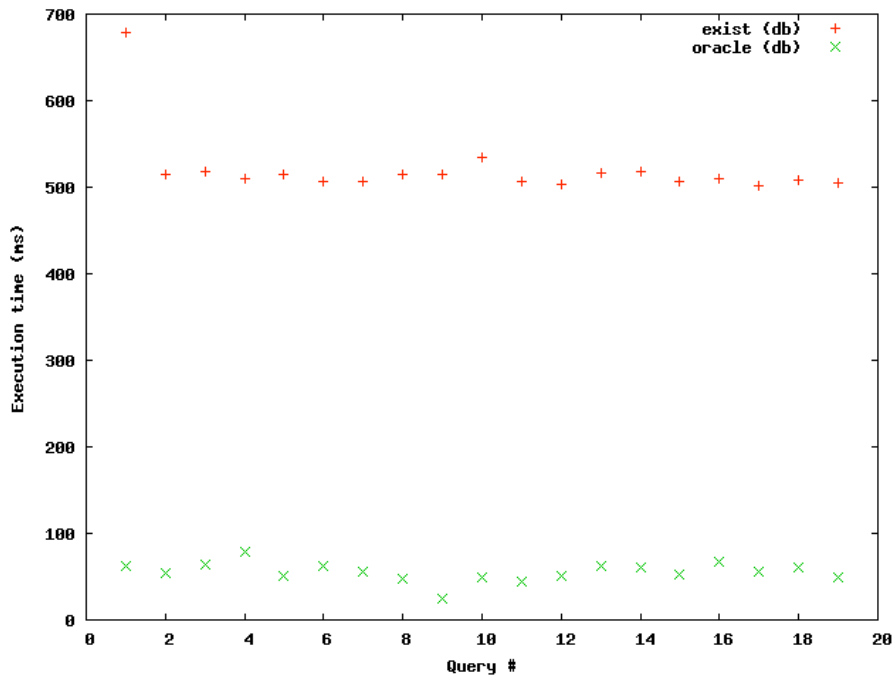
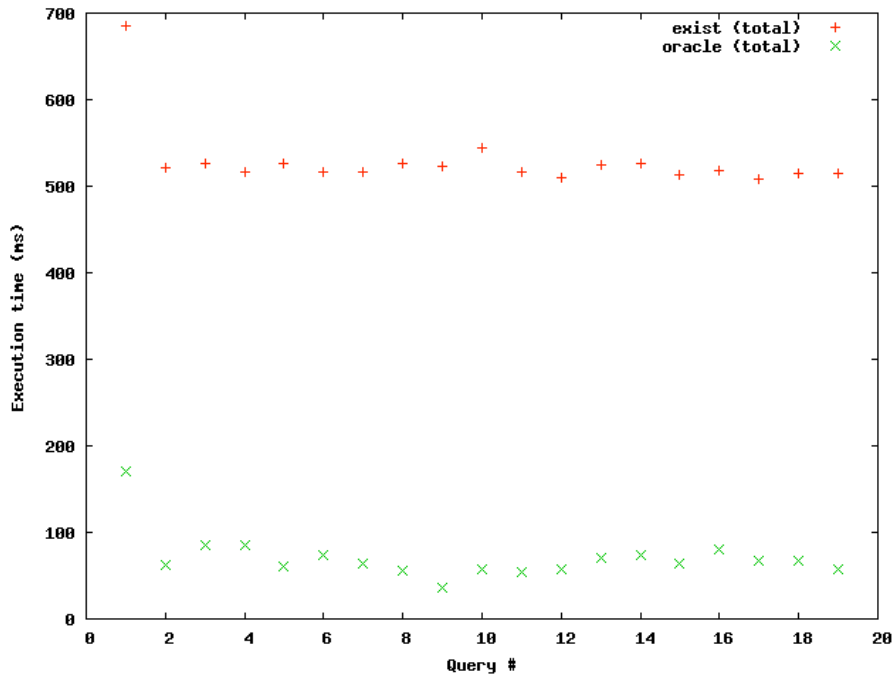
- 1 QUERY XML select * from First where var.CS299 in V:
(time: 613 db: 586)
- 2 QUERY XML select * from First where var.CS205 in V:
(time: 500 db: 492)
- 3 QUERY XML select * from First where var.CS286 in V:
(time: 504 db: 493)
- 4 QUERY XML select * from First where var.CS268 in V:
(time: 505 db: 497)
- 5 QUERY XML select * from Second where var.CS211 in V:
(time: 519 db: 505)
- 6 QUERY XML select * from Second where var.CS141 in V:
(time: 502 db: 491)
- 7 QUERY XML select * from Second where var.CS279 in V:
(time: 494 db: 484)
- 8 QUERY XML select * from Second where var.CS269 in V:
(time: 502 db: 490)
- 9 QUERY XML select * from Second where var.CS251 in V:
(time: 498 db: 489)
- 10 QUERY XML select * from Third where var.CS127 in V:
(time: 495 db: 488)
- 11 QUERY XML select * from Third where var.CS200 in V:
(time: 490 db: 483)
- 12 QUERY XML select * from Third where var.CS112 in V:
(time: 503 db: 492)
- 13 QUERY XML select * from Third where var.CS224 in V:
(time: 499 db: 488)
- 14 QUERY XML select * from Third where var.CS239 in V:
(time: 512 db: 505)
- 15 QUERY XML select * from Fourth where var.CS124 in V:
(time: 494 db: 485)
- 16 QUERY XML select * from Fourth where var.CS157 in V:
(time: 498 db: 488)
- 17 QUERY XML select * from Fourth where var.CS270 in V:
(time: 498 db: 490)
- 18 QUERY XML select * from Fourth where var.CS129 in V:
(time: 492 db: 483)

19 QUERY XML select * from Fourth where var.CS166 in V:
(time: 500 db: 481)

B.2.5 Select on Table

Comparison





Oracle Data

```
1 QUERY XML select * from First where tbl.CS239=B: (
    time: 171 db: 62)
2 QUERY XML select * from First where tbl.CS255=A: (
    time: 63 db: 54)
3 QUERY XML select * from First where tbl.CS211=A: (
    time: 85 db: 64)
4 QUERY XML select * from First where tbl.CS217=B: (
    time: 86 db: 78)
5 QUERY XML select * from First where tbl.CS219=A: (
    time: 60 db: 51)
6 QUERY XML select * from Second where tbl.CS292=B: (
    time: 73 db: 63)
7 QUERY XML select * from Second where tbl.CS107=A: (
    time: 64 db: 55)
8 QUERY XML select * from Second where tbl.CS244=B: (
    time: 55 db: 47)
9 QUERY XML select * from Second where tbl.CS279=B: (
    time: 36 db: 24)
10 QUERY XML select * from Second where tbl.CS183=B: (
    time: 58 db: 50)
11 QUERY XML select * from Third where tbl.CS277=A: (
    time: 54 db: 45)
12 QUERY XML select * from Third where tbl.CS228=B: (
    time: 58 db: 51)
13 QUERY XML select * from Third where tbl.CS298=A: (
    time: 71 db: 63)
14 QUERY XML select * from Third where tbl.CS281=B: (
    time: 73 db: 60)
15 QUERY XML select * from Third where tbl.CS183=B: (
    time: 64 db: 52)
16 QUERY XML select * from Fourth where tbl.CS136=A: (
    time: 81 db: 67)
17 QUERY XML select * from Fourth where tbl.CS251=B: (
    time: 67 db: 56)
18 QUERY XML select * from Fourth where tbl.CS277=B: (
    time: 68 db: 60)
19 QUERY XML select * from Fourth where tbl.CS249=B: (
    time: 58 db: 50)
```

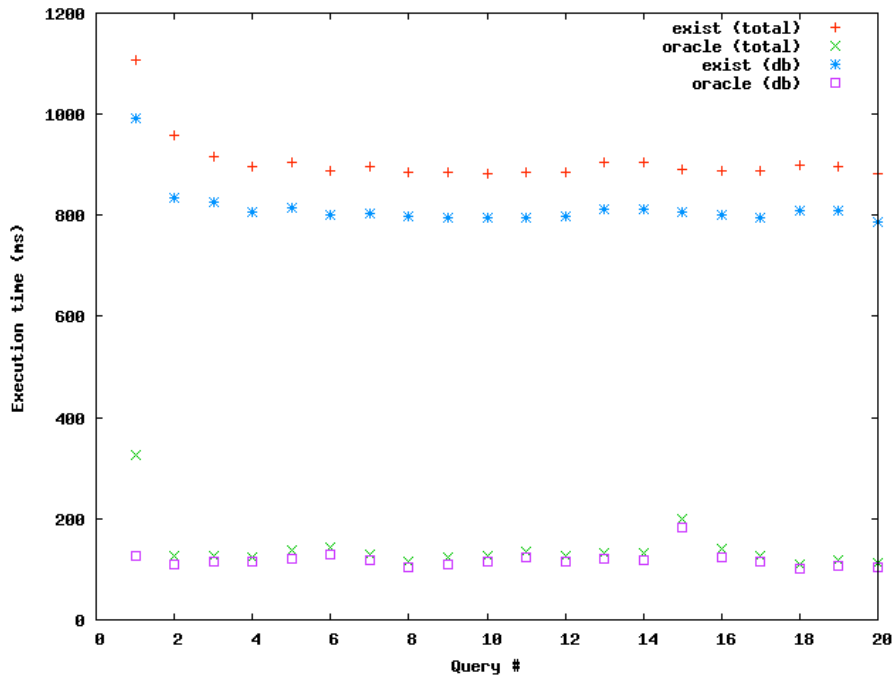
Exist Data

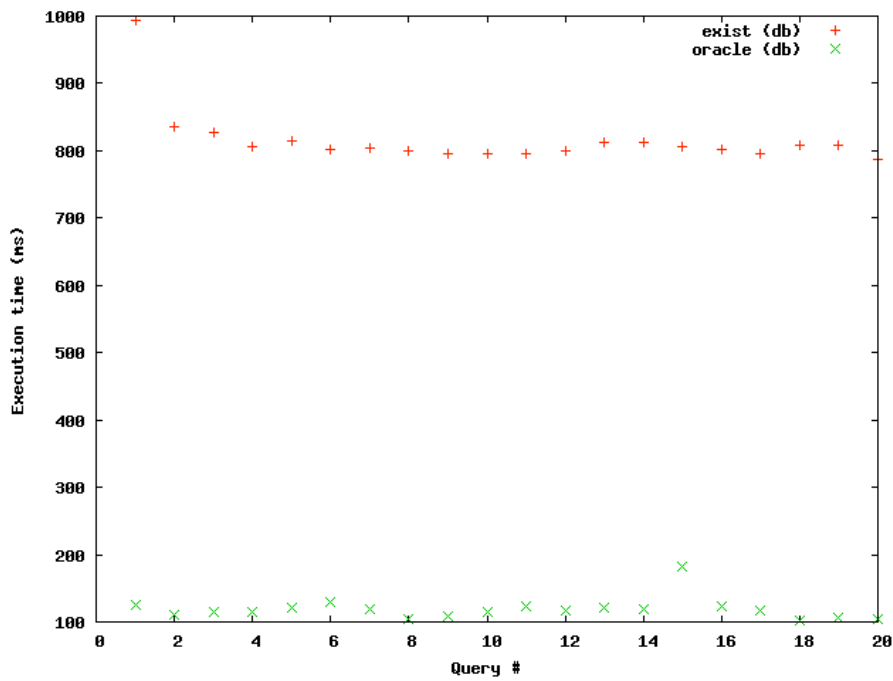
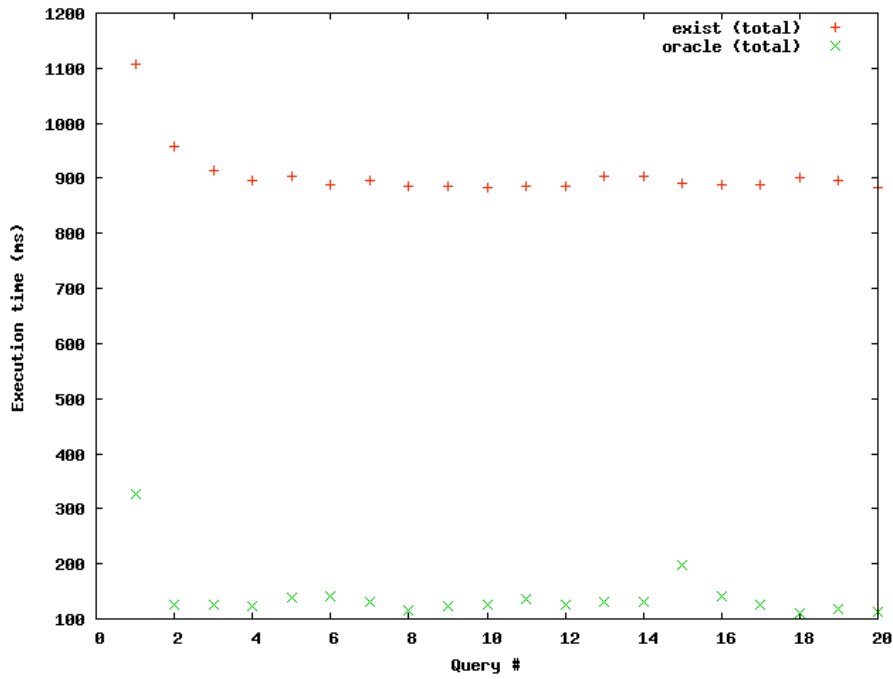
- 1 QUERY XML select * from First where tbl.CS239=B: (time: 686 db: 679)
- 2 QUERY XML select * from First where tbl.CS255=A: (time: 521 db: 514)
- 3 QUERY XML select * from First where tbl.CS211=A: (time: 527 db: 518)
- 4 QUERY XML select * from First where tbl.CS217=B: (time: 517 db: 510)
- 5 QUERY XML select * from First where tbl.CS219=A: (time: 527 db: 515)
- 6 QUERY XML select * from Second where tbl.CS292=B: (time: 516 db: 507)
- 7 QUERY XML select * from Second where tbl.CS107=A: (time: 516 db: 507)
- 8 QUERY XML select * from Second where tbl.CS244=B: (time: 526 db: 514)
- 9 QUERY XML select * from Second where tbl.CS279=B: (time: 523 db: 515)
- 10 QUERY XML select * from Second where tbl.CS183=B: (time: 545 db: 534)
- 11 QUERY XML select * from Third where tbl.CS277=A: (time: 516 db: 506)
- 12 QUERY XML select * from Third where tbl.CS228=B: (time: 510 db: 504)
- 13 QUERY XML select * from Third where tbl.CS298=A: (time: 525 db: 517)
- 14 QUERY XML select * from Third where tbl.CS281=B: (time: 527 db: 518)
- 15 QUERY XML select * from Third where tbl.CS183=B: (time: 513 db: 506)
- 16 QUERY XML select * from Fourth where tbl.CS136=A: (time: 518 db: 510)
- 17 QUERY XML select * from Fourth where tbl.CS251=B: (time: 509 db: 501)
- 18 QUERY XML select * from Fourth where tbl.CS277=B: (time: 515 db: 508)

19 QUERY XML select * from Fourth where tbl.CS249=B: (
time: 514 db: 505)

B.2.6 Project on Context

Comparison





Oracle Data

```
1 QUERY XML select cnt.year from First: (time: 327 db:
  126)
2 QUERY XML select cnt.college from First: (time: 126
  db: 110)
3 QUERY XML select cnt.comments from First: (time: 127
  db: 114)
4 QUERY XML select cnt.major from First: (time: 124 db
  : 115)
5 QUERY XML select cnt.semester from First: (time: 138
  db: 121)
6 QUERY XML select cnt.instructor from Second: (time:
  142 db: 129)
7 QUERY XML select cnt.comments from Second: (time:
  130 db: 118)
8 QUERY XML select cnt.semester from Second: (time:
  115 db: 104)
9 QUERY XML select cnt.college from Second: (time: 123
  db: 109)
10 QUERY XML select cnt.major from Second: (time: 126
  db: 115)
11 QUERY XML select cnt.year from Third: (time: 135 db:
  123)
12 QUERY XML select cnt.major from Third: (time: 127 db
  : 116)
13 QUERY XML select cnt.instructor from Third: (time:
  131 db: 121)
14 QUERY XML select cnt.college from Third: (time: 132
  db: 118)
15 QUERY XML select cnt.comments from Third: (time: 199
  db: 182)
16 QUERY XML select cnt.instructor from Fourth: (time:
  140 db: 124)
17 QUERY XML select cnt.year from Fourth: (time: 127 db
  : 116)
18 QUERY XML select cnt.comments from Fourth: (time:
  111 db: 102)
19 QUERY XML select cnt.major from Fourth: (time: 117
  db: 107)
```

20 QUERY XML select cnt.college from Fourth: (time: 113
db: 104)

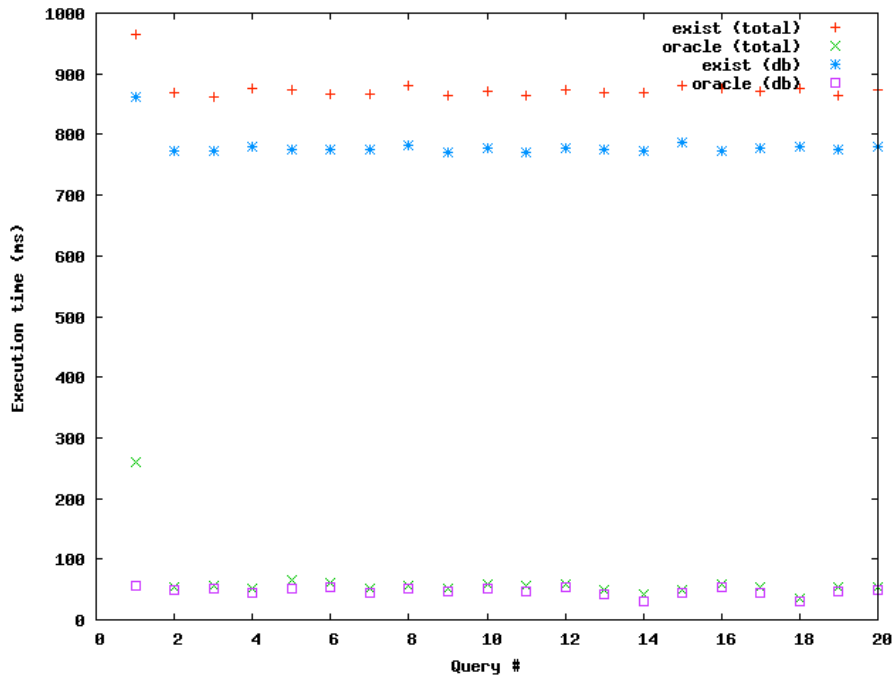
Exist Data

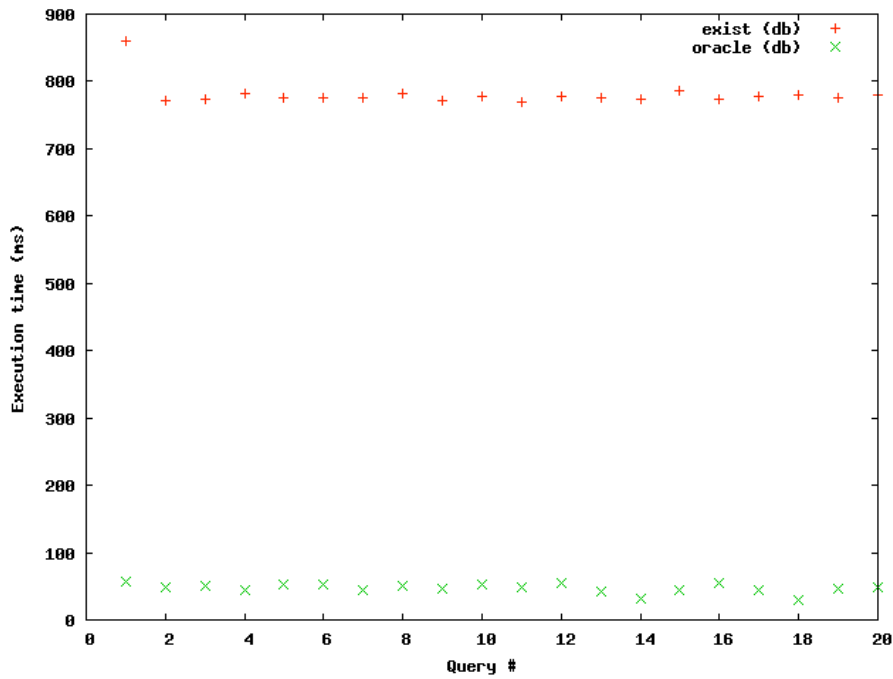
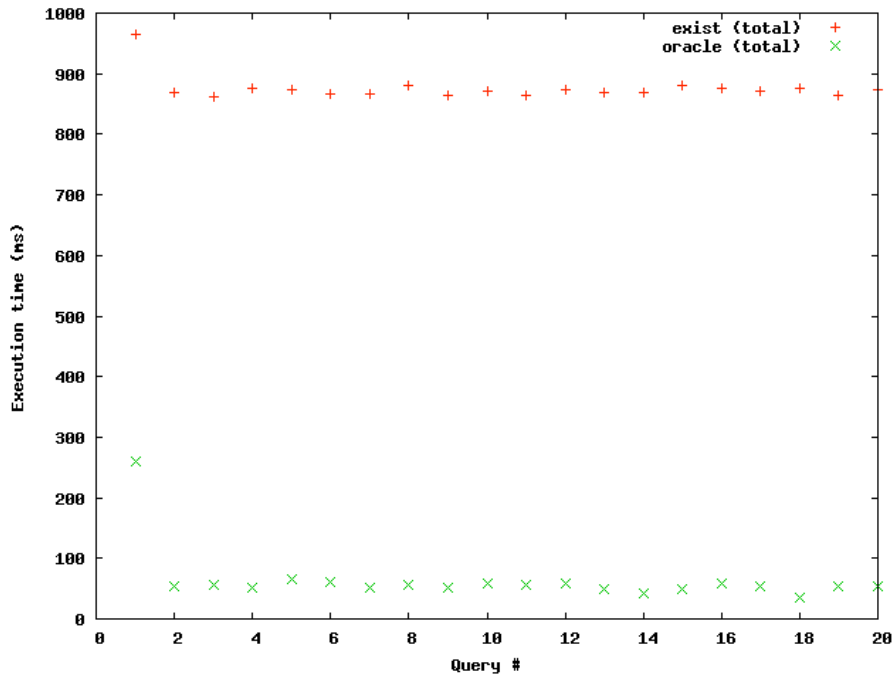
- 1 QUERY XML select cnt.year from First: (time: 1107 db
: 993)
- 2 QUERY XML select cnt.college from First: (time: 958
db: 836)
- 3 QUERY XML select cnt.comments from First: (time: 915
db: 827)
- 4 QUERY XML select cnt.major from First: (time: 896 db
: 807)
- 5 QUERY XML select cnt.semester from First: (time: 904
db: 814)
- 6 QUERY XML select cnt.instructor from Second: (time:
889 db: 802)
- 7 QUERY XML select cnt.comments from Second: (time:
897 db: 804)
- 8 QUERY XML select cnt.semester from Second: (time:
886 db: 799)
- 9 QUERY XML select cnt.college from Second: (time: 885
db: 795)
- 10 QUERY XML select cnt.major from Second: (time: 882
db: 795)
- 11 QUERY XML select cnt.year from Third: (time: 885 db:
796)
- 12 QUERY XML select cnt.major from Third: (time: 886 db
: 799)
- 13 QUERY XML select cnt.instructor from Third: (time:
905 db: 813)
- 14 QUERY XML select cnt.college from Third: (time: 905
db: 813)
- 15 QUERY XML select cnt.comments from Third: (time: 891
db: 806)
- 16 QUERY XML select cnt.instructor from Fourth: (time:
888 db: 802)
- 17 QUERY XML select cnt.year from Fourth: (time: 889 db
: 796)

- 18 QUERY XML select cnt.comments from Fourth: (time: 900 db: 809)
- 19 QUERY XML select cnt.major from Fourth: (time: 896 db: 808)
- 20 QUERY XML select cnt.college from Fourth: (time: 882 db: 787)

B.2.7 Project on Conditional

Comparison





Oracle Data

```
1 QUERY XML select cnd.CS115 from First: (time: 260 db
  : 57)
2 QUERY XML select cnd.CS146 from First: (time: 55 db:
  49)
3 QUERY XML select cnd.CS105 from First: (time: 57 db:
  51)
4 QUERY XML select cnd.CS134 from First: (time: 51 db:
  44)
5 QUERY XML select cnd.CS130 from First: (time: 65 db:
  52)
6 QUERY XML select cnd.CS132 from Second: (time: 60 db
  : 53)
7 QUERY XML select cnd.CS146 from Second: (time: 52 db
  : 45)
8 QUERY XML select cnd.CS109 from Second: (time: 56 db
  : 51)
9 QUERY XML select cnd.CS100 from Second: (time: 51 db
  : 46)
10 QUERY XML select cnd.CS122 from Second: (time: 59 db
  : 52)
11 QUERY XML select cnd.CS105 from Third: (time: 56 db:
  48)
12 QUERY XML select cnd.CS117 from Third: (time: 58 db:
  54)
13 QUERY XML select cnd.CS147 from Third: (time: 49 db:
  42)
14 QUERY XML select cnd.CS122 from Third: (time: 41 db:
  31)
15 QUERY XML select cnd.CS123 from Third: (time: 49 db:
  44)
16 QUERY XML select cnd.CS111 from Fourth: (time: 59 db
  : 54)
17 QUERY XML select cnd.CS119 from Fourth: (time: 53 db
  : 45)
18 QUERY XML select cnd.CS128 from Fourth: (time: 36 db
  : 30)
19 QUERY XML select cnd.CS132 from Fourth: (time: 53 db
  : 47)
```

20 QUERY XML select cnd.CS116 from Fourth: (time: 54 db
: 49)

Exist Data

1 QUERY XML select cnd.CS115 from First: (time: 965 db
: 861)

2 QUERY XML select cnd.CS146 from First: (time: 869 db
: 772)

3 QUERY XML select cnd.CS105 from First: (time: 862 db
: 774)

4 QUERY XML select cnd.CS134 from First: (time: 875 db
: 781)

5 QUERY XML select cnd.CS130 from First: (time: 874 db
: 776)

6 QUERY XML select cnd.CS132 from Second: (time: 866
db: 776)

7 QUERY XML select cnd.CS146 from Second: (time: 867
db: 776)

8 QUERY XML select cnd.CS109 from Second: (time: 881
db: 783)

9 QUERY XML select cnd.CS100 from Second: (time: 864
db: 771)

10 QUERY XML select cnd.CS122 from Second: (time: 871
db: 777)

11 QUERY XML select cnd.CS105 from Third: (time: 864 db
: 770)

12 QUERY XML select cnd.CS117 from Third: (time: 873 db
: 777)

13 QUERY XML select cnd.CS147 from Third: (time: 869 db
: 776)

14 QUERY XML select cnd.CS122 from Third: (time: 868 db
: 774)

15 QUERY XML select cnd.CS123 from Third: (time: 880 db
: 787)

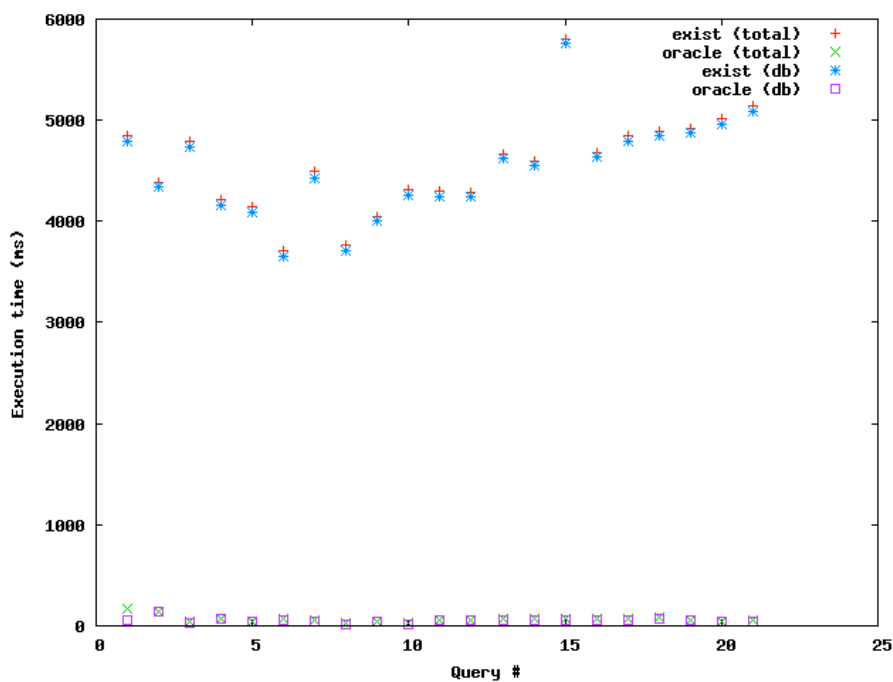
16 QUERY XML select cnd.CS111 from Fourth: (time: 875
db: 773)

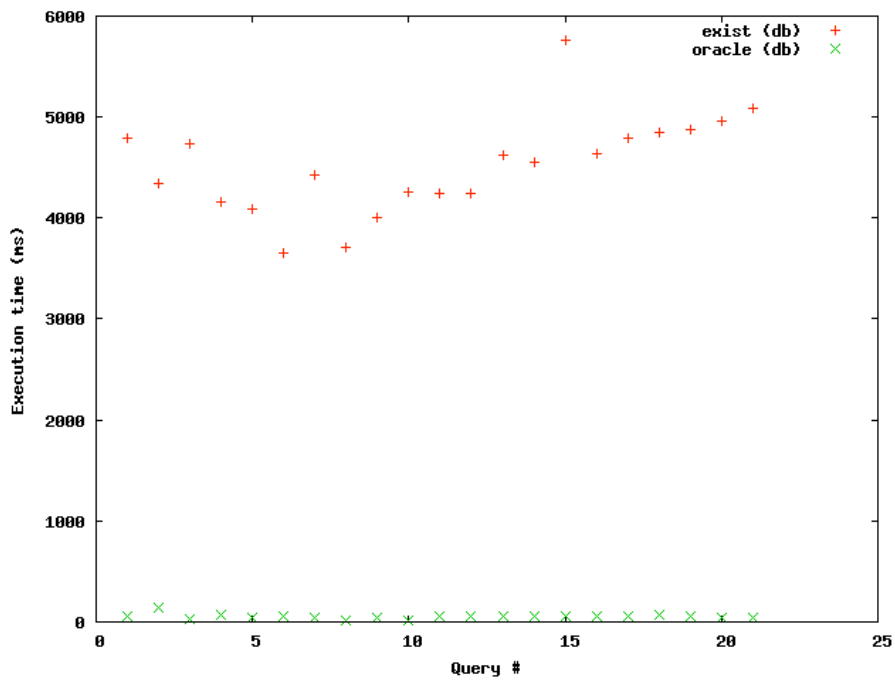
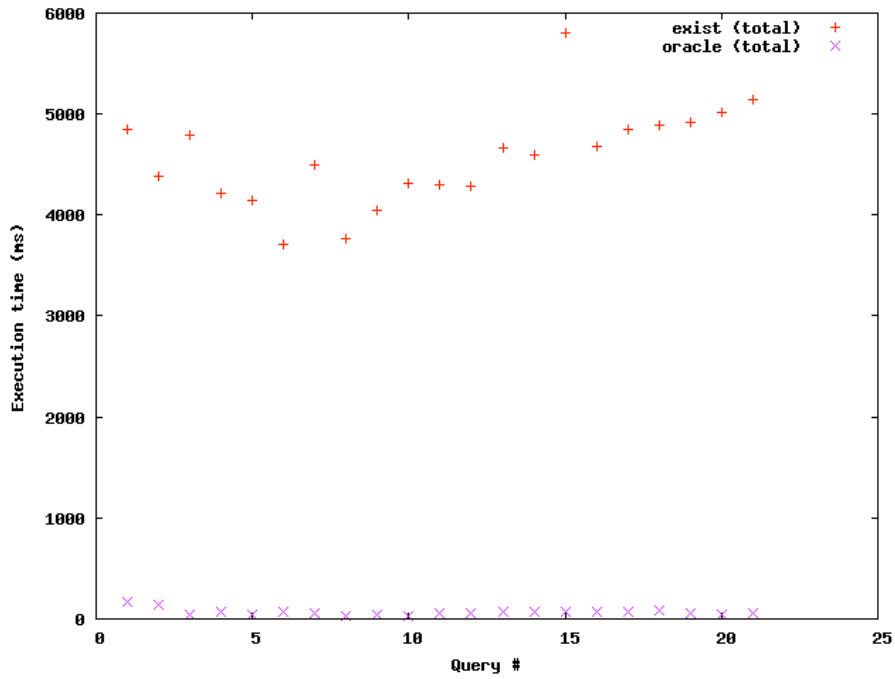
17 QUERY XML select cnd.CS119 from Fourth: (time: 872
db: 778)

- 18 QUERY XML select cnd.CS128 from Fourth: (time: 875 db: 779)
- 19 QUERY XML select cnd.CS132 from Fourth: (time: 865 db: 776)
- 20 QUERY XML select cnd.CS116 from Fourth: (time: 873 db: 779)

B.2.8 Project on Variables

Comparison





Oracle Data

```
1 QUERY XML select var.CS214 from First: (time: 171 db
  : 61)
2 QUERY XML select var.CS249 from First: (time: 145 db
  : 139)
3 QUERY XML select var.CS239 from First: (time: 38 db:
  31)
4 QUERY XML select var.CS282 from First: (time: 76 db:
  69)
5 QUERY XML select var.CS299 from First: (time: 41 db:
  36)
6 QUERY XML select var.CS297 from Second: (time: 67 db
  : 60)
7 QUERY XML select var.CS265 from Second: (time: 56 db
  : 49)
8 QUERY XML select var.CS265 from Second: (time: 25 db
  : 18)
9 QUERY XML select var.CS244 from Second: (time: 45 db
  : 39)
10 QUERY XML select var.CS244 from Second: (time: 24 db
  : 19)
11 QUERY XML select var.CS128 from Second: (time: 58 db
  : 53)
12 QUERY XML select var.CS244 from Third: (time: 61 db:
  56)
13 QUERY XML select var.CS299 from Third: (time: 64 db:
  57)
14 QUERY XML select var.CS113 from Third: (time: 65 db:
  58)
15 QUERY XML select var.CS192 from Third: (time: 66 db:
  60)
16 QUERY XML select var.CS169 from Third: (time: 68 db:
  60)
17 QUERY XML select var.CS126 from Fourth: (time: 64 db
  : 59)
18 QUERY XML select var.CS210 from Fourth: (time: 83 db
  : 76)
19 QUERY XML select var.CS164 from Fourth: (time: 56 db
  : 50)
```

20 QUERY XML select var.CS270 from Fourth: (time: 46 db
: 39)

21 QUERY XML select var.CS213 from Fourth: (time: 56 db
: 49)

Exist Data

1 QUERY XML select var.CS214 from First: (time: 4844
db: 4795)

2 QUERY XML select var.CS249 from First: (time: 4384
db: 4338)

3 QUERY XML select var.CS239 from First: (time: 4798
db: 4741)

4 QUERY XML select var.CS282 from First: (time: 4214
db: 4153)

5 QUERY XML select var.CS299 from First: (time: 4149
db: 4092)

6 QUERY XML select var.CS297 from Second: (time: 3714
db: 3654)

7 QUERY XML select var.CS265 from Second: (time: 4491
db: 4431)

8 QUERY XML select var.CS265 from Second: (time: 3764
db: 3716)

9 QUERY XML select var.CS244 from Second: (time: 4044
db: 3999)

10 QUERY XML select var.CS244 from Second: (time: 4309
db: 4254)

11 QUERY XML select var.CS128 from Second: (time: 4305
db: 4240)

12 QUERY XML select var.CS244 from Third: (time: 4284
db: 4240)

13 QUERY XML select var.CS299 from Third: (time: 4666
db: 4620)

14 QUERY XML select var.CS113 from Third: (time: 4595
db: 4546)

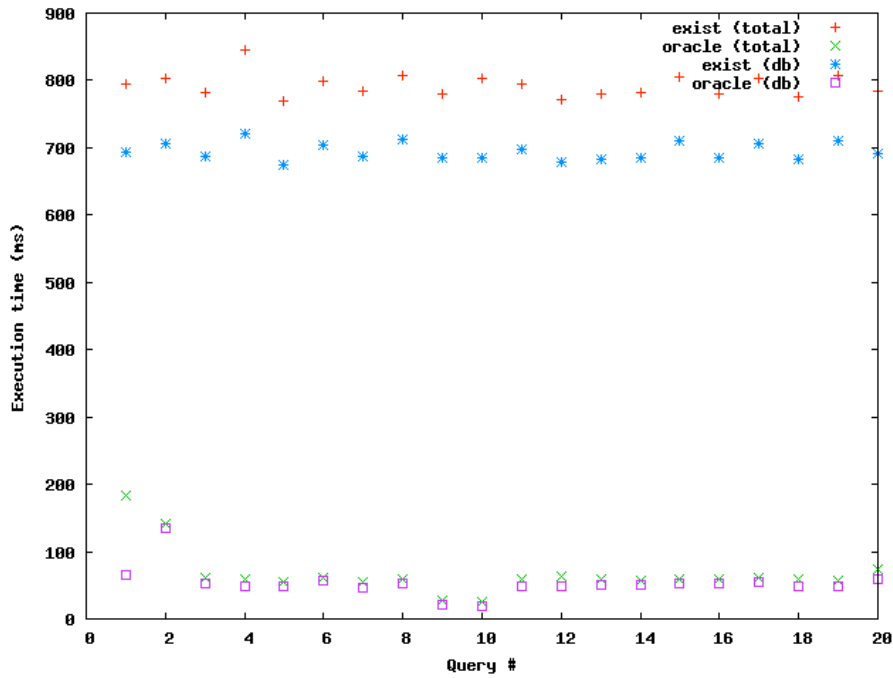
15 QUERY XML select var.CS192 from Third: (time: 5810
db: 5767)

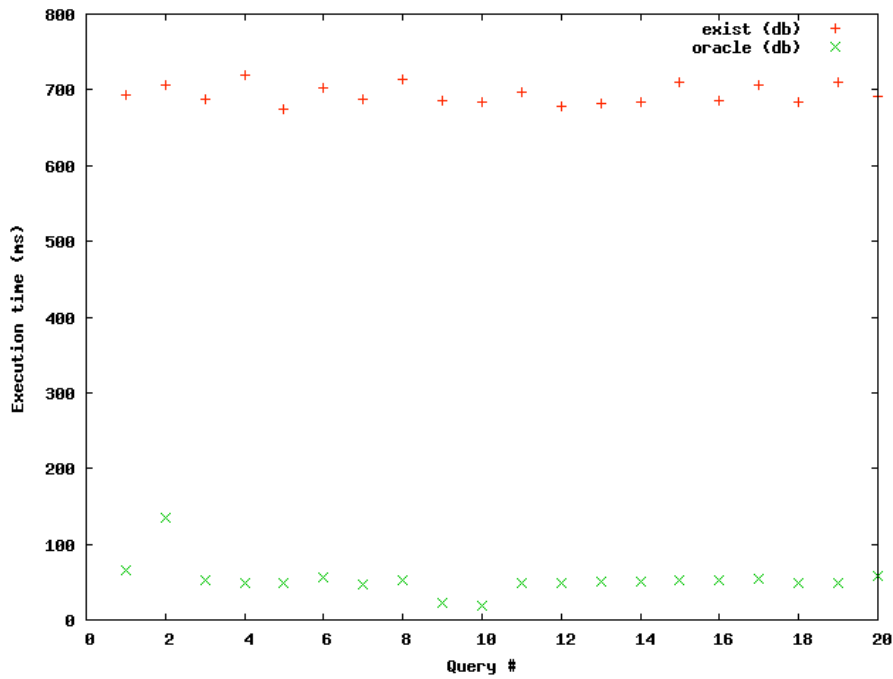
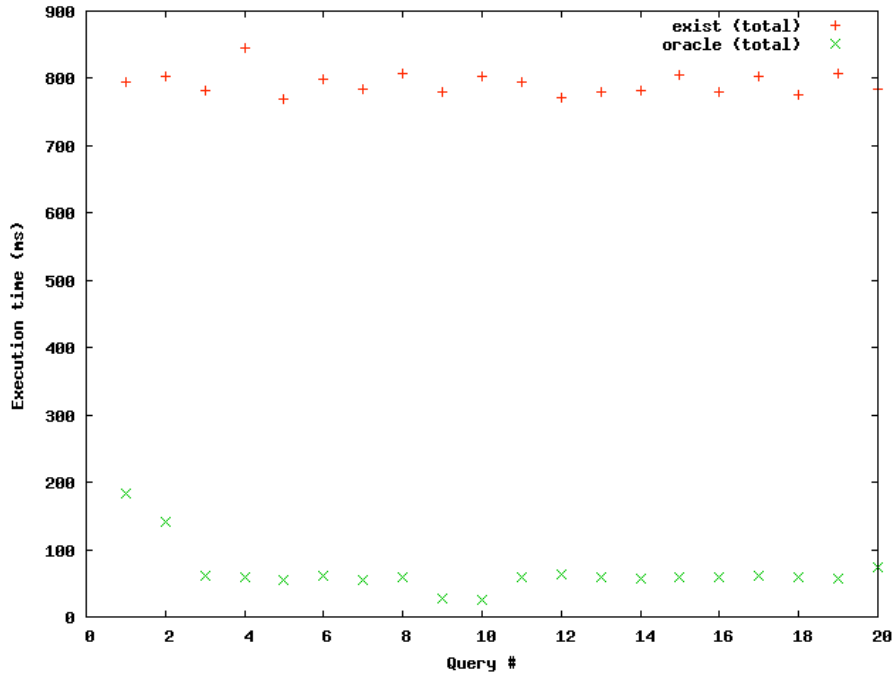
16 QUERY XML select var.CS169 from Third: (time: 4679
db: 4635)

- 17 QUERY XML select var.CS126 from Fourth: (time: 4843
db: 4798)
- 18 QUERY XML select var.CS210 from Fourth: (time: 4894
db: 4851)
- 19 QUERY XML select var.CS164 from Fourth: (time: 4919
db: 4870)
- 20 QUERY XML select var.CS270 from Fourth: (time: 5014
db: 4965)
- 21 QUERY XML select var.CS213 from Fourth: (time: 5137
db: 5091)

B.2.9 Conditional

Comparison





Oracle Data

```
1 QUERY XML select * from First conditional var.CS172=
  A: (time: 184 db: 65)
2 QUERY XML select * from First conditional var.CS113=
  A: (time: 142 db: 135)
3 QUERY XML select * from First conditional var.CS229=
  B: (time: 61 db: 53)
4 QUERY XML select * from First conditional var.CS251=
  B: (time: 58 db: 48)
5 QUERY XML select * from First conditional var.CS256=
  B: (time: 55 db: 49)
6 QUERY XML select * from Second conditional var.CS233
  =A: (time: 62 db: 56)
7 QUERY XML select * from Second conditional var.CS219
  =B: (time: 55 db: 47)
8 QUERY XML select * from Second conditional var.CS190
  =B: (time: 58 db: 52)
9 QUERY XML select * from Second conditional var.CS297
  =B: (time: 28 db: 22)
10 QUERY XML select * from Second conditional var.CS297
  =B: (time: 26 db: 18)
11 QUERY XML select * from Third conditional var.CS138=
  A: (time: 58 db: 49)
12 QUERY XML select * from Third conditional var.CS251=
  A: (time: 64 db: 48)
13 QUERY XML select * from Third conditional var.CS202=
  B: (time: 60 db: 51)
14 QUERY XML select * from Third conditional var.CS144=
  A: (time: 57 db: 51)
15 QUERY XML select * from Third conditional var.CS240=
  B: (time: 59 db: 53)
16 QUERY XML select * from Fourth conditional var.CS112
  =B: (time: 60 db: 53)
17 QUERY XML select * from Fourth conditional var.CS206
  =A: (time: 62 db: 55)
18 QUERY XML select * from Fourth conditional var.CS217
  =A: (time: 59 db: 49)
19 QUERY XML select * from Fourth conditional var.CS275
  =A: (time: 56 db: 49)
```

20 QUERY XML select * from Fourth conditional var.CS111
=A: (time: 73 db: 59)

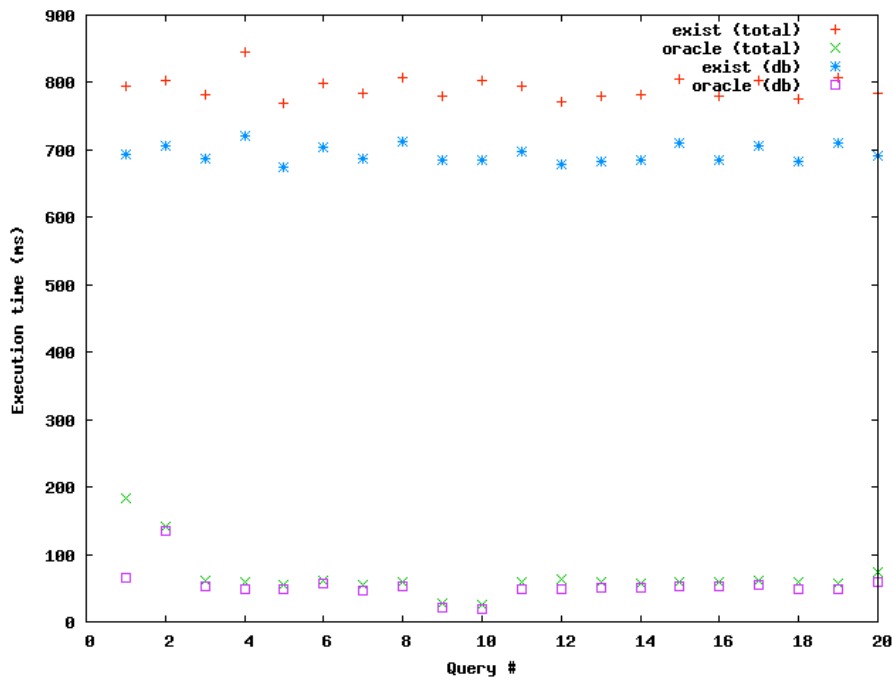
Exist Data

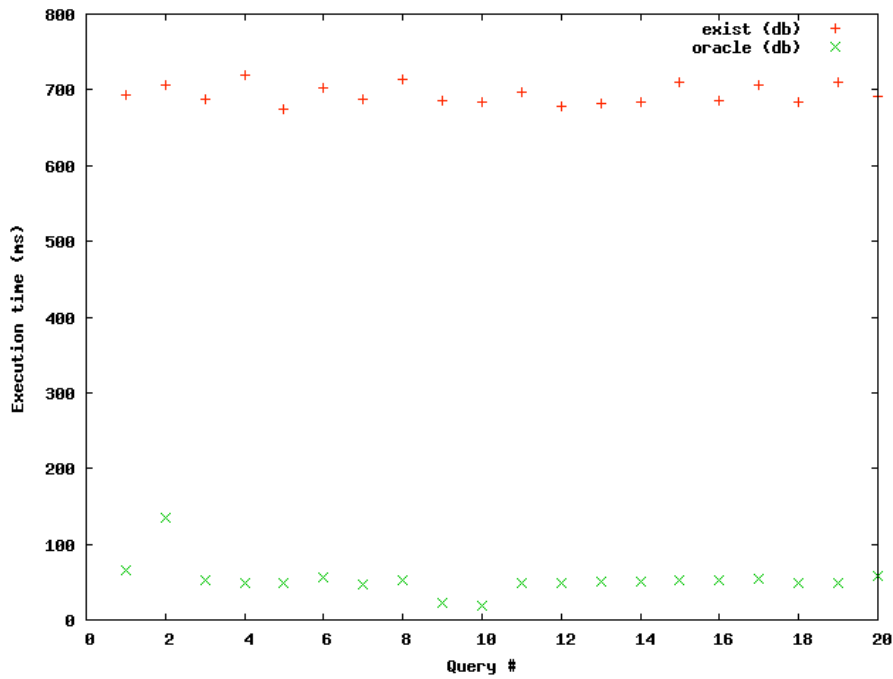
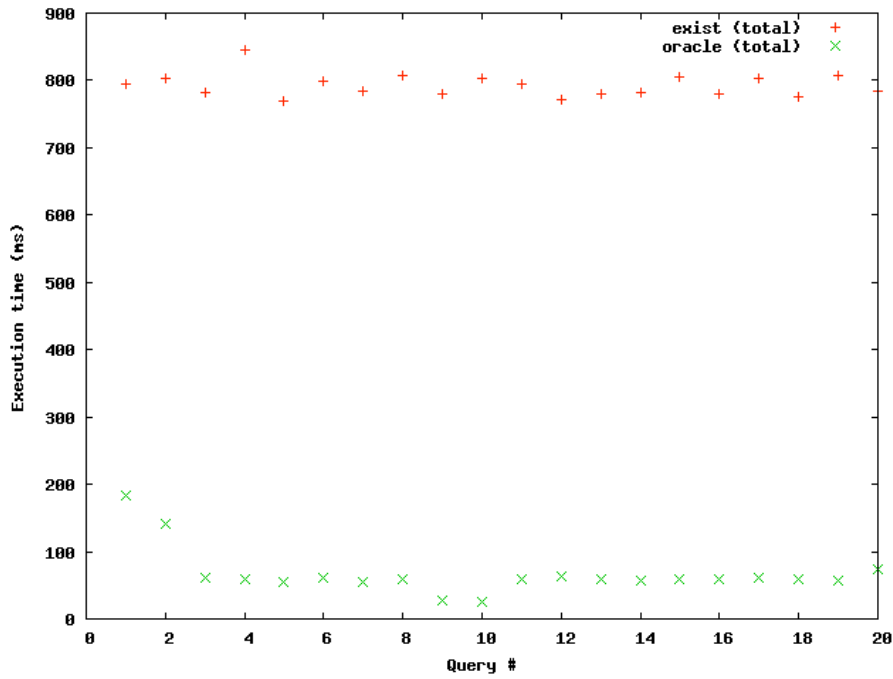
1 QUERY XML select * from First conditional var.CS172=
A: (time: 795 db: 694)
2 QUERY XML select * from First conditional var.CS113=
A: (time: 802 db: 707)
3 QUERY XML select * from First conditional var.CS229=
B: (time: 781 db: 688)
4 QUERY XML select * from First conditional var.CS251=
B: (time: 846 db: 720)
5 QUERY XML select * from First conditional var.CS256=
B: (time: 770 db: 675)
6 QUERY XML select * from Second conditional var.CS233
=A: (time: 799 db: 703)
7 QUERY XML select * from Second conditional var.CS219
=B: (time: 785 db: 688)
8 QUERY XML select * from Second conditional var.CS190
=B: (time: 807 db: 713)
9 QUERY XML select * from Second conditional var.CS297
=B: (time: 780 db: 686)
10 QUERY XML select * from Second conditional var.CS297
=B: (time: 803 db: 684)
11 QUERY XML select * from Third conditional var.CS138=
A: (time: 794 db: 697)
12 QUERY XML select * from Third conditional var.CS251=
A: (time: 772 db: 679)
13 QUERY XML select * from Third conditional var.CS202=
B: (time: 779 db: 682)
14 QUERY XML select * from Third conditional var.CS144=
A: (time: 782 db: 684)
15 QUERY XML select * from Third conditional var.CS240=
B: (time: 805 db: 710)
16 QUERY XML select * from Fourth conditional var.CS112
=B: (time: 780 db: 686)
17 QUERY XML select * from Fourth conditional var.CS206
=A: (time: 803 db: 707)

- 18 QUERY XML select * from Fourth conditional var.CS217
=A: (time: 776 db: 683)
- 19 QUERY XML select * from Fourth conditional var.CS275
=A: (time: 807 db: 710)
- 20 QUERY XML select * from Fourth conditional var.CS111
=A: (time: 785 db: 691)

B.2.10 Conditional

Comparison





Oracle Data

```
1 QUERY XML select * from First conditional var.CS172=
  A: (time: 184 db: 65)
2 QUERY XML select * from First conditional var.CS113=
  A: (time: 142 db: 135)
3 QUERY XML select * from First conditional var.CS229=
  B: (time: 61 db: 53)
4 QUERY XML select * from First conditional var.CS251=
  B: (time: 58 db: 48)
5 QUERY XML select * from First conditional var.CS256=
  B: (time: 55 db: 49)
6 QUERY XML select * from Second conditional var.CS233
  =A: (time: 62 db: 56)
7 QUERY XML select * from Second conditional var.CS219
  =B: (time: 55 db: 47)
8 QUERY XML select * from Second conditional var.CS190
  =B: (time: 58 db: 52)
9 QUERY XML select * from Second conditional var.CS297
  =B: (time: 28 db: 22)
10 QUERY XML select * from Second conditional var.CS297
  =B: (time: 26 db: 18)
11 QUERY XML select * from Third conditional var.CS138=
  A: (time: 58 db: 49)
12 QUERY XML select * from Third conditional var.CS251=
  A: (time: 64 db: 48)
13 QUERY XML select * from Third conditional var.CS202=
  B: (time: 60 db: 51)
14 QUERY XML select * from Third conditional var.CS144=
  A: (time: 57 db: 51)
15 QUERY XML select * from Third conditional var.CS240=
  B: (time: 59 db: 53)
16 QUERY XML select * from Fourth conditional var.CS112
  =B: (time: 60 db: 53)
17 QUERY XML select * from Fourth conditional var.CS206
  =A: (time: 62 db: 55)
18 QUERY XML select * from Fourth conditional var.CS217
  =A: (time: 59 db: 49)
19 QUERY XML select * from Fourth conditional var.CS275
  =A: (time: 56 db: 49)
```

20 QUERY XML select * from Fourth conditional var.CS111
=A: (time: 73 db: 59)

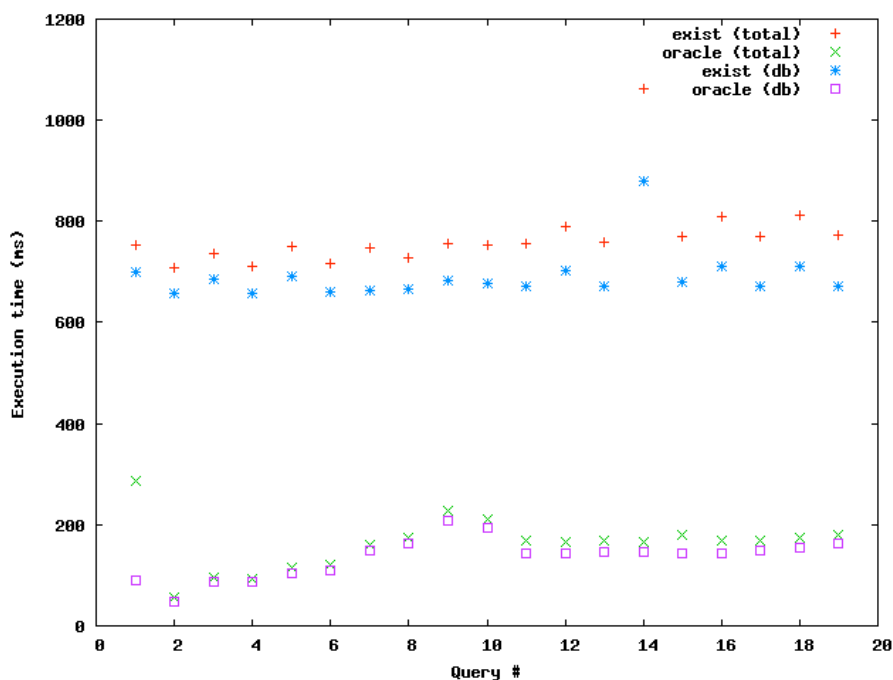
Exist Data

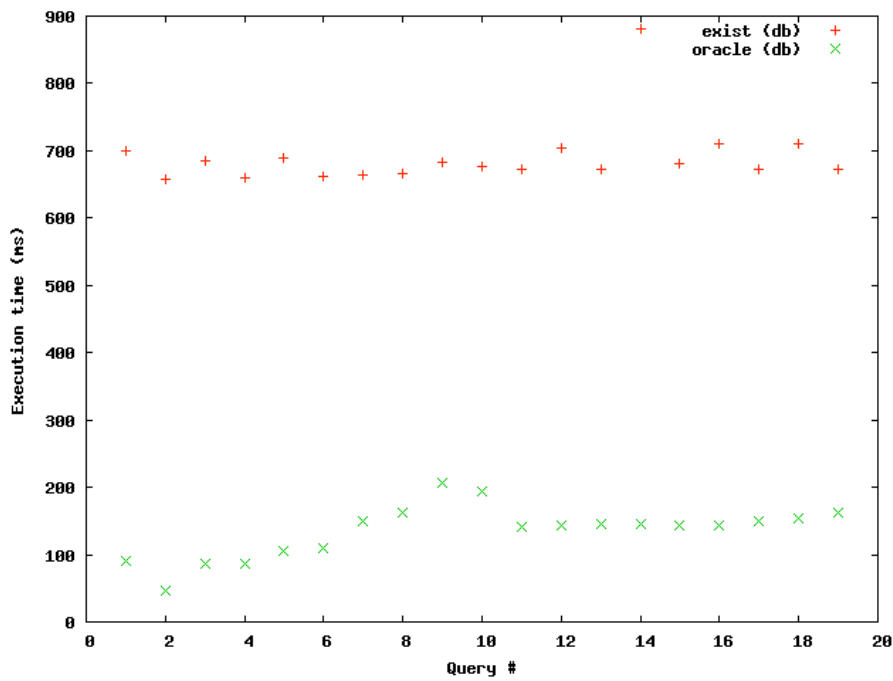
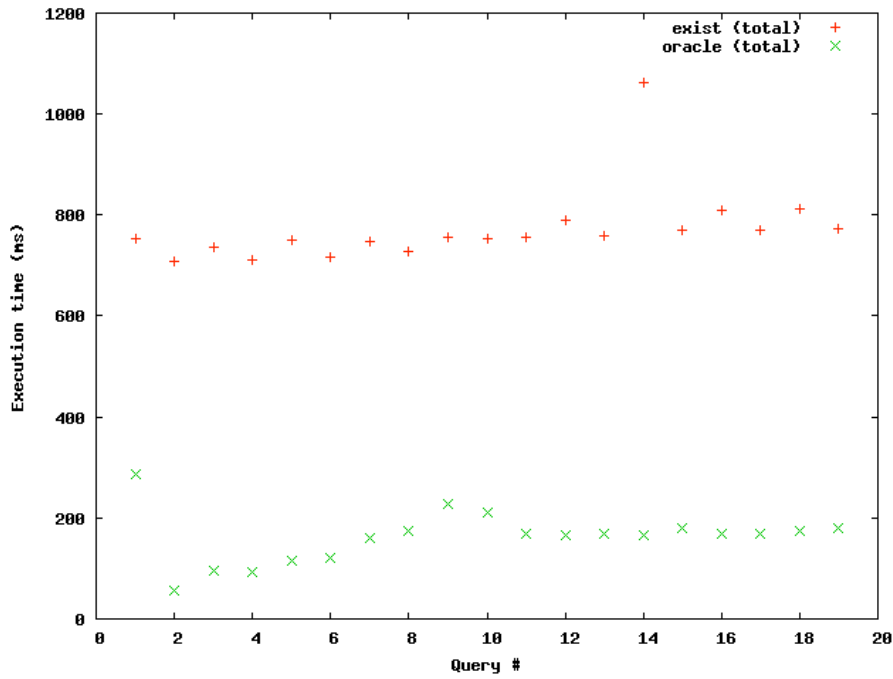
1 QUERY XML select * from First conditional var.CS172=
A: (time: 795 db: 694)
2 QUERY XML select * from First conditional var.CS113=
A: (time: 802 db: 707)
3 QUERY XML select * from First conditional var.CS229=
B: (time: 781 db: 688)
4 QUERY XML select * from First conditional var.CS251=
B: (time: 846 db: 720)
5 QUERY XML select * from First conditional var.CS256=
B: (time: 770 db: 675)
6 QUERY XML select * from Second conditional var.CS233
=A: (time: 799 db: 703)
7 QUERY XML select * from Second conditional var.CS219
=B: (time: 785 db: 688)
8 QUERY XML select * from Second conditional var.CS190
=B: (time: 807 db: 713)
9 QUERY XML select * from Second conditional var.CS297
=B: (time: 780 db: 686)
10 QUERY XML select * from Second conditional var.CS297
=B: (time: 803 db: 684)
11 QUERY XML select * from Third conditional var.CS138=
A: (time: 794 db: 697)
12 QUERY XML select * from Third conditional var.CS251=
A: (time: 772 db: 679)
13 QUERY XML select * from Third conditional var.CS202=
B: (time: 779 db: 682)
14 QUERY XML select * from Third conditional var.CS144=
A: (time: 782 db: 684)
15 QUERY XML select * from Third conditional var.CS240=
B: (time: 805 db: 710)
16 QUERY XML select * from Fourth conditional var.CS112
=B: (time: 780 db: 686)
17 QUERY XML select * from Fourth conditional var.CS206
=A: (time: 803 db: 707)

- 18 QUERY XML select * from Fourth conditional var.CS217
=A: (time: 776 db: 683)
- 19 QUERY XML select * from Fourth conditional var.CS275
=A: (time: 807 db: 710)
- 20 QUERY XML select * from Fourth conditional var.CS111
=A: (time: 785 db: 691)

B.2.11 Select on Probability

Comparison





Oracle Data

```
1 QUERY XML select * from First where tbl.prob>0.64: (
    time: 286 db: 91)
2 QUERY XML select * from First where tbl.prob>0.60: (
    time: 56 db: 47)
3 QUERY XML select * from First where tbl.prob>0.56: (
    time: 95 db: 87)
4 QUERY XML select * from First where tbl.prob>0.52: (
    time: 94 db: 86)
5 QUERY XML select * from First where tbl.prob>0.48: (
    time: 115 db: 105)
6 QUERY XML select * from First where tbl.prob>0.44: (
    time: 120 db: 110)
7 QUERY XML select * from First where tbl.prob>0.40: (
    time: 159 db: 149)
8 QUERY XML select * from First where tbl.prob>0.36: (
    time: 175 db: 162)
9 QUERY XML select * from First where tbl.prob>0.28: (
    time: 228 db: 207)
10 QUERY XML select * from First where tbl.prob>0.24: (
    time: 212 db: 194)
11 QUERY XML select * from First where tbl.prob>0.20: (
    time: 169 db: 142)
12 QUERY XML select * from First where tbl.prob>0.18: (
    time: 165 db: 144)
13 QUERY XML select * from First where tbl.prob>0.16: (
    time: 168 db: 145)
14 QUERY XML select * from First where tbl.prob>0.14: (
    time: 166 db: 146)
15 QUERY XML select * from First where tbl.prob>0.12: (
    time: 180 db: 144)
16 QUERY XML select * from First where tbl.prob>0.10: (
    time: 169 db: 144)
17 QUERY XML select * from First where tbl.prob>0.08: (
    time: 170 db: 149)
18 QUERY XML select * from First where tbl.prob>0.04: (
    time: 174 db: 154)
19 QUERY XML select * from First where tbl.prob>0.02: (
    time: 180 db: 163)
```

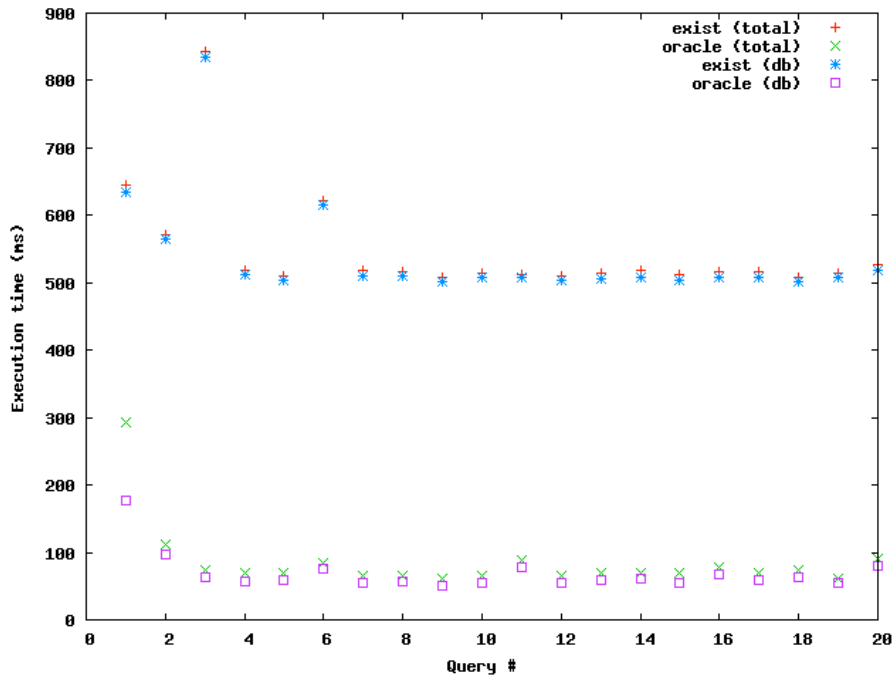
Exist Data

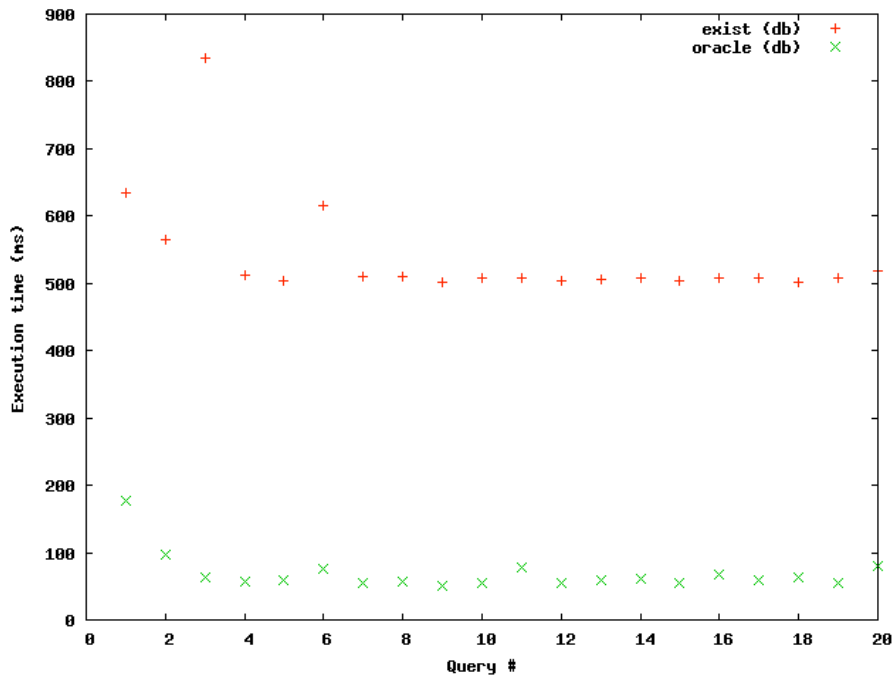
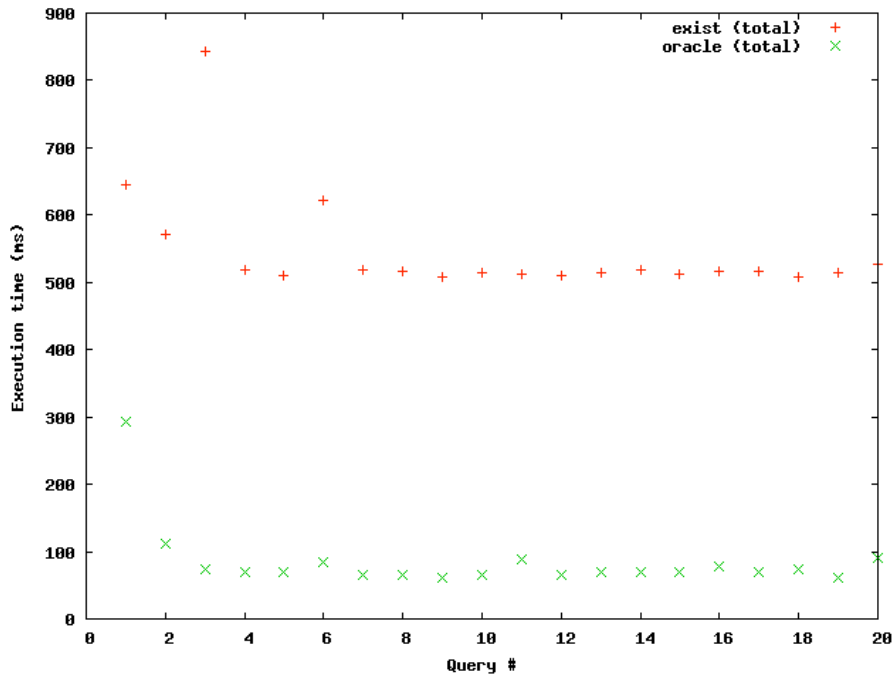
```
1 QUERY XML select * from First where tbl.prob>0.64: (
    time: 754 db: 700)
2 QUERY XML select * from First where tbl.prob>0.60: (
    time: 708 db: 658)
3 QUERY XML select * from First where tbl.prob>0.56: (
    time: 737 db: 685)
4 QUERY XML select * from First where tbl.prob>0.52: (
    time: 711 db: 659)
5 QUERY XML select * from First where tbl.prob>0.48: (
    time: 750 db: 690)
6 QUERY XML select * from First where tbl.prob>0.44: (
    time: 716 db: 661)
7 QUERY XML select * from First where tbl.prob>0.40: (
    time: 748 db: 664)
8 QUERY XML select * from First where tbl.prob>0.36: (
    time: 728 db: 666)
9 QUERY XML select * from First where tbl.prob>0.28: (
    time: 757 db: 682)
10 QUERY XML select * from First where tbl.prob>0.24: (
    time: 753 db: 676)
11 QUERY XML select * from First where tbl.prob>0.20: (
    time: 757 db: 672)
12 QUERY XML select * from First where tbl.prob>0.18: (
    time: 790 db: 703)
13 QUERY XML select * from First where tbl.prob>0.16: (
    time: 758 db: 672)
14 QUERY XML select * from First where tbl.prob>0.14: (
    time: 1062 db: 881)
15 QUERY XML select * from First where tbl.prob>0.12: (
    time: 771 db: 681)
16 QUERY XML select * from First where tbl.prob>0.10: (
    time: 808 db: 711)
17 QUERY XML select * from First where tbl.prob>0.08: (
    time: 769 db: 673)
18 QUERY XML select * from First where tbl.prob>0.04: (
    time: 813 db: 710)
```

19 QUERY XML select * from First where tbl.prob>0.02: (time: 773 db: 672)

B.2.12 Complex Select Conditions

Comparison





Oracle Data

```
1 QUERY XML Select * from First where First.tbl.CS115
    ="A" and First.tbl.prob>0.01: (time: 292 db: 177)
2 QUERY XML Select * from First where First.tbl.CS242
    ="B" and First.tbl.prob>0.03: (time: 111 db: 98)
3 QUERY XML Select * from First where First.tbl.CS219
    ="A" and First.tbl.prob>0.06: (time: 74 db: 63)
4 QUERY XML Select * from First where First.tbl.CS211
    ="B" and First.tbl.prob>0.08: (time: 70 db: 56)
5 QUERY XML Select * from First where First.tbl.CS172
    ="B" and First.tbl.prob>0.09: (time: 70 db: 60)
6 QUERY XML Select * from Second where Second.tbl.
    CS292="B" and Second.tbl.prob>0.01: (time: 84 db:
    76)
7 QUERY XML Select * from Second where Second.tbl.
    CS107="A" and Second.tbl.prob>0.03: (time: 65 db:
    54)
8 QUERY XML Select * from Second where Second.tbl.
    CS244="B" and Second.tbl.prob>0.06: (time: 66 db:
    56)
9 QUERY XML Select * from Second where Second.tbl.
    CS279="B" and Second.tbl.prob>0.08: (time: 61 db:
    50)
10 QUERY XML Select * from Second where Second.tbl.
    CS183="B" and Second.tbl.prob>0.09: (time: 65 db:
    55)
11 QUERY XML Select * from Third where Third.tbl.CS277
    ="A" and Third.tbl.prob>0.01: (time: 89 db: 79)
12 QUERY XML Select * from Third where Third.tbl.CS228
    ="B" and Third.tbl.prob>0.03: (time: 65 db: 54)
13 QUERY XML Select * from Third where Third.tbl.CS298
    ="A" and Third.tbl.prob>0.06: (time: 70 db: 60)
14 QUERY XML Select * from Third where Third.tbl.CS281
    ="B" and Third.tbl.prob>0.08: (time: 70 db: 62)
15 QUERY XML Select * from Third where Third.tbl.CS183
    ="B" and Third.tbl.prob>0.09: (time: 69 db: 54)
16 QUERY XML Select * from Fourth where Fourth.tbl.
    CS136="A" and Fourth.tbl.prob>0.09: (time: 77 db:
    68)
```

- 17 QUERY XML Select * from Fourth where Fourth.tbl.
CS251="B" and Fourth.tbl.prob>0.09: (time: 70 db:
59)
- 18 QUERY XML Select * from Fourth where Fourth.tbl.
CS277="B" and Fourth.tbl.prob>0.09: (time: 73 db:
63)
- 19 QUERY XML Select * from Fourth where Fourth.tbl.
CS249="B" and Fourth.tbl.prob>0.09: (time: 62 db:
54)
- 20 QUERY XML Select * from Fourth where Fourth.tbl.
CS271="A" and Fourth.tbl.prob>0.09: (time: 90 db:
81)

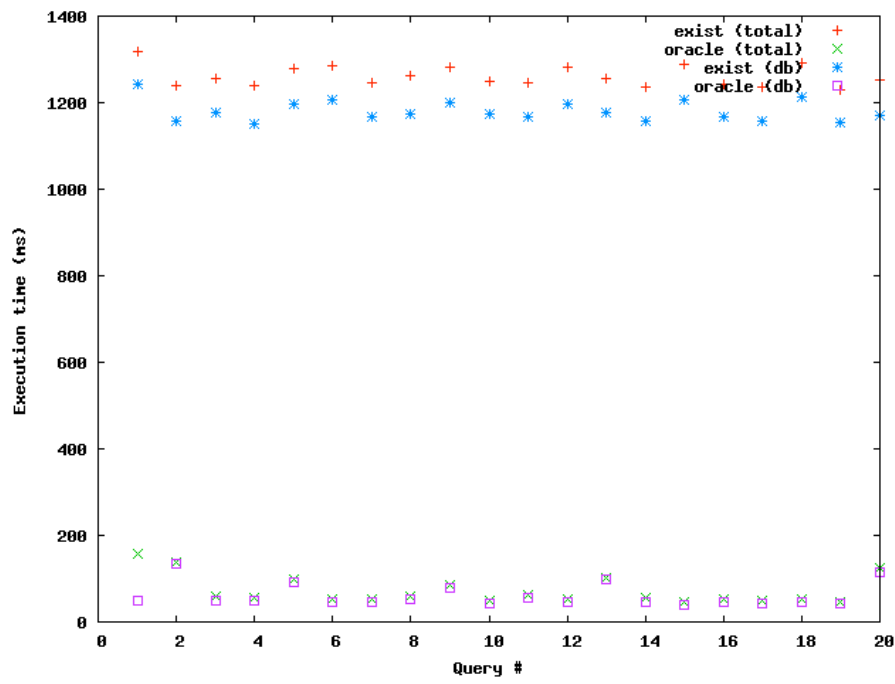
Exist Data

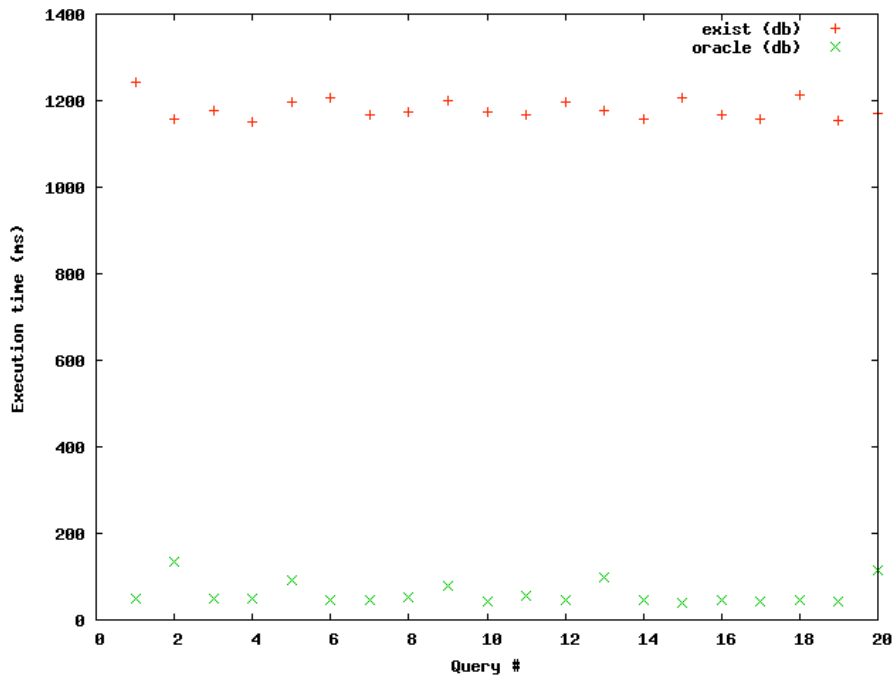
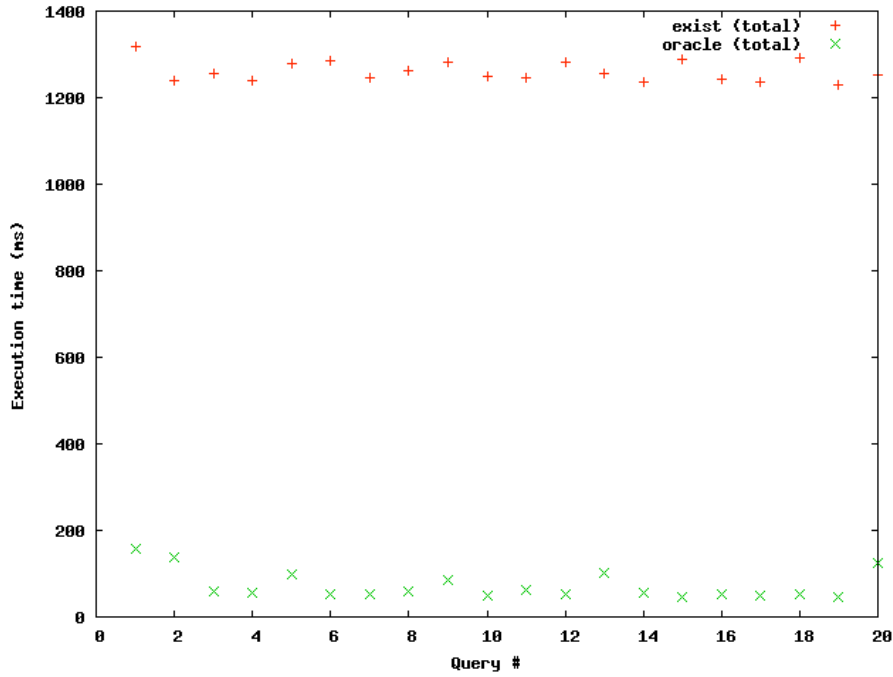
- 1 QUERY XML Select * from First where First.tbl.CS115
="A" and First.tbl.prob>0.01: (time: 645 db: 635)
- 2 QUERY XML Select * from First where First.tbl.CS242
="B" and First.tbl.prob>0.03: (time: 571 db: 564)
- 3 QUERY XML Select * from First where First.tbl.CS219
="A" and First.tbl.prob>0.06: (time: 844 db: 834)
- 4 QUERY XML Select * from First where First.tbl.CS211
="B" and First.tbl.prob>0.08: (time: 518 db: 512)
- 5 QUERY XML Select * from First where First.tbl.CS172
="B" and First.tbl.prob>0.09: (time: 511 db: 504)
- 6 QUERY XML Select * from Second where Second.tbl.
CS292="B" and Second.tbl.prob>0.01: (time: 622 db
: 615)
- 7 QUERY XML Select * from Second where Second.tbl.
CS107="A" and Second.tbl.prob>0.03: (time: 518 db
: 511)
- 8 QUERY XML Select * from Second where Second.tbl.
CS244="B" and Second.tbl.prob>0.06: (time: 517 db
: 510)
- 9 QUERY XML Select * from Second where Second.tbl.
CS279="B" and Second.tbl.prob>0.08: (time: 509 db
: 502)
- 10 QUERY XML Select * from Second where Second.tbl.
CS183="B" and Second.tbl.prob>0.09: (time: 514 db

- : 507)
- 11 QUERY XML Select * from Third where Third.tbl.CS277="A" and Third.tbl.prob>0.01: (time: 513 db: 507)
 - 12 QUERY XML Select * from Third where Third.tbl.CS228="B" and Third.tbl.prob>0.03: (time: 510 db: 503)
 - 13 QUERY XML Select * from Third where Third.tbl.CS298="A" and Third.tbl.prob>0.06: (time: 514 db: 506)
 - 14 QUERY XML Select * from Third where Third.tbl.CS281="B" and Third.tbl.prob>0.08: (time: 518 db: 509)
 - 15 QUERY XML Select * from Third where Third.tbl.CS183="B" and Third.tbl.prob>0.09: (time: 512 db: 504)
 - 16 QUERY XML Select * from Fourth where Fourth.tbl.CS136="A" and Fourth.tbl.prob>0.09: (time: 516 db: 509)
 - 17 QUERY XML Select * from Fourth where Fourth.tbl.CS251="B" and Fourth.tbl.prob>0.09: (time: 516 db: 507)
 - 18 QUERY XML Select * from Fourth where Fourth.tbl.CS277="B" and Fourth.tbl.prob>0.09: (time: 509 db: 502)
 - 19 QUERY XML Select * from Fourth where Fourth.tbl.CS249="B" and Fourth.tbl.prob>0.09: (time: 514 db: 508)
 - 20 QUERY XML Select * from Fourth where Fourth.tbl.CS271="A" and Fourth.tbl.prob>0.09: (time: 526 db: 519)

B.2.13 Complex Project Conditions

Comparison





Oracle Data

```
1 QUERY XML Select First.cnt.college,First.cnd.CS101
   from First: (time: 159 db: 50)
2 QUERY XML Select First.cnt.comments,First.cnd.CS106
   from First: (time: 139 db: 133)
3 QUERY XML Select First.cnt.year,First.cnd.CS143 from
   First: (time: 58 db: 50)
4 QUERY XML Select First.cnt.semester,First.cnd.CS103
   from First: (time: 55 db: 49)
5 QUERY XML Select First.cnt.instructor,First.cnd.
   CS121 from First: (time: 100 db: 93)
6 QUERY XML Select Second.cnt.semester,Second.cnd.
   CS143 from Second: (time: 53 db: 46)
7 QUERY XML Select Second.cnt.major,Second.cnd.CS146
   from Second: (time: 53 db: 46)
8 QUERY XML Select Second.cnt.major,Second.cnd.CS135
   from Second: (time: 58 db: 51)
9 QUERY XML Select Second.cnt.year,Second.cnd.CS102
   from Second: (time: 85 db: 78)
10 QUERY XML Select Second.cnt.semester,Second.cnd.
   CS127 from Second: (time: 50 db: 44)
11 QUERY XML Select Third.cnt.major,Third.cnd.CS144
   from Third: (time: 63 db: 55)
12 QUERY XML Select Third.cnt.instructor,Third.cnd.
   CS138 from Third: (time: 54 db: 47)
13 QUERY XML Select Third.cnt.semester,Third.cnd.CS112
   from Third: (time: 102 db: 97)
14 QUERY XML Select Third.cnt.comments,Third.cnd.CS140
   from Third: (time: 55 db: 46)
15 QUERY XML Select Third.cnt.comments,Third.cnd.CS122
   from Third: (time: 47 db: 40)
16 QUERY XML Select Fourth.cnt.comments,Fourth.cnd.
   CS117 from Fourth: (time: 54 db: 47)
17 QUERY XML Select Fourth.cnt.major,Fourth.cnd.CS124
   from Fourth: (time: 48 db: 42)
18 QUERY XML Select Fourth.cnt.major,Fourth.cnd.CS107
   from Fourth: (time: 51 db: 46)
19 QUERY XML Select Fourth.cnt.major,Fourth.cnd.CS128
   from Fourth: (time: 46 db: 41)
```

20 QUERY XML Select Fourth.cnt.semester,Fourth.cnd.
CS130 from Fourth: (time: 124 db: 116)

Exist Data

- 1 QUERY XML Select First.cnt.college,First.cnd.CS101
from First: (time: 1319 db: 1241)
- 2 QUERY XML Select First.cnt.comments,First.cnd.CS106
from First: (time: 1239 db: 1158)
- 3 QUERY XML Select First.cnt.year,First.cnd.CS143 from
First: (time: 1255 db: 1177)
- 4 QUERY XML Select First.cnt.semester,First.cnd.CS103
from First: (time: 1238 db: 1151)
- 5 QUERY XML Select First.cnt.instructor,First.cnd.
CS121 from First: (time: 1280 db: 1198)
- 6 QUERY XML Select Second.cnt.semester,Second.cnd.
CS143 from Second: (time: 1284 db: 1205)
- 7 QUERY XML Select Second.cnt.major,Second.cnd.CS146
from Second: (time: 1245 db: 1167)
- 8 QUERY XML Select Second.cnt.major,Second.cnd.CS135
from Second: (time: 1262 db: 1175)
- 9 QUERY XML Select Second.cnt.year,Second.cnd.CS102
from Second: (time: 1283 db: 1199)
- 10 QUERY XML Select Second.cnt.semester,Second.cnd.
CS127 from Second: (time: 1250 db: 1174)
- 11 QUERY XML Select Third.cnt.major,Third.cnd.CS144
from Third: (time: 1246 db: 1167)
- 12 QUERY XML Select Third.cnt.instructor,Third.cnd.
CS138 from Third: (time: 1282 db: 1196)
- 13 QUERY XML Select Third.cnt.semester,Third.cnd.CS112
from Third: (time: 1255 db: 1176)
- 14 QUERY XML Select Third.cnt.comments,Third.cnd.CS140
from Third: (time: 1237 db: 1156)
- 15 QUERY XML Select Third.cnt.comments,Third.cnd.CS122
from Third: (time: 1288 db: 1206)
- 16 QUERY XML Select Fourth.cnt.comments,Fourth.cnd.
CS117 from Fourth: (time: 1243 db: 1166)
- 17 QUERY XML Select Fourth.cnt.major,Fourth.cnd.CS124
from Fourth: (time: 1237 db: 1156)

- 18 QUERY XML Select Fourth.cnt.major,Fourth.cnd.CS107
from Fourth: (time: 1293 db: 1213)
- 19 QUERY XML Select Fourth.cnt.major,Fourth.cnd.CS128
from Fourth: (time: 1229 db: 1153)
- 20 QUERY XML Select Fourth.cnt.semester,Fourth.cnd.
CS130 from Fourth: (time: 1253 db: 1169)

C SPO Schema Document

This is the Structured Probability Object XSD, used to validate XML describing an SPO. The original document is available at <http://www.csr.uky.edu/wtw/schema/spo.xsd>. This document was *not* created as part of this thesis.

```

1 <?xml version="1.0"?>
2 <xsd:schema xmlns:xsd="http://www.w3.org/2001/
   XMLSchema"
3   targetNamespace="http://www.csr.uky.edu/wtw/schema
   "
4   elementFormDefault="qualified"
5   xmlns:spo="http://www.csr.uky.edu/wtw/schema">
6 <!-- xmlns:xsi="http://www.w3.org/1999/XMLSchema-
   instance" -->
7
8
9 <xsd:annotation><xsd:documentation>
10 This is the definition for SPO.
11 </xsd:documentation></xsd:annotation>
12 <!--
13 xmlns:xsd="http://www.w3.org/2000/10/XMLSchema"
14 <xsd:simpleType name="nameType">
15     <xsd:restriction base = "xsd:string">
16     </xsd:restriction>
17     <xsd:attribute name="ID" type="xsd:ID" use="
   optional"/>
18 </xsd:simpleType>
19 -->
20
21 <xsd:complexType name="nameType">
22     <xsd:simpleContent>
23         <xsd:extension base="xsd:string">
24             <!-- <xsd:attribute name="ID" type="xsd:
   NCName" use="optional"/> -->
25         </xsd:extension>
26     </xsd:simpleContent>

```

```
27 </xsd:complexType>
28
29 <xsd:simpleType name="valType">
30     <xsd:restriction base = "xsd:string">
31 <!--
32         <xsd:pattern value="[A-Z|0-9]+"\>
33 -->
34     </xsd:restriction>
35 </xsd:simpleType>
36
37 <xsd:complexType name="elemType">
38     <xsd:sequence>
39         <xsd:element name="name" minOccurs="1" maxOccurs
40             ="1" type="spo:nameType"/>
41         <xsd:element name="val" minOccurs="1" maxOccurs
42             ="unbounded" type="spo:valType"/>
43     </xsd:sequence>
44     <xsd:attribute name="IDREF" type="xsd:NCName" use
45         ="optional"/>
46 </xsd:complexType>
47
48 <xsd:complexType name="contextType">
49     <xsd:sequence>
50         <xsd:element name="elem" minOccurs="0" maxOccurs
51             ="unbounded" type="spo:elemType"/>
52     </xsd:sequence>
53 </xsd:complexType>
54
55 <xsd:complexType name="variableType">
56     <xsd:sequence>
57         <xsd:element name="name" minOccurs="1" maxOccurs
58             ="unbounded" type="spo:nameType"/>
59     </xsd:sequence>
60 </xsd:complexType>
61
62 <xsd:simpleType name="probabilityType">
63     <xsd:restriction base="xsd:float">
64         <xsd:minInclusive value="0.0000000"/>
65     </xsd:restriction>
66 </xsd:simpleType>
```

```
60     <xsd:maxInclusive value="1.0000000"/>
61   </xsd:restriction>
62 </xsd:simpleType>
63
64 <xsd:complexType name="rowType">
65   <xsd:sequence>
66     <xsd:element name="val" minOccurs="1" maxOccurs
        ="unbounded" type="spo:valType"/>
67     <xsd:element name="P" type="spo:probabilityType
        "/>
68   </xsd:sequence>
69 </xsd:complexType>
70
71 <xsd:complexType name="tableType">
72   <xsd:sequence>
73     <xsd:element name="variable" type="spo:
        variableType"/>
74     <xsd:element name="row" minOccurs="1"
        maxOccurs="unbounded" type="spo:rowType
        "/>
75   </xsd:sequence>
76 </xsd:complexType>
77
78 <xsd:complexType name="conditionalType">
79   <xsd:sequence>
80     <xsd:element name="elem" minOccurs="0" maxOccurs
        ="unbounded" type="spo:elemType"/>
81   </xsd:sequence>
82 </xsd:complexType>
83
84 <xsd:simpleType name="pathType">
85   <xsd:restriction base="xsd:string"/>
86 </xsd:simpleType>
87
88 <xsd:complexType name="spoType">
89   <xsd:sequence>
90     <xsd:element name="context" type="spo:
        contextType"/>
```



```
91     <xsd:element name="table" type="spo:tableType
92         "/>
93     <xsd:element name="conditional" type="spo:
94         conditionalType"/>
95     </xsd:sequence>
96     <xsd:attribute name="path" type="spo:pathType"
97         use="required"/>
98 </xsd:complexType>
99
100 <xsd:complexType name="sposType">
101     <xsd:sequence>
102     <xsd:element name="spo" minOccurs="1" maxOccurs
103         ="unbounded" type="spo:spoType">
104         <xsd:key name="myId">
105             <xsd:selector xpath="./table/variable"/>
106             <xsd:field xpath="name"/>
107         </xsd:key>
108         <xsd:keyref name="myIdref" refer="spo:myId">
109             <!-- Note spo:myId -->
110             <xsd:selector xpath="./context/ele"/>
111             <xsd:field xpath="@IDREF"/>
112         </xsd:keyref>
113     </xsd:element>
114     </xsd:sequence>
115 </xsd:complexType>
116
117 <xsd:element name="spos" type="spo:sposType">
118 </xsd:element>
119
120 </xsd:schema>
```