

# Data Integration under the Schema Tuple Query Assumption

Daniel Hedlund  
*c97dhd@cs.umu.se*

Version 1.1

November 7, 2003

This page intentionally contains only this sentence.

## **Abstract**

Typically data integration systems have significant gaps of coverage over the global (or mediated) schema they purport to cover. Given this reality, users are interested in knowing exactly which part of their query is supported by the available data sources. This thesis introduces a set of assumptions which enable users to obtain intensional descriptions of the certain, uncertain and missing answers to their queries given the available data sources. The general assumption is that query and source descriptions are written as tuple relational queries which return only whole schema tuples as answers. More specifically, queries and source descriptions must be within an identified sub-class of these ‘schema tuple queries’ which is closed over syntactic query difference. Because this identified query class is decidable for satisfiability, query containment and equivalence are also decidable. Sidestepping the schema tuple query assumption, the identified query class is more expressive than conjunctive queries with negated subgoals. The ability to directly express members of the query class in standard SQL makes this work immediately applicable in a wide variety of contexts. Also, a software implementation of syntactic query difference operators and schema tuple query predicates is presented.

This page intentionally contains only this sentence.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Answering user queries . . . . .	2
1.2	Broadening the System Response . . . . .	3
1.3	Schema Tuple Queries . . . . .	6
1.4	Organization of this Thesis . . . . .	7
<b>2</b>	<b>The Languages <math>\mathcal{L}</math> and <math>\mathcal{Q}</math></b>	<b>9</b>
2.1	Foundations . . . . .	9
2.2	$\mathcal{L}$ and its Decidability . . . . .	11
2.3	The Language $\mathcal{Q}$ . . . . .	13
<b>3</b>	<b>Reasoning over <math>\mathcal{L}</math> and <math>\mathcal{Q}</math></b>	<b>15</b>
3.1	Syntactic Query Difference over $\mathcal{L}$ . . . . .	15
3.2	Closure of Intersection over $\mathcal{L}$ . . . . .	16
3.3	Closure of Difference and Intersection over $\mathcal{Q}$ . . . . .	17
3.4	Expressing Integrity Constraints . . . . .	18
3.5	Containment and Equivalence over $\mathcal{Q}$ . . . . .	19
<b>4</b>	<b>SQD in Data Ingegration</b>	<b>21</b>
4.1	Querying in the Idealized Case . . . . .	21
4.2	Querying $\mathcal{L}$ Data Sources . . . . .	24
4.3	Querying Collections of $\mathcal{L}_{pos}$ Sources . . . . .	25
4.4	Querying Collections of $\mathcal{Q}_{pos}$ Sources . . . . .	29

<b>5</b>	<b>An Explorative Implementation</b>	<b>31</b>
5.1	The Syntactic Query Differencing Engine . . . . .	31
5.1.1	Syntax . . . . .	31
5.1.2	User-level Global Variables . . . . .	33
5.1.3	User-level Functions . . . . .	33
5.1.4	Work Order . . . . .	34
<b>6</b>	<b>Discussion</b>	<b>35</b>
6.1	Related Work . . . . .	35
6.2	Future Directions . . . . .	38
6.2.1	More expressive queries and schemas . . . . .	38
6.2.2	Annotating source advertisements . . . . .	39
6.2.3	Simplification and natural language generation . . . . .	39
<b>7</b>	<b>Conclusions</b>	<b>41</b>
<b>8</b>	<b>Acknowledgments</b>	<b>43</b>
	<b>Bibliography</b>	<b>44</b>
<b>A</b>	<b>Source Code</b>	<b>49</b>
A.1	/sqd/ . . . . .	49
A.1.1	load.lisp . . . . .	49
A.1.2	sqd-engine.lisp . . . . .	49
A.1.3	tuple-query-utils.lisp . . . . .	50
A.1.4	utilities.lisp . . . . .	51
A.2	/sqd/lsat/ . . . . .	60
A.2.1	load.lisp . . . . .	60
A.2.2	constants2fol.lisp . . . . .	60
A.2.3	fds2fol.lisp . . . . .	61
A.2.4	fol-sat.lisp . . . . .	62
A.2.5	fol-utils.lisp . . . . .	64
A.2.6	l2fol.lisp . . . . .	66

A.2.7	normal.lisp . . . . .	70
A.2.8	unify.lisp . . . . .	72

This page intentionally contains only this sentence.



# List of Figures

1.1	A typical data integration architecture . . . . .	1
1.2	An integrated movie system . . . . .	4
1.3	A query partitioning the universe of tuples . . . . .	5
4.1	The search tree for the example of chapter 1 . . . . .	26

This page intentionally contains only this sentence.

# Chapter 1

## Introduction

The need to answer queries over integrated databases has been a long explored topic. Normally the architecture governing such an effort is similar to that of figure 1.1.

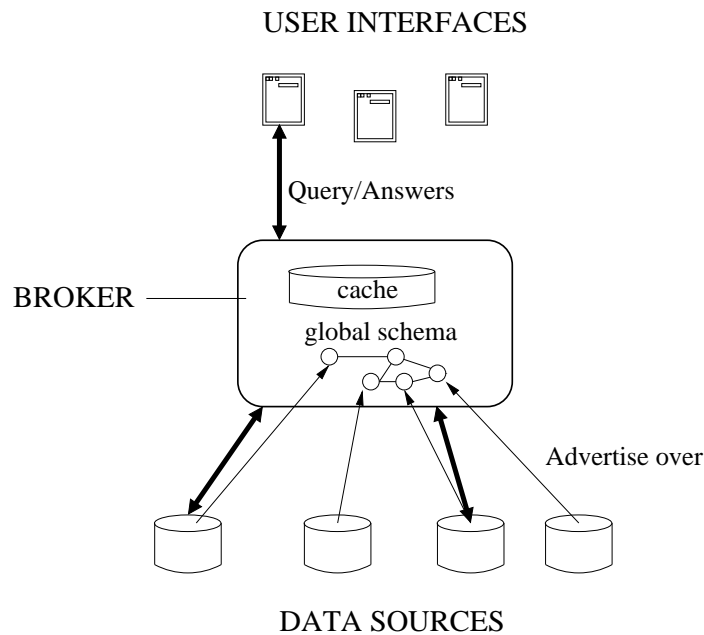


Figure 1.1: A typical data integration architecture

Here a *global schema* models the domain (e.g. movies, medical records, experimental findings, etc.) and autonomous *data sources* dynamically advertise their contents to a *broker* (or mediator) as views over the *global schema* (or mediated schema). A

user may then query the broker using the relations of the global schema and receive answers over a large number of separate data sources. This eliminates the burden of manually visiting each data source, launching queries over local schemas and combining answers into a globally coherent answer. Furthermore, a cache within the broker limits needless network transmission of commonly accessed tuples.

We assume that data sources advertise their contents as views over the global schema (known as *local-as-view*). This is in contrast to the less extensible method of defining global schema relations as views over source database relations (known as *global-as-view*). Under *local-as-view*, data sources may dynamically join the data integration system without necessitating a reorganization of the global schema. See [28][17] for further discussion.

The global schemas we shall consider are relational and are constrained under sets of functional dependencies. We also assume that the data sources provide correct and *locally complete* views. Locally complete sources contain all of the information described by their views. Thus one may invoke the closed world assumption[26] over the space defined by a source's view on the global schema. Finally we assume that the local sources provide globally standard attribute value encodings through their views.

## 1.1 Answering user queries

Given that source descriptions are views over the global schema, the user's query may be answered by rewriting it as a query using just the 'views' defining the sources. This relates to the general problem of answering queries using views. In the context of data integration systems, the following steps are normally carried out[17]:

- 1.) Obtain a set of *maximally contained rewritings* of the user's query over the 'views' defining the sources.
- 2.) Execute query rewritings to obtain a maximal set of answers from the sources.

Informally step 1 means that we rewrite the user's query in terms of the source descriptions such that the resulting query is logically contained within the user's

query. Furthermore such a contained rewriting must be maximal. That is there is no more general rewriting, with respect to a given query language, that is also logically contained by the user's query. When the data sources collectively cover the user's query, this reduces to finding an equivalent rewriting of the query using the data source descriptions. Normally, however, there will be gaps in coverage and we must obtain a set of maximally contained rewritings.

Step 2 considers how we actually obtain answers from the sources given the maximally contained rewritings. Certainly we may obtain tuples by simply executing the rewritings over the sources. However the notion of *certain* answers, introduced in [1], singles out a class of tuples that are of particular interest. Informally a tuple is a certain answer if it is an answer in all possible database instances (under the global schema) consistent with the contents of the given sources. See [1][17] for the formal definition of certain answers. As we shall see, the approach within this thesis eases the task of identifying certain answers considerably.

## 1.2 Broadening the System Response

Answers are generally presumed to be tuples. However because of gaps in coverage, users might be misled by a purely extensional response. Users might require an accompanying *intensional* response that puts the extensional answer within proper context<sup>1</sup>. Consider an example dialog over the data integration system pictured in figure 1.2.

**User:** give the films of the 90's that are available at video stores in Umeå.

Normally a system would be expected to present the certain set of tuples **A**. However, a broader response is more appropriate.

---

<sup>1</sup>The terms extensional and intensional denote two basic ways of describing a set of objects. The extensional description of a set is simply an enumeration of the items within the set. Thus the set  $A$  may be extensionally described as  $\{\text{'Alien'}, \text{'ET'}, \dots\}$ . An intensional description of a set describes the set conceptually, by summarizing the properties that members have. Thus an intensional description of set  $A$  may be "Science fiction films of the 1980's".

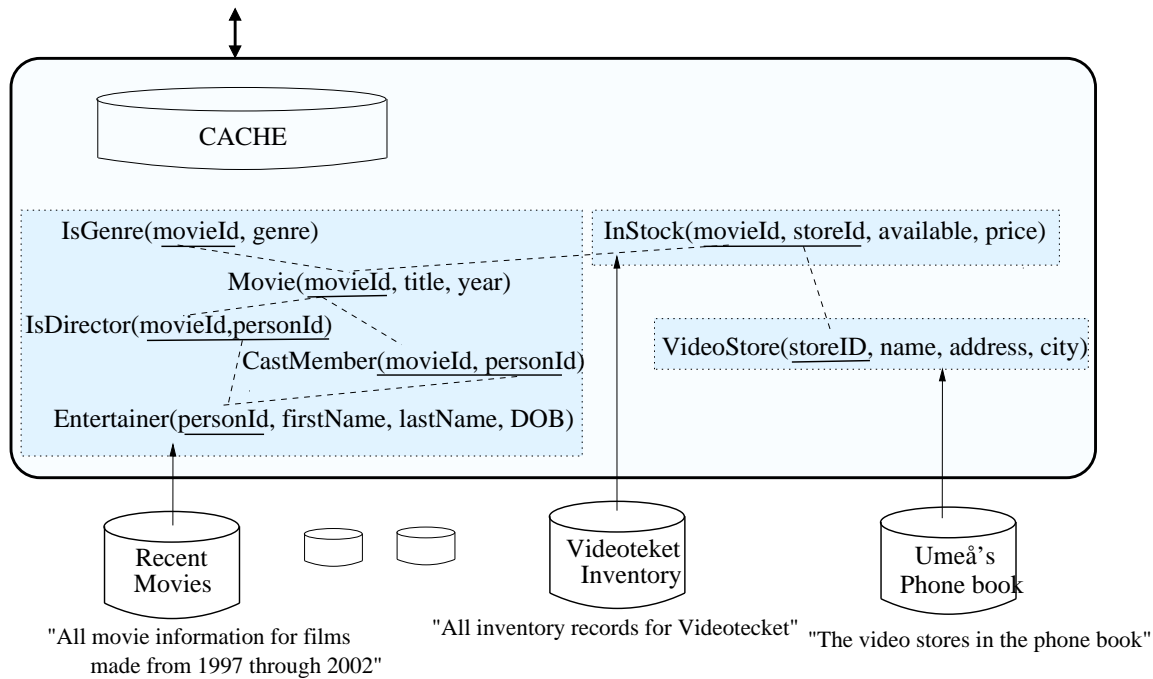


Figure 1.2: An integrated movie system

**System:** There are no sources that provide information on movies made before 1997. **A** is the set of movies made from 1997 through 1999 which are in the inventory of ‘Videoteket’. **B** is the set of movies made from 1997 through 1999 that are not at ‘Videoteket’, but may be available at video stores in Umeå that are not providing inventory information.

In this case **A** is the *certain* answer set and **B** is termed the *uncertain* answer set. The certain answer set, in correspondence with [1], is the set of tuples that are answers to the query no matter what additional data sources are added to the system. The uncertain answer set is the set of tuples that could be rendered in or out of the answer set if more information is made available. Both **A** and **B** are extensional, consisting of sets of movie tuples. The intensional description, “no sources that provide information on movies made before 1997,” is the *missing* answers description. This summarizes the set of movies for which no information whatsoever is provided. The intensional description, “the set of movies made from 1997 through 1999 that are not at ‘Videoteket’, but may be available at video stores in Umeå

that are not providing inventory information,” is the intensional description of the *uncertain answers*. Finally the description “the set of movies made from 1997 through 1999 which are in the inventory of ‘Videoteket’,” provides the *certain answer* set description.

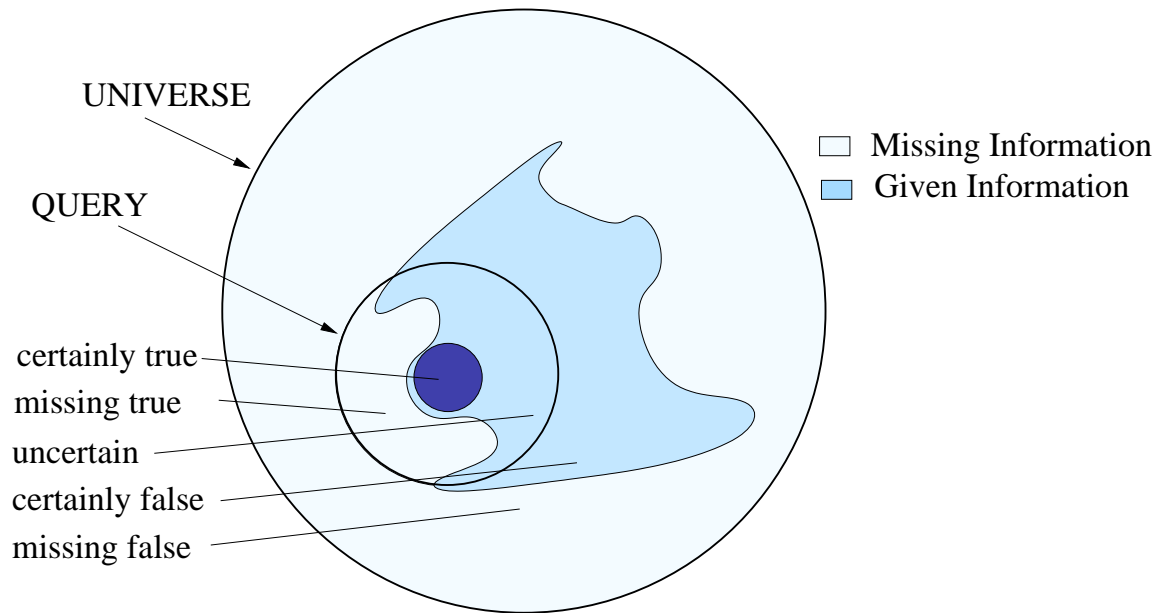


Figure 1.3: A query partitioning the universe of tuples

Figure 1.3 shows how the universe of tuples within a data integration environment is partitioned by a query. The five disjoint sets are the ‘certainly true’, ‘uncertain’, ‘missing true’, ‘certainly false’ and the ‘missing false’ tuples. We would focus on generating correct and complete descriptions of the first three of these sets and we shall refer to these three sets by the shorter names of the ‘certain’, ‘uncertain’ and ‘missing’ answers.

Thus, in summary, we propose a third step to the process of answering queries over data integration systems.

- 3.) Present intensional descriptions of the certain, uncertain and missing answers to a user’s query with respect to the available data sources.

### 1.3 Schema Tuple Queries

To cleanly represent the partitioning of the universe of tuples of figure 1.3, we restrict the language in which user queries and data source descriptions may be expressed. The type of queries (and views) we consider are the *schema tuple queries*<sup>2</sup>. Such queries are tuple relational queries[10] that return only whole tuples from relations of the schema, not arbitrary combinations of projected attributes. Thus, to obtain movies of the year 1999, we would write  $\{m \mid \text{Movie}(m) \wedge m.\text{year} = 1999\}$  rather than  $\{(m.\text{MovieId}, m.\text{Title}, m.\text{Year}) \mid \text{Movie}(m) \wedge m.\text{year} = 1999\}$ .

In this thesis schema tuple queries are further restricted by limiting their specification to the languages:  $\mathcal{L}$ ,  $\mathcal{L}_{pos}$ ,  $\mathcal{Q}$  and  $\mathcal{Q}_{pos}$ . Formulas in the language  $\mathcal{L}$  are signed quantifier sequences  $\exists^*$  over conjunctions of predicates over a single free tuple variable. Formulas in the language  $\mathcal{Q}$  are finite disjunctions of formulas within  $\mathcal{L}$ . The languages  $\mathcal{L}_{pos}$  and  $\mathcal{Q}_{pos}$  are corresponding languages which disallow negation of existential quantifier sequences.

The restriction to return only whole tuples from relations of the schema, not arbitrary combinations of projected attributes, is in contrast to the flexibility of normal tuple relational queries and relational algebra. However schema tuple queries specified in  $\mathcal{L}$  and  $\mathcal{Q}$  may have their set differences calculated syntactically. This enables one to perform set theoretic operations over query answer sets without having to materialize tuples. In addition the complexity of generating understandable intensional descriptions of schema tuple queries is much less formidable than with classical tuple relational calculus and relational algebra. Finally the limited flexibility of schema tuple queries may often be overcome by considering highly decomposed schemas.

Given that we may invoke the schema tuple query assumption, the main theoretical result of this thesis is to show that if the data sources are described using the language  $\mathcal{Q}_{pos}$  and the user's query is defined using the language  $\mathcal{L}$ , then one may: 1.) maximally rewrite queries using source descriptions; 2.) retrieve certain answers over the integrated system; 3.) generate intensional descriptions (within the language  $\mathcal{Q}$ ) of the certain, the uncertain and the missing answers to the user's query over the

---

<sup>2</sup>The term 'schema tuple query' is being newly introduced in this thesis. To the author's knowledge the notion is being newly introduced as well.



available data sources.

## 1.4 Organization of this Thesis

Chapter 2 formally defines the language  $\mathcal{L}$  and its closure over disjunction,  $\mathcal{Q}$ . Both  $\mathcal{L}$  and  $\mathcal{Q}$  are proven decidable for satisfiability; that is one may determine if there exists a database state for which the query will return a tuple. Chapter 3 shows that queries built using  $\mathcal{Q}$  are closed over syntactic query difference and intersection. Based on the decidability of  $\mathcal{Q}$ , we may thus decide query containment and equivalence over  $\mathcal{Q}$ . Chapter 4 gives an algorithm for how to use  $\mathcal{Q}$  within a data integration system to generate descriptions of the certain, uncertain and missing answers of the type in section 1.2. Chapter 5 gives an overview of the syntactic query differencing engine. Chapter 6 relates this work to prior work and gives future directions. Chapter 7 gives conclusions.

This page intentionally contains only this sentence.

# Chapter 2

## The Languages $\mathcal{L}$ and $\mathcal{Q}$

### 2.1 Foundations

We assume the existence of two disjoint, countable sets:  $\mathcal{U}$ , the *universal domain of atomic values*, and  $\mathcal{P}$ , *predicate names*. Let  $U$  be a distinct symbol representing the type of  $\mathcal{U}$ . A *relation schema*  $R$  is an  $n$ -tuple  $[U, \dots, U]$  where  $n \geq 1$  is called the *arity* of  $R$ . A *database schema*  $D$  is a sequence  $\langle P_1 : R_1, \dots, P_m : R_m \rangle$ , where  $m \geq 1$ ,  $P_i$ 's are distinct predicate names and  $R_i$ 's are relation schemas. A *relation instance*  $r$  of  $R$  with arity  $n$  is a finite subset of  $\mathcal{U}^n$ . A *database instance*  $d$  of  $D$  is a sequence  $\langle P_1 : r_1, \dots, P_m : r_m \rangle$ , where  $r_i$  is an instance of  $R_i$  for  $i \in [1..m]$ .

**Definition 1** (*Schema tuples*)

A *schema tuple*  $\tau$  of the database instance  $d$  is the pair  $\langle P_i : \mu \rangle$ , where  $1 \leq i \leq m$  and  $\mu \in r_i$ .

We say that the *type* of  $\tau$  is  $P_i$  and we say the components of  $\tau$  are the components of  $\mu$ . The schema tuple  $\tau_1$  is equal to the schema tuple  $\tau_2$  if and only if  $\tau_1$  and  $\tau_2$  match on type and they match on all components. Thus if  $\tau_1 = \langle \text{IsDirector} : ['367434', '43252'] \rangle$  and  $\tau_2 = \langle \text{CastMember} : ['367434', '43252'] \rangle$  then  $\tau_1 \neq \tau_2$ . The positional access operator is extended to the schema tuples to mirror the standard tuple relational calculus. Thus  $\tau_1[2] = '43252'$ . Furthermore we shall assume that tuple components may be accessed through attribute names (e.g.  $\tau_1.\text{personId}$ ).

Finally we shall assume that  $\mathcal{U}$  is totally ordered so that arithmetic comparison operators ( $=, >, <, \geq, \leq$  and  $\neq$ ) are well defined.

We now recursively define the set of *tuple relational formulas*. *Atomic formulas* provide the base for the inductive definition. The atomic formula  $P(z)$ , where  $P$  is a predicate name and  $z$  is a tuple variable, means that the tuple referred to by  $z$  is a schema tuple of type  $P$ . We term such formulas *range conditions*.  $X\theta Y$  is an atomic formula where  $X$  and  $Y$  are either constants or component references (of the form  $z.a$ ) and  $\theta$  is one of the arithmetic comparison operators. We term such formulas to be *simple conditions* if either  $X$  or  $Y$  are constants and to be *join conditions* if both  $X$  and  $Y$  are component references. Lastly we include atomic formula  $X\in C$  where  $X$  is a component reference,  $C$  is a set of constants and  $\in$  is a set membership operator ( $\in$  and  $\notin$ ). We term such formulas *set conditions*. Finally if  $F_1$  and  $F_2$  are tuple relational formula, where  $F_1$  has some free variable  $z$ , then  $F_1 \wedge F_2$ ,  $F_1 \vee F_2$ ,  $\neg F_1$ ,  $(\exists z)F_1$  and  $(\forall z)F_1$  are also tuple relational formulas.

We now define the notion of a schema tuple query.

**Definition 2** (*Schema tuple queries*)

A *schema tuple query* is an expression of the form  $\{x|\varphi\}$ , where  $\varphi$  is a tuple relational formula over the single free tuple variable  $x$ .

When we write the expression  $\{x|\varphi\}$  we normally assume that the expression  $\varphi$  is over the free variable  $x$ . For the schema tuple  $\tau$ ,  $\tau \in \{x|\varphi\}$  iff  $\{\tau/x\}(\varphi)$  where  $\{t/x\}\varphi$  means to substitute the term  $t$  in place of  $x$  in  $\varphi$ . Thus a *schema tuple query* is simply the intensional description of a schema tuple set. The actual tuples within this set are the extensional answers to the query.

We shall now turn our attention to several interesting sub-languages of the tuple relational formula that may be used to specify safe schema tuple queries.

## 2.2 $\mathcal{L}$ and its Decidability

A *basic query component* built over the free tuple variable  $x$  is a formula of the form  $(\exists y_1)(\exists y_2)\dots(\exists y_n)\Psi$ , where the  $\Psi$  is a conjunction of range conditions, simple conditions, set conditions and join conditions using exactly the tuple variables  $x, y_1, \dots, y_n$ . A *signed basic query component* is either a basic query component or the formula  $\neg\Theta$  where  $\Theta$  is a basic query component.

**Definition 3** (*The language  $\mathcal{L}$* )

The language  $\mathcal{L}$  consists of all formulas of the form:

$$P(x) \wedge (\bigwedge_{i=1}^k \Phi_i)$$

where  $P(x)$  restricts the free tuple variable  $x$  to range over  $P$  and  $\Phi_i$  is a signed basic query component over  $x$ .

The following two queries are built using  $\mathcal{L}$ :

- An SPJ query: “The movies directed by Lucas”

$$\begin{aligned} \{m_1|\ell_1\} = \{m_1| & Movie(m_1) \wedge \\ & (\exists y_1)(\exists y_2)( \\ & \quad IsDirector(y_1) \wedge Entertainer(y_2) \wedge \\ & \quad y_2.lastName = 'Lucas' \wedge \\ & \quad m_1.movieId = y_1.movieId \wedge \\ & \quad y_1.personId = y_2.personId) \} \end{aligned}$$

- A query with negation: “The films made in the year 2000 that are not in stock in a video store in the city Umeå”

$$\begin{aligned} \{m_2|\ell_2\} = \{m_2| & Movie(m_2) \wedge \\ & \neg(\exists y_3)(\exists y_4)( \\ & \quad VideoStore(y_3) \wedge InStock(y_4) \wedge \\ & \quad y_3.city = 'Umeå') \} \end{aligned}$$

$$\begin{aligned} & m_2.movieId = y_4.movieId \wedge \\ & y_4.storeId = y_3.storeId) \wedge \\ & m_2.year = 2000 \} \end{aligned}$$

Note that  $\mathcal{L}$  does not allow for any mixed quantifier sequences within basic query components. Thus a query such as “The movies that are not in stock in some video store in Umeå” is not expressible using  $\mathcal{L}$ .

We shall also identify the sub-language  $\mathcal{L}_{pos} \subset \mathcal{L}$  where no quantifier sequences are negated ( $s_i$  is positive for  $i$ ,  $1 \leq i \leq k$ ). When  $\ell \in \mathcal{L}_{pos}$  is expressed using a single quantifier sequence ( $k = 1$ ), we say the query is in *outward* form. When the number of variables is minimized, but  $k$  is maximized, we say that the formula is in *inward* form.

We now arrive at the main result of this section.

**Theorem 1** ( *$\mathcal{L}$  is decidable for satisfiability*)

*For all  $\ell \in \mathcal{L}$  over the free variable  $x$ , we may determine if there exists a database instance  $d$  where for some schema tuple  $\tau$ ,  $\{\tau/x\}\ell$ .*

Proof:

Because the test will be for satisfiability, the free variable of  $\ell$  is considered existentially quantified:  $\ell' = \exists x\ell$ . The formula  $\ell'$  may then be converted into an equivalent closed formula written in domain calculus: Range conditions are expanded into  $n$ -ary relations with introduced variables for each attribute of the corresponding relation. Tuple variable quantifications are replaced with corresponding sets of domain variable quantifications. Set conditions are re-written as disjunctions of simple conditions, and then all simple conditions are rewritten in their domain calculus equivalent form. All join conditions are also written in domain calculus using built in predicates such as  $> (X, Y)$  or  $= (X, Y)$ . This transformation builds a logically equivalent expression to  $\ell'$  within domain calculus.

The domain calculus expression then may then be converted into conjunctive normal form, CNF, in the standard way. All variables are existentially quantified in the input expression, but those that have a negation proceeding them, are converted to universal quantifiers in the standard way:  $\neg(\exists y_1)\dots(\exists y_m)(\psi)$  is equivalent to

$(\forall y_1) \dots (\forall y_m) \neg(\psi)$ . Thus all quantifier prefixes over variables are of  $\exists^* \forall^*$  form. Thus the Skolemization process will only introduce Skolem *constants* for variables.

The final function-free CNF expression may then be fed to a resolution theorem prover. Because the CNF has a finite Herbrand Universe, resolution will saturate after a finite number of steps. The CNF and hence  $\ell'$  is unsatisfiable if and only if the empty clause is generated. This has been implemented, and is briefly described in chapter 5.  $\square$

## 2.3 The Language $\mathcal{Q}$

We define the language  $\mathcal{Q}$  as the set of formula which are disjunctions of formula within  $\mathcal{L}$ . Not surprisingly  $\mathcal{Q}$  is decidable for satisfiability as well.

**Definition 4** (*The language  $\mathcal{Q}$* )

$q \in \mathcal{Q}$  if  $q$  is written as a finite expression  $\ell_1 \vee \dots \vee \ell_k$  where each  $\ell \in \mathcal{L}$  and  $\ell$  is over the free tuple variable  $x$ .

We also identify  $\mathcal{Q}_{pos} \subset \mathcal{Q}$  where every formula in the disjunction is in  $\mathcal{L}_{pos}$ .

**Theorem 2** ( *$\mathcal{Q}$  is decidable for satisfiability*)

*For all  $q \in \mathcal{Q}$  over the free variable  $x$ , we may determine if there exists a database instance  $d$  where for some schema tuple  $\tau$ ,  $\{\tau/x\}q$ .*

Proof:

A satisfiability test may be conducted for each disjunct of  $q \in \mathcal{Q}$ . Based on theorem 1 each such test is decidable. Since the number of disjuncts is finite, this decision process will terminate with a correct answer.  $\square$

This page intentionally contains only this sentence.



# Chapter 3

## Reasoning over $\mathcal{L}$ and $\mathcal{Q}$

We now cover several important properties that hold for queries built over the languages  $\mathcal{L}$  and  $\mathcal{Q}$ .

### 3.1 Syntactic Query Difference over $\mathcal{L}$

This theorem states that we may describe the set difference of two queries built over  $\mathcal{L}$  with a query built over  $\mathcal{Q}$ .

**Theorem 3** (*Syntactic query difference over  $\mathcal{L}$  is in  $\mathcal{Q}$* )

Let  $l_1 \in \mathcal{L}$  and  $l_2 \in \mathcal{L}$ . Then there is a  $q \in \mathcal{Q}$  with the property that for all database instances  $\{x_1|l_1\} - \{x_2|l_2\} = \{x_1|q\}$

Proof:

Assuming that  $l_1$  and  $l_2$  do not share any variable names in common,

$$\begin{aligned} \{x_1|l_1\} - \{x_2|l_2\} &= \{x_1|l_1 \wedge \{x_1/x_2\}\neg l_2\} = \\ \{x_1|l_1 \wedge \{x_1/x_2\}\neg(P(x_2) \wedge (\bigwedge_{i=1}^k \Phi_i))\} &= \\ \{x_1|(l_1 \wedge \neg P(x_1)) \vee & \\ (l_1 \wedge \neg \Phi_1) & \\ \vee \dots \vee & \\ (l_1 \wedge \neg \Phi_k)\} &= \end{aligned}$$

$$\{x_1 | \ell'_0 \vee \dots \vee \ell'_k\} = \\ \{x_1 | q\} \text{ where } q = \ell'_0 \vee \dots \vee \ell'_k. \square$$

The following illustrates syntactic query difference between the two example queries of section 2.2.

$$\{m_1 | \ell_1\} - \{m_2 | \ell_2\} = \{m_1 | \ell_1 \wedge \{m_1/m_2\} \neg \ell_2\} = \\ \{m_1 | \\ (Movie(m_1) \wedge \\ (\exists y_1)(\exists y_2)( \\ \quad IsDirector(y_1) \wedge Entertainer(y_2) \wedge \\ \quad y_2.lastName = 'Lucas' \wedge \\ \quad m_1.movieId = y_1.movieId \wedge \\ \quad y_1.personId = y_2.personId) \wedge \\ (\exists y_3)(\exists y_4) \\ \quad (VideoStore(y_3) \wedge InStock(y_4) \wedge \\ \quad y_3.city = 'Umeå' \wedge \\ \quad m_1.movieId = y_4.movieId \wedge \\ \quad y_4.storeId = y_3.storeId)) \\ \vee \\ (Movie(m_1) \wedge \\ (\exists y_1)(\exists y_2)( \\ \quad IsDirector(y_1) \wedge Entertainer(y_2) \wedge \\ \quad y_2.lastName = 'Lucas' \wedge \\ \quad m_1.movieId = y_1.movieId \wedge \\ \quad y_1.personId = y_2.personId) \wedge \\ \quad m_1.year \neq 2000)\}$$

## 3.2 Closure of Intersection over $\mathcal{L}$

The following, although obvious, is of use later in this thesis.

**Lemma 1** (*The closure of  $\mathcal{L}$  over intersection*)

Let  $\ell_1 \in \mathcal{L}$  and  $\ell_2 \in \mathcal{L}$ . Then there is an  $\ell_3 \in \mathcal{L}$  with the property that for all database instances  $\{x_1|\ell_1\} \cap \{x_2|\ell_2\} = \{x_1|\ell_3\}$

Proof:

By definition.  $\square$

### 3.3 Closure of Difference and Intersection over $\mathcal{Q}$

We now show that  $\mathcal{Q}$  stays closed over syntactic difference and intersection. This means that we may describe the set difference (or intersection) of two queries built over  $\mathcal{Q}$  with a third query built over  $\mathcal{Q}$ .

**Theorem 4** (*Syntactic query difference is closed over  $\mathcal{Q}$* )

Let  $q_1 \in \mathcal{Q}$  and  $q_2 \in \mathcal{Q}$ . Then there is a  $q_3 \in \mathcal{Q}$  such that for all database instances  $\{x_1|q_1\} - \{x_2|q_2\} = \{x_1|q_3\}$

Proof:

By induction on the number of disjuncts in  $q_2$ .

*Base case:* Assume that  $q_1$  consists of a finite number of disjuncts and  $q_2$  consists of  $n = 1$  disjunct. The result of set difference will simply be each disjunct of  $q_1$  minus the single disjunct of  $q_2$ . Because  $q_2 \in \mathcal{L}$ , this may be carried out for each disjunct of  $q_1$  as in theorem 3. The resulting formula, when all generated disjuncts are grouped together, is within  $\mathcal{Q}$ . Therefore the theorem holds in the base case.

*Induction hypothesis:* Assume the theorem holds when  $q_1$  consists of a finite number of disjuncts and the number of disjuncts in  $q_2$  is  $n = k$ .

*Induction step:* Consider  $q'_2$ , written as  $q_2 \vee \ell$  where  $\ell \in \mathcal{L}$  and  $q_2$  consists of  $k$  disjuncts.  $\{x_1|q_1\} - \{x_2|q'_2\} = \{x_1|q_1\} - (\{x_2|q_2\} \cup \{x_2|\ell\}) = (\{x_1|q_1\} - \{x_2|q_2\}) - \{x_2|\ell\}$ . From the induction hypothesis  $(\{x_1|q_1\} - \{x_2|q_2\} = \{x_1|q_3\})$  where  $q_3 \in \mathcal{Q}$ . From the base case  $\{x_1|q_3\} - \{x_2|\ell\} = \{x_1|q'_3\}$  where  $q'_3 \in \mathcal{Q}$ . Therefore  $\{x_1|q_1\} - \{x_2|q'_2\} = \{x_1|q'_3\}$  where  $q'_3 \in \mathcal{Q}$  for  $q'_2$  of  $k + 1$  disjuncts.

Thus the theorem has been proven by induction.  $\square$

**Theorem 5** (*Syntactic query intersection is closed over  $\mathcal{Q}$* )

Let  $q_1 \in \mathcal{Q}$  and  $q_2 \in \mathcal{Q}$ . Then there is a  $q_3 \in \mathcal{Q}$  such that for all database instances  $\{x_1|q_1\} \cap \{x_2|q_2\} = \{x_1|q_3\}$

Proof:

Since  $q_1 \wedge q_2 \equiv q_1 \wedge \neg(q_1 \wedge \neg q_2)$ , the truth of theorem 4 establishes the truth of this theorem.  $\square$

### 3.4 Expressing Integrity Constraints

Since functional dependencies constrain the set of legal database states, an expression that necessarily violates a functional dependency must be ruled to be unsatisfiable. We shall represent functional dependencies as universally quantified formulas[2]. In turn these formulas can well be included in the resolution process that decides satisfiability of queries.

**Definition 5** (*Functional dependencies as formulas*)

The functional dependency  $W \rightarrow V$  over the relation  $R$  where  $W$  is a set of  $m$  attributes and  $V$  is a set of  $n$  attributes is expressed as the universally quantified formula:

$$\begin{aligned} & (\forall x)(\forall y)(P(x) \wedge P(y) \wedge \\ & \quad y.w_1 = x.w_1 \wedge \dots \wedge y.w_m = x.w_m \Rightarrow \\ & \quad y.v_1 = x.v_1 \wedge \dots \wedge y.v_n = x.v_n) \end{aligned}$$

As an example, the functional dependency  $movieId \rightarrow title\ year$  on the relation  $Movie$  is represented by:

$$\begin{aligned} & (\forall m_1)(\forall m_2)(Movie(m_1) \wedge Movie(m_2) \wedge \\ & \quad m_1.movieId = m_2.movieId \Rightarrow \\ & \quad m_1.title = m_2.title \wedge m_1.year = m_2.year). \end{aligned}$$

Note also that other integrity constraints may be encoded as universally quantified formula. For example the constraint that all film years are greater than 1927 but less than or equal to 2002 could be encoded:

$$(\forall m_1)(Movie(m_1) \Rightarrow \\ m_1.year \geq 1927 \wedge m_1.year \leq 2002)$$

We shall use the formula  $\Gamma$  to signify all of the integrity constraints of the domain.  $\Gamma$  is simply a conjunction of universally quantified formulas and thus may be converted to CNF without Skolemization. When determining the satisfiability of  $\ell \in \mathcal{L}$ , the CNF representing  $\Gamma$  may simply be conjoined with the CNF representing  $\ell$  before we start the resolution procedure.

### 3.5 Containment and Equivalence over $\mathcal{Q}$

**Theorem 6** (*Decidability of  $\subseteq, =$  and disjointness over  $\mathcal{Q}$* )

*if  $q_1 \in \mathcal{Q}$ ,  $q_2 \in \mathcal{Q}$  and  $\Gamma$  expresses the integrity constraints over the domain, then there exists sound and complete inference mechanisms to decide if the three predicates:*

- 1.)  $\{x_1|q_1\} \subseteq \{x_2|q_2\}$
- 2.)  $\{x_1|q_1\} = \{x_2|q_2\}$
- 3.)  $\{x_1|q_1\} \cap \{x_2|q_2\} = \emptyset$ .

*are necessarily true over the set of all database instances for which  $\Gamma$  holds.*

Proof:

Predicate 1:  $\{x_1|q_1\} \subseteq \{x_2|q_2\}$  under the constraints  $\Gamma$  iff  $\{x_1|q_2\} - \{x_2|q_2\} = \emptyset$  under the constraints  $\Gamma$ . Based on the closure of  $\mathcal{Q}$  over difference this reduces to a satisfiability test enriched with  $\Gamma$  over a query built over  $\mathcal{Q}$ . Because the CNF representing  $\mathcal{Q}$  and  $\Gamma$  has a finite Herbrand universe, we may decide predicate 1. Predicate 2 follows directly from the ability to decide predicate 1. Predicate 3 is decided by deciding the satisfiability of  $\{x_1|q_2\} \cap \{x_2|q_2\}$ . Again this is decidable because the CNF representing  $\mathcal{Q}$  and  $\Gamma$  has a finite Herbrand universe.  $\square$

This page intentionally contains only this sentence.

# Chapter 4

## Syntactic Query Differencing in Data Integration

### 4.1 Querying in the Idealized Case

Let us return to the problem of answering queries over data integration systems. We shall assume that the user's query is  $\{x|\ell\}$  where  $\ell \in \mathcal{L}$  over the free variable  $x$ . Furthermore we shall assume that there exist  $n$  source databases where the  $i$ -th source database is represented by the 'view expression'  $\{x_{s_i}|desc(s_i)\}$  where  $desc(s_i) \in \mathcal{Q}$ . This 'view expression' describes the *complete* information that the source  $s_i$  has over the global schema. Because we assume that the information in the source is complete, we may invoke the closed world assumption over the space of tuples satisfying  $desc(s_i)$ . In addition the broker holds a universally quantified formula  $\Gamma$  which states the functional dependencies and other constraints of the domain. Finally we assume that the broker maintains a cache, which may be described based on the history of queries issued to sources.

In the example of figure 1.2 the three data sources have advertised the following 'views' over the global schema:

$s_1$ : **Recent Movies:**

$\{m|Movie(m)\wedge$

$$\begin{aligned}
& m.year \leq 2002 \wedge m.year \geq 1997 \} \cup \\
& \{g | IsGenre(g) \wedge (\exists m')(Movie(m') \wedge \\
& \quad m'.year \leq 2002 \wedge m'.year \geq 1997 \wedge \\
& \quad g.movieId = m'.movieId) \} \cup \\
& \{d | IsDirector(d) \wedge (\exists m')(Movie(m') \wedge \\
& \quad m'.year \leq 2002 \wedge m'.year \geq 1997 \wedge \\
& \quad d.movieId = m'.movieId) \} \cup \\
& \{c | CastMember(c) \wedge (\exists m')(Movie(m') \wedge \\
& \quad m'.year \leq 2002 \wedge m'.year \geq 1997 \wedge \\
& \quad c.movieId = m'.movieId) \} \cup \\
& \{e | Entertainer(e) \wedge (\exists d')(\exists m') \\
& \quad (IsDirector(d') \wedge Movie(m') \wedge \\
& \quad m'.year \leq 2002 \wedge m'.year \geq 1997 \wedge \\
& \quad d'.movieId = m'.movieId \wedge \\
& \quad e.personId = d'.personId) \} \cup \\
& \{e | Entertainer(e) \wedge (\exists c')(\exists m') \\
& \quad (CastMember(c') \wedge Movie(m') \wedge \\
& \quad m'.year \leq 2002 \wedge m'.year \geq 1997 \wedge \\
& \quad c'.movieId = c'.movieId \wedge \\
& \quad e.personId = c'.personId) \}
\end{aligned}$$

$s_2$ : **Videoteket Inventory**:

$$\begin{aligned}
& \{s | InStock(s) \wedge (\exists v')(VideoStore(v') \wedge \\
& \quad s.StoreId = v'.StoreId \wedge \\
& \quad v'.Name = 'Videoteket' \wedge v'.city = 'Umeå' \}
\end{aligned}$$

$s_3$ : **Umeå's Phone book**: The video store listings

$$\{v | VideoStore(v) \wedge v.city = 'Umeå'\}$$

In the ideal case a data source may answer any query about a tuple in its possession. Thus, for example,  $s_2$  could answer queries such as, “give the inventory records



for horror movies of the year 2000.” As a more extreme example,  $s_1$  could answer queries such as, “give the cast members of horror films that are in the inventory of Videoteket.” Though the idealized case is rather unrealistic, it provides a good starting point for discussion.

The pseudo code below gives the certain answers to the user’s query under the ideal assumption. The method  $ask(source, query)$  returns the answers to the  $query$  posed over the  $source$ . Because data sources are assumed to be able to evaluate any condition over a tuple in their possession, such answers will be certain. The operation  $A \leftarrow B$  inserts the tuples of set  $B$  into the set  $A$ . All uses of  $-$  and  $\cap$  are syntactic operations over query expressions.

*Process(query)*

**begin**

$q = query$ ;

$ANSWERS = \emptyset$ ;

**for**  $i = 1$  **to**  $n$  **do**

$ANSWERS \leftarrow ask(s_i, q \cap \{x_{s_i} | desc(s_i)\})$ ;

$q = q - (q \cap \{x_{s_i} | desc(s_i)\})$ ;

**od**

$certain = query - q$ ;

$missing = q$ ;

**end**

$ANSWERS$  thus contains the tuples that satisfy the query over the space described by  $certain$ . Missing tuples are described with the expression  $missing$ .

Now we shall describe querying over three successively more complex scenarios where data sources do not have perfect knowledge about the tuples they possess. First we consider querying sources described with formula in  $\mathcal{L}$ . Second we consider the case of querying collections of  $\mathcal{L}$  sources. Unfortunately the results break down for data sources using the full  $\mathcal{L}$  language and we must limit the source descriptions to  $\mathcal{L}_{pos}$ . Finally we shall consider the full case of querying collections of  $\mathcal{Q}_{pos}$  sources.

## 4.2 Querying $\mathcal{L}$ Data Sources

We now consider how to obtain answers to  $\{x|\ell\}$  over a data source  $s$  described by  $\{x_s|desc(s)\}$  where  $desc(s) \in \mathcal{L}$ . Just as in the idealized case, the query is intersected with the data source description and the resulting query is sent to the data source. However, unlike the ideal case, a ‘real’ data source is not always able to evaluate the conditions of this intersected query. For example assume we issue the following query over the following source:

$\ell_1$ : **Give year 2000 horror movies:**

$$\{x|Movie(x) \wedge x.year = 2000 \wedge$$

$$(\exists y_1)(Genre(y_1) \wedge y_1.type = \text{‘horror’} \wedge$$

$$y_1.movieID = x.movieID)\}$$

$s'_1$ : **Just movie information from  $s_1$ :**

$$\{m|Movie(m) \wedge$$

$$m.year \geq 1997 \wedge m.year \leq 2002\}$$

Even though the source contains all of the tuples satisfying the query, it is not able to distinguish which of them are horror films. Thus the best the source may do is to drop the conditions it can not evaluate and return uncertain answers. The process is termed *query truncation* and it is carried out by simply *dropping* all non-free variables of the query formula  $\ell$  from the formula  $\ell \wedge desc(s)$ . Dropping a variable simply means that the variable and any corresponding range conditions, simple conditions, set conditions and join conditions within the query are removed. The result of this process is signified  $\{x|drop_\ell(\ell \wedge desc(s))\}$ . For the example above the query evaluated over  $s'_1$  is:

$$\{m|drop_{\ell_1}(\{m/x\}\ell_1 \wedge desc(s'_1))\} =$$

$$\{m|Movie(m) \wedge m.year = 2000\}.$$

All of the conditions in  $\{x_s|drop_\ell(\{x_s/x\}\ell \wedge desc(s))\}$  may be evaluated by the source  $s$ . Moreover, if  $drop_\ell(\{x_s/x\}\ell \wedge desc(s)) \equiv \{x_s/x\}\ell \wedge desc(s)$  then we know

that truncation did not generalize the query sent to the source and thus the answers over the source will be certain with respect to the original query. Alternatively if  $drop_\ell(\{x_s/x\}\ell \wedge desc(s)) \not\equiv \{x_s/x\}\ell \wedge desc(s)$  then the query was generalized and thus the answers to the query over the source will be uncertain. In either case  $\{x_s | drop_\ell(\{x_s/x\}\ell \wedge desc(s))\}$  is posed over the source to obtain answers.

As a more complex example consider the following query posed over the following source:

$\ell_2$ : **Horror or action films not available in Umeå:**

$$\begin{aligned} & \{x | Movie(x) \wedge \\ & (\exists g)(Genre(g) \wedge g.genre \in \{\text{'horror'}, \text{'action'}\} \wedge \\ & \quad g.movieId = m.movieId) \wedge \\ & \neg(\exists s)(\exists v)(InStock(s) \wedge VideoStore(v) \wedge \\ & \quad s.movieId = m.movieId \wedge \\ & \quad s.storeId = v.storeId \wedge v.city = \text{'Umeå'}) \} \end{aligned}$$

$s_4$ : **Horror films not shown anywhere**

$$\begin{aligned} & \{m | desc(s_4)\} \equiv \{m | Movie(m) \wedge \\ & (\exists g)(Genre(g) \wedge g.genre = \text{'horror'} \wedge \\ & \quad g.movieId = m.movieId) \wedge \\ & \neg(\exists s)(InStock(s) \wedge s.movieId = m.movieId) \} \end{aligned}$$

In this case one may verify that certain answers to the query are returned.

### 4.3 Querying Collections of $\mathcal{L}_{pos}$ Sources

Now we shall consider posing the query over a collection of  $n$  data sources, where the  $i$ -th source is described with  $\{x_i | desc(s_i)\}$  where  $desc(s_i) \in \mathcal{L}_{pos}$ . As opposed to the simple case above, this time the  $n$  sources may be combined in an effort to provide certain answers<sup>1</sup>. Each combination is considered a *plan*, and a plan is a formula within  $\mathcal{L}$ .

<sup>1</sup>Note that  $n$  may indeed be 1 and the case will be different than that in section 4.2. Specifically section 4.2 did not propose a strategy for how a source answers self joining queries.

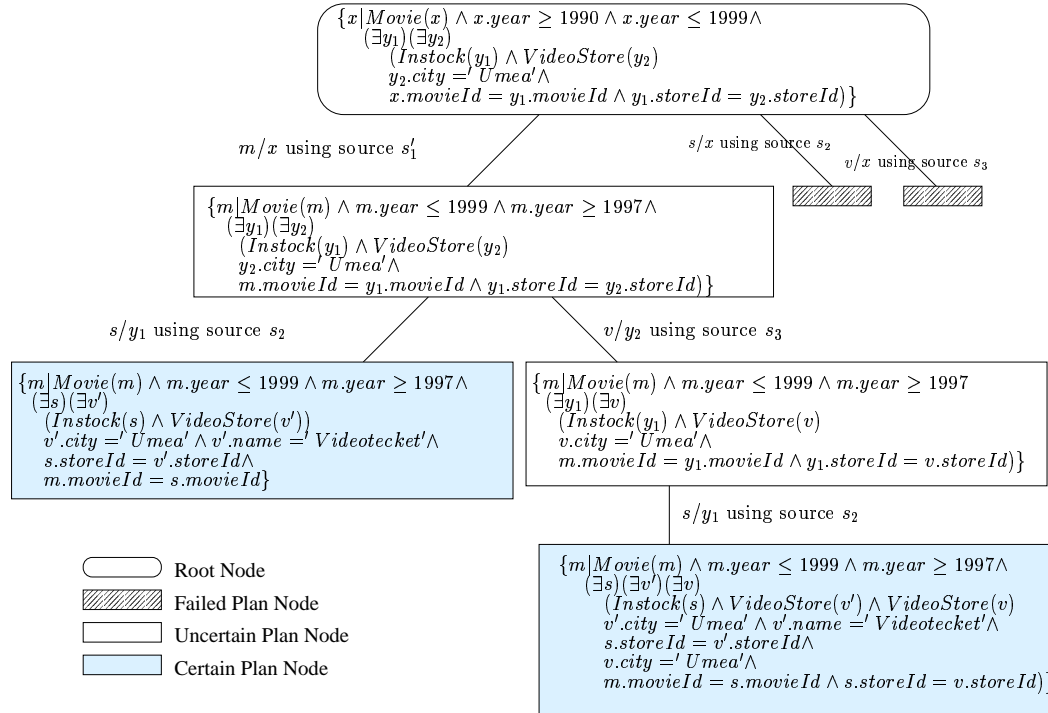


Figure 4.1: The search tree for the example of chapter 1

Our strategy shall be to search the space of potential plans looking for the set of the most general plans that return certain answers to the query based on the process defined in section 4.2. Although computationally intensive, the manner in which this is achieved is straightforward. A search tree is built where each node contains a plan. The original query formula  $\ell$  is the ‘plan’ within the root of the search tree. Nodes are expanded, by binding a source description formula to an unbound variable within the node’s plan. This process results in a finite tree of nodes, each a different plan to answer the query. This tree is traversed and the set of certain (and uncertain) answer generating plans are obtained.

These certain (and uncertain) answer generating plans are executed to retrieve needed tuples from the sources. Using the retrieved tuples and tuples already within the cache, certain and uncertain answer sets are identified. These answer sets, along with the certain and uncertain answer generating plans, provide the basis for a mixed extensional/intensional response of the type described in chapter 1.

We shall now describe the critical parts of this process in more detail. Figure 4.1 shall be referred to during this discussion. Figure 4.1 shows the search tree built for the query of chapter 1 over the sources  $s'_1$ ,  $s_2$  and  $s_3$ .  $s'_1$  is defined in section 4.2.

### The root node and its initial expansion

The root note contains the user query  $\ell$  as its plan with all of the variables within  $\ell$  marked *unbound*<sup>2</sup>. The root node is expanded to  $n$  depth level one nodes, each containing the plan  $\{x_s/x\}\ell \wedge desc(s_i)$ . By lemma 1, these plans are in  $\mathcal{L}$ . In figure 4.1 we can see that only the plan over  $s'_1$  is satisfiable.

### Expansion of non-root nodes

Now we show how search tree nodes of depth level one and beyond are expanded. A single unbound variable  $y$  within the plan of the node is matched with a single source  $s_i$ . This match is the basis of the expansion of the node. Assume that  $desc(s_i)$  is written in outward form and  $\ell'$  is the plan within the node to be expanded. The new node will have the plan  $\ell'' \in \mathcal{L}$  where  $\ell''$  is  $\ell'$  with the variable  $y$  replaced by the free variable of  $desc(s_i)$ , all of the existential variables of  $desc(s_i)$  placed within the quantifier sequence under which  $y$  had appeared, and the single conjunct of  $desc(s_i)$  conjoined with the conjunct over which  $y$  had ranged.

Since all combinations are considered, the branching factor of the tree is on the order of  $O(nm)$  where  $n$  is the number of sources and  $m$  is the number of variables in  $\ell$ . Figure 4.1 shows only the matches that make type sense.

### Termination tests

There are three conditions that signal that a search tree node does not need to be expanded:

- 1.) The node's plan is unsatisfiable.

---

<sup>2</sup>Figure 4.1 uses the variable names  $x$  and  $y_i$  to signify unbounded variables.

- 2.) The node's plan, with unbound variables dropped, could be used to give certain answers to the query according to the procedure of section 4.2.
- 3.) There are no remaining unbound variables within the node's plan.

If a node meets the first condition, it is marked *failed*. If a node meets the second condition, the unbound variables are dropped from its plan and it is marked *certain*. The third condition guarantees termination by bounding the size of the search tree to  $O((nm)^m)$ . In practice the size of the search tree is expected to be much smaller than this worst case.

### Obtaining certain, uncertain and missing answer descriptions

At the end of this process the tree is traversed to obtain descriptions of the certain, uncertain and missing answers. The certain answers are the disjunction of the plans obtained from the nodes that are marked certain. Thus  $certain \in \mathcal{Q}$ . The uncertain answers are the disjunction of the non-failed plans at level one minus the description of the certain answers. Thus  $uncertain \in \mathcal{Q}$ . Finally the missing answers is the original query minus both  $certain$  and  $uncertain$ . Thus  $missing \in \mathcal{Q}$ . From figure 4.1 one may verify that this results in the queries described in the example dialog of chapter 1.

### Issuing queries to the sources

Once the certain and uncertain answer set descriptions are obtained, the records of source/variable bindings are used to materialize the needed tuples from the sources. This is achieved by building a query for each variable/source binding where the variable is free and the conditions are the simple and type conditions over the variable within the plan. Such queries are prepared for dispatch to their corresponding sources.

In figure 4.1:  $\{m|Movie(m) \wedge m.year \leq 1999 \wedge m.year \geq 1999\}$  would be prepared for  $s'_1$ ,  $\{s|InStock(s)\}$  would be prepared for  $s_2$  and  $\{v|VideoStore(v)\}$  would be prepared for  $s_3$ . The reason that these queries are prepared, but not immediately sent, is because their answers, in part or whole, may already reside within the cache.

### Using and maintaining the cache

A cache over all the relations of the global schema is maintained within the broker. The description of what is within the cache is  $cache \in \mathcal{Q}_{pos}$ . Before the prepared queries of the previous subsection are issued to sources, the information within  $cache$  is syntactically subtracted from them. The remaining portions of the prepared queries are then sent to the sources and, upon return of answers, the cache description is augmented to reflect the additional tuples.

Once all of the source/variable bindings have been used to materialize relevant tuples into the cache, the certain answer set description is applied to obtain the extension of the certain answers.

The uncertain answer set description has all of its non free variables dropped and is likewise applied to the cache. This materializes all of the uncertain and certain tuples. From this the certain answers are set subtracted to yield an approximation of the uncertain answer extension. Note that this is the only place in the whole system where a classical set difference operation is applied over an actual extension. To obtain the exact extension of uncertain answers, the certain answers from the complement of the original query must be solved for and in turn subtracted from the approximation of the uncertain answer extension.

## 4.4 Querying Collections of $\mathcal{Q}_{pos}$ Sources

At a logical level, extending the system to the full example where we have a set of sources each described by an expression in  $\mathcal{Q}_{pos}$ , does not significantly alter the strategy of section 4.3. The  $\mathcal{Q}_{pos}$  expressions are merely pooled into a very large  $\mathcal{Q}_{pos}$  expression that fits the above case by treating each portion of a source as an ‘independent source’.

For performance however, one would like to insure that joins that could be performed at the sources are performed there and that only the relevant portions of the resulting answers are sent to the broker. In addition we might wish to enable a principled way to include semi-join optimizations into the framework. No doubt a

more sophisticated cache use and maintenance strategy will also be required. These considerations, however, are beyond the scope of this current work and would be the subject of future work.



# Chapter 5

## An Explorative Implementation

To better understand the possibilities and show the viability of the schema tuple query approach, we have also implemented the core functionality for syntactic reasoning over queries in the  $\mathcal{L}$  and  $\mathcal{Q}$  query languages.

### 5.1 The Syntactic Query Differencing Engine

We choose to develop the engine in Common LISP, because this language is mature, widely available, is very expressive, and at the same time allows for very rapid prototyping. Furthermore, using Common LISP allowed us to easily incorporate already existing automated reasoning code, and it also that the implementation will be easy to port and extend.

#### 5.1.1 Syntax

**q-query** = '(or ', {*l*-query}, ')';

**l-query** = '(, free variable, type, {signed query component} *or* range condition *or* simple condition *or* set condition, ')';

**signed query component** = '(, sign, '(, {existential variable}, ')', '(, {condition}, ')';

**condition** = range condition *or* simple condition *or* set condition *or* join condition;

**range condition** = '(' , type, existential variable, ')';

**simple condition** = '(' , operator, variable, role, value, ')';

**operator** = '=' *or* '<>' *or* '>' *or* '<' *or* '>=' *or* '<=';

**set condition** = '(' , set operator, variable, role, set, ')';

**set operator** = 'in' *or* 'not\_in';

**set** = '(' , {element}, ')';

**join condition** = '(' , operator, variable, role, variable, role, ')';

**sign** = '+' *or* '-';

**functional dependency** = '(' , type, lhs, rhs, ')';

**lhs** = '(' , {attribute}, ')';

**rhs** = '(' , {attribute}, ')';

**schema** = '(' , {relation}, ')';

**relation** = '(' , relation name, '(' , {attribute}, 2 \* ')';

**type** = relation name;

### Schema Example

```
((Movie (movieId title year rating))
 (IsGenre (movieId genre))
 (IsDirector (personId movieId))
 (CastMember (movieId personId role))
 (Entertainer (personId firstName lastName dob gender))
 (VideoStore (storeId name address city))
 (InStock (storeId movieId quantity available price))))
```

### Functional Dependencies Example

```
((Movie (movieId) (title year rating))
 (Entertainer (personId) (firstName lastName dob gender))
 (Movie (movieId) (title year rating))
 (VideoStore (storeId) (name address city))
 (InStock (storeId movieId) (quantity available price))))
```

### Query Example

```
(m Movie
  (- (y3 y4) (
    (VideoStore y3)
    (InStock y4)
    (= y3 city "Umea")
    (= m movieId y4 movieId)
    (= y4 storeId y3 storeId))
  )
  (= m year 2000)))
```

### 5.1.2 User-level Global Variables

**\*schema\*** Holds the current schema. Default is `nil`.

**\*simplify\*** `t`: Perform on-the-fly simplification. `nil`: Do not. Default is `t`.

### 5.1.3 User-level Functions

`setGlobalFDs fds` Given functional dependencies in the appropriate form, sets them as the current functional dependencies in the system.

*In these functions  $q$ ,  $q_1$ , and  $q_2$  are either  $\mathcal{L}$  or  $\mathcal{Q}$  queries, or a mixture of the two:*

`minus q1 q2` Computes  $q_1 - q_2$ .

`intersect` `q1` `q2` Computes  $q_1 \cap q_2$ .

`subsumes?` `q1` `q2` Computes whether  $q_1$  *subsumes*  $q_2$ .

`equiv?` `q1` `q2` Computes whether  $q_1$  and  $q_2$  are *equivalent*.

`sat?` `q` Computes whether  $q$  is *satisfiable*.

### 5.1.4 Work Order

1. Load into a Common LISP interpreter, e.g. CLISP, when standing in the `.../sqd/` directory and loading `load.lisp`
2. Set the current schema by assigning a value to `*schema*` via `(setf *schema* schema)`
3. Set the current functional dependencies by `(setGlobalFDs fds)`
4. Issue function or predicate calls
5. Return to step 2 or 3 if one wants to to change schema or functional dependencies

# Chapter 6

## Discussion

### 6.1 Related Work

This work extends [23] which assumed a Universal Relation[15] for the global schema and an extended form of relational algebra for the query language. The main improvement here is to lift these results to arbitrary relational schemas and to better specify the exact query form in standard logical notation.

We have identified a family of specification language for schema tuple queries. Since there is a direct translation from schema tuple queries to the domain relational calculus, this work may directly access the vast body of work on decidability classes for first order formulas[13]. The language  $\mathcal{Q}$  essentially falls within the Schönfinkel-Bernays class[4]. This is the class of function free formula that contain quantifier sequences of only the  $\exists^*\forall^*$  form. The extension of  $\mathcal{Q}$  toward more expressive schema tuple query languages, will no doubt borrow heavily from [13] and its sister collection [21] which chronicles interesting undecidable classes of first order formula.

As expected, the general problem of equivalence between relational algebra[10] expressions was shown to be undecidable[3]. Other work[8] singled out conjunctive queries (the select-project-join queries of the relational algebra) as a special case where query equivalence is decidable. Subsequent work specified a *first order query hierarchy* over query languages[7]. Roughly speaking,  $\mathcal{L}$  is contained within the class of conjunctive queries closed over complementation and  $\mathcal{L}_{pos}$  is contained within the

class of conjunctive queries closed over composition. Clearly neither  $\mathcal{L}$  nor  $\mathcal{L}_{pos}$  are *relationally complete*.

Recently most work around the question of query containment has adopted Datalog notion to represent queries. A *conjunctive query* ( $CQ$ ) in Datalog is simply a query where each predicate in the body of a rule references an extensional database relation. To decide if  $q_1 \subseteq q_2$  one first *freezes*[25]  $q_1$  by replacing the variables of its body and head with constants. Then if  $q_2$  includes the frozen head of  $q_1$  in its answer set when applied over the *canonical* database consisting of just the frozen predicates of  $q_1$ , we may conclude that  $q_1 \subseteq q_2$ . When no predicate appears more than once in the body of the rule, a linear time algorithm exists to decide containment, otherwise the decision problem is *NP*-complete. This approach may be enriched to decide containment between conjunctive queries with inequalities[18] ( $CQ^\neq$ ). Containment between conjunctive queries with negation of extensional predicates within their bodies ( $CQ^\neg$ ) may also be decided[20]. The complexity of deciding containment over  $CQ^\neq$  and  $CQ^\neg$  is  $\Pi_P^2$ . Containment between Datalog programs (Support recursion, but not negation) is undecidable[27]. Containment of a Datalog program within a conjunctive query is doubly exponential[9], while the converse question is easier.

Though there has been a lot of effort to chart languages over which containment and equivalence may be decided, to our knowledge no prior work has invoked the ‘schema tuple query’ assumption. Certainly one could imagine a regime in which the schema tuple assumption would be enforced over non-recursive Datalog. In such a case  $\mathcal{Q}$  would be contained in the class of non-recursive Datalog programs with negation in the rule bodies. The class of conjunctive queries with negation bears resemblance to the class  $\mathcal{L}$ , however  $\mathcal{L}$  is not contained within  $CQ^\neg$ , because negation in  $\mathcal{L}$  may span more than one predicate. For example a single query in  $CQ^\neg$ , could not express the second query of section 2.2.

While we have made no effort to gauge the complexity of query containment over  $\mathcal{Q}$ , it is clearly NP-hard in terms of query lengths. In fact, based on the fact that general Schönfinkle-Bernays satisfiability is NEXP, it is likely that the general containment problem is quite complex. Still since the complexity is worst case and is in query length, we anticipate that the approach will scale in real world environments.

It should be noted that the description logics[5][12] community has investigated ‘query’ containment under the name of *concept subsumption*. Description logics use unary and binary predicates which formalize traditional semantic network role/filler systems. As such they are interesting fragments of logic over predicates of at most two variables. The notion of syntactic difference between the concepts  $A(x)$  and  $B(x)$  may be represented as  $A(x) \sqcap \neg B(x)$ . Thus the focus here is on the ‘decidable’ description logics which have the conjunction ( $\sqcap$ ) and complement operator ( $\neg$ ). These are the description logics of the type *ALC* and beyond[12] which have sound and complete subsumption procedures. A serious limitation of description logics however, is their restriction to one and two place predicates. This has limited their impact in database environments. The recent introduction of *DLR* based description logics which are based on a unary concept and  $n$ -ary relationships[6] may hold some promise. Still, as an example, it is difficult to envision current description logics representing concepts such as “the movies whose director plays an acting role.”

Most work that considers the task of obtaining maximally contained rewritings assumes conjunctive queries for both the user query as well as the source descriptions. The problem of finding maximally contained rewritings under such assumptions is NP-complete[11]. However since it is often the case that the the number of conditions within a query is bounded, rewriting algorithms are still able to scale up to large numbers of sources. The bucket algorithm[19], the inverse-rules approach[14] and recently the mini-con algorithm[24] are all efforts toward scalable rewriting algorithms for conjunctive queries over conjunctive views. The work here considers user queries including general negation over sources described with conjunctive views.

The problem of identifying certain answers over data integration systems has some rather negative complexity results[1]. Unfortunately the complexity measures are in data complexity (e.g. number of provided tuples), rather than query complexity. The only class that gives polynomial complexity is that of conjunctive views under open world assumption of the database. Even inequalities in the conjunctive queries pushes the system into intractability.

The status of the schema tuple query assumption with regard to these issues is still not settled. It does stand to reason, however, that complexity may be in terms of

query size rather than data size. We have shown a method within chapter 4 to obtain a description of the certain answers to a query over a data integration environment. This yields a method by which to retrieve the full set of certain answers when the sources operate under a closed world assumption.

## 6.2 Future Directions

### 6.2.1 More expressive queries and schemas

Our contained rewriting algorithm of section 4.3 is unable to handle sources described using the full language  $\mathcal{L}$  because substitution of source descriptions using negation results in plans outside of  $\mathcal{L}$ . We seek decidable schema tuple query languages that satisfy the same closure properties as  $\mathcal{Q}$  which stay closed over such rewritings.

Another issue we would consider, is extending the types of schemas that we may cover. The inclusion of union types seems relatively straightforward. For example, if we extend the example of figure 2 to include private video collections, we could introduce the union type of *HasVideo* of which the relations *InStock* and *InPersonalCollection* would be members. Such a union type would enable queries to be posed over the more general category using common attributes. Tuples, however, would still be members of one, and only one, member predicate. Formally queries using union types may always be rewritten in a form within  $\mathcal{Q}$  that only references base level, non union type, schema relations.

Modeling IS-A hierarchies through inclusion dependencies would violate the assumption that tuples are drawn from a single schema relation. Since union types may be handled, and since the decidability of containment enables subsumption hierarchies to be built, we conjecture that modeling true IS-A hierarchies is unnecessary for most domains. Finally it will be interesting to see how cardinality and foreign key constraints could be integrated into this work.



### 6.2.2 Annotating source advertisements

We should be able to enable sources to claim open-world in addition to closed world coverage over their ‘views’ of the global schema. This would extend the breadth of application of our approach and would enrich intensional responses. Specifically responses would need to distinguish between the *certain and complete* and the *certain and partial* answers to the user’s query. Certainly the extension to open world data sources will have a significant complicating impact on the process of identifying certain answers.

We also acknowledge that the assumption that data sources provide only correct information is too optimistic. We would explore how to discover and summarize differences of opinion between the data sources relative to a specific user query.

### 6.2.3 Simplification and natural language generation

Intensional descriptions within  $\mathcal{Q}$  must be simplified before they are presented. Such simplification may be envisioned as an informed search process where utility is inversely correlated with expression complexity and operators split or fuse the disjuncts of the expression.

Assuming that the complexity measure is simply the number of disjuncts within an expression, the fact that there are a finite number of constants within the input expression, means that there are a finite number of reasonable splitting and fusion operations. Given these considerations, it seems likely that modern search techniques may be brought to bear on the problem of minimizing the number of disjunction within an expression in  $\mathcal{Q}$ . Given this, we shall focus on generating descriptions of each disjunct in the simplified query. Thus we focus on the problem of generating descriptions of expressions in  $\mathcal{L}$ .

An initial approach to generating natural language description from expressions within  $\mathcal{L}$  has been proposed in [22]. The approach uses a phrasal lexicon and exploits the decidability of containment over  $\mathcal{L}$ .

This page intentionally contains only this sentence.

# Chapter 7

## Conclusions

This thesis has introduced the sub-classes of *schema tuple queries* specified in the languages  $\mathcal{L}$ ,  $\mathcal{L}_{pos}$ ,  $\mathcal{Q}$  and  $\mathcal{Q}_{pos}$ . Formulas in the language  $\mathcal{L}$  are signed quantifier sequences ( $\exists^*$ ) over conjunctions of predicates over a single free tuple variable. Formulas in the language  $\mathcal{Q}$  are finite disjunctions of formulas within  $\mathcal{L}$ . The languages  $\mathcal{L}_{pos}$  and  $\mathcal{Q}_{pos}$  are corresponding languages which disallow negation of existential quantifier sequences.

Within the class of queries specified using  $\mathcal{Q}$ , one may express the difference between two queries as a third query within the class. This ability to calculate syntactic query difference enables the generation of intensional descriptions of query differences and intersections. Though the approach is limited to queries returning full schema tuples, such limitations may often be overcome by considering a virtual, highly decomposed version of the schema. Thus the schema tuple query assumption may be appropriate for many, if not most, real world schemas.

Sidestepping the schema tuple query restriction, the queries specified using  $\mathcal{L}$  are more expressive than conjunctive queries with negated subgoals, but are less expressive than non-recursive Datalog with negated subgoals. However, the query classes here are more naturally expressed in tuple relational calculus than in Datalog. Resting the form of these query classes within tuple calculus also enables us to directly access a large body of classical work on decidable classes of logical formulas.

There is a very direct correspondence between queries specified in  $\mathcal{L}$  (and  $\mathcal{Q}$ ) and

standard SQL. Queries in  $\mathcal{L}$  mirror the SQL of the form:

```
SELECT *  
FROM R AS x  
WHERE  
  [NOT] EXISTS (sub-query) ... ;
```

where the *sub-query* is of the same form but may not use NOT. Naturally simple and join conditions may be added to such queries and there may be more than one sub-query. Queries specified  $\mathcal{Q}$  correspond to the UNION of such SQL expressions.

The main practical/theoretical result of this thesis is to show that if data sources are described using the language  $\mathcal{Q}_{pos}$  and the user's query is specified using the language  $\mathcal{L}$ , then one may: 1.) maximally rewrite queries using source descriptions; 2.) retrieve *certain answers* over the integrated system; 3.) generate intensional descriptions (within the language  $\mathcal{Q}$ ) of the certain, the uncertain and the missing answers to the user's query over the available data sources. The main practical result is a LISP-system that is capable of computing the syntactic difference and intersection of queries in the languages  $\mathcal{L}$  and  $\mathcal{Q}$ , and able to compute the predicates *subsumes*, *equality*, and *disjointness*.

It seems reasonable to suppose that other problems will have interesting simplifications when restricted to the schema tuple queries specified using  $\mathcal{L}$  and  $\mathcal{Q}$ . Of particular interest are problems in updates over views, semantic query optimization and cooperative query answering[16].

# Chapter 8

## Acknowledgments

I would like to thank Michael Minock for his useful theoretical pointers and valuable support in the work process, as well as his substantial contribution of writing support of this thesis. I would also like to thank the Department of Computing Science at the University of Umeå for providing a pleasant environment in which to conduct this work.

This page intentionally contains only this sentence.

# Bibliography

- [1] S. Abiteboul and O. Duschka. Complexity of answering queries using materialized views. In *Sym. on Principles of Database Systems*, pages 254–263, 1998.
- [2] S. Abiteboul, R. Viannu, and V. Hull. *Foundations of Database Systems 3rd edition*. Addison Wesley, 1995.
- [3] A. Aho, Y. Sagiv., and J. Ullman. Equivalences of relational expressions. *SIAM Journal on Computing*, 8(2):218–246, 1979.
- [4] P. Bernays and M. Schönfinkel. Zum entscheidungsproblem der mathematischen logik. *M.A.*, 99:342–372, 1928.
- [5] R. Brachman and J. Schmolze. An overview of the KL-ONE knowledge representation system. *Cognitive Science*, 9(2):171–216, 1985.
- [6] D. Calvanese, G. De Giacomo, M. Lenzerini, D. Nardi, and R. Rosati. Description logic framework for information integration. In *Principles of Knowledge Representation and Reasoning*, pages 2–13, 1998.
- [7] A. Chandra and D. Harel. Structure and complexity of relational queries. *Journal of Computer Systems and Sciences*, 25(1):99–128, 1982.
- [8] A. Chandra and P. Merlin. Optimal implementation of conjunctive queries in relational databases. In *Proc. 9 of the ACM Sym. on the Theory of Computing*, pages 77–90., 1977.

- 
- [9] S. Chaudhuri and M. Vardi. On the equivalence of recursive and nonrecursive datalog programs. In *Sym. on Principles of Database Systems*, pages 55–66, 1992.
- [10] E. Codd. Relational completeness of data base sublanguages. In R. Rustin, editor, *Database Systems*, pages 33–64. Prentice-Hall, 1972.
- [11] J. Van den Busshe. Two remarks on the complexity of answering queries using views. In *Information Processing Letters*, 2000.
- [12] F. Donini, M. Lenzerini, D. Nardi, and A. Schaerf. Reasoning in description logics. In G. Brewka, editor, *Studies in Logic, Language and Information*, pages 193–238, 1996.
- [13] B. Dreben and W.D. Goldfarb. *The decision problem. Solvable classes of quantificational formulas*. Addison Wesley, 1979.
- [14] O. Duschka, M. Genesereth, and A. Levy. Recursive query plans for data integration. *Journal of Logic Programming*, 43(1):49–73, 2000.
- [15] R. Fagin, A. Mendelzon, and J. Ullman. A simplified universal relation assumption and its properties. *ACM Transactions on Database Systems*, 7, 1982.
- [16] T. Gaasterland, P. Godfrey, and J. Minker. An overview of cooperative answering. *Intelligent Information Systems*, 1(2):127–157, 1992.
- [17] A. Halevy. Answering queries using views: A survey. *VLDB Journal* 10(4): 270–294, 2001.
- [18] A. Klug. On conjunctive queries containing inequalities. *Journal of the ACM*, 35(1):146–160, 1988.
- [19] A. Levy, A. Rajaraman, and J. Ordille. Query-answering algorithms for information agents. In *Proc. of the Thirteenth National Conference on Artificial Intelligence*, pages 40–47, 1996.



- 
- [20] A. Levy and Y. Sagiv. Queries independent of updates. In *Proc. of VLDB*, pages 171–181, 1993.
- [21] H. Lewis. *Unsolvable Classes of Quantificational Formulas*. Addison Wesley, 1979.
- [22] M. Minock. A phrasal generator for describing relational database queries. In *Proc. of the 9th EACL workshop on natural language generation*, Budapest, Hungary, April 2003.
- [23] M. Minock, M. Rusinkiewicz, and B. Perry. The identification of missing information resource agents by using the query difference operator. In *COOPIS '99*. IEEE Computer Society Press, 1999.
- [24] R. Pottinger and A. Halevy. A scalable algorithm for answering queries using views. In *The VLDB Journal*, pages 484–495, 2000.
- [25] R. Ramakrishnan., Y. Sagiv, J. Ullman, and M. Vardi. Proof tree transformation theorems and their applications. In *Sym. on Principles of Database Systems*, pages 172–181, 1989.
- [26] R. Reiter. *Logic and Databases*, chapter On Closed World Databases. Plenum Press, 1978.
- [27] O. Shmueli. Decidability and expressiveness of logic queries. In *Sym. on Principles of Database Systems*, pages 237–249, 1987.
- [28] J. Ullman. Information integration using logical views. *Theoretical Computer Science*, 239(2):189–210, 2000.

This page intentionally contains only this sentence.

# Appendix A

## Source Code

### A.1 /sqd/

#### A.1.1 load.lisp

```
(defvar *sat-debug* nil)

(compile-file "lsat/fol-utils") (load "lsat/fol-utils.lisp")
(compile-file "lsat/l2fol.lisp") (load "lsat/l2fol.lisp")
(compile-file "lsat/constants2fol.lisp") (load "lsat/constants2fol.lisp")
(compile-file "lsat/fds2fol.lisp") (load "lsat/fds2fol.lisp")
(compile-file "lsat/unify") (load "lsat/unify")
(compile-file "lsat/normal") (load "lsat/normal")
(compile-file "lsat/fol-sat") (load "lsat/fol-sat")
```

#### A.1.2 sqd-engine.lisp

```
(defvar *schema* '()) ;; The schema
(defvar *fds* '()) ;; The functional dependencies
(defvar *simplify* t)

;; *** General methods - take either an 'l' or a 'q' query ***

(defun minus (q1 q2)
  (cond ((and (l-query? q1) (l-query? q2))
        (cons 'or (l1-minus q1 q2)))
        ((and (q-query? q1) (l-query? q2))
        (cons 'or (q1-minus (rest q1) q2)))
        ((and (l-query? q1) (q-query? q2))
        (cons 'or (qq-minus (list q1) (rest q2))))
        (t
         (cons 'or (qq-minus (rest q1) (rest q2))))))

(defun intersect (q1 q2)
  (minus q1 (minus q1 q2)))
```

```

(defun subsumes? (q1 q2)
  (not (q-sat? (minus q2 q1))))

(defun equiv? (q1 q2)
  (and (subsumes? q1 q2)
       (subsumes? q2 q1)))

(defun sat? (q)
  (cond ((q-query? q) (q-sat? q))
        (t (l-sat? q))))

;; *** Typed variants - take a specific form: either l or q. Note that ;; the 'OR has been striped
from 'q'. ***

(defun q-sat? (query)
  (cond ((null query) nil)
        ((equal (first query) 'OR)
         (q-sat? (rest query)))
        (t
         (or (l-sat? (first query)) (q-sat? (rest query))))))

(defun l-sat? (l)
  (let ((una-constants (unas (multiple-value-list (get-constants l))))
        (l-fol (2fol2 l)))
    (fol-sat? l-fol (append '(and) una-constants (getRelevantFDs l)))))

(defun qq-minus (q1 q2)
  (cond ((null q2) q1)
        ((null q1) '())
        (t
         (qq-minus
          (ql-minus q1 (first q2))
          (rest q2)))))

(defun ql-minus (q1 l2)
  (cond ((null q1) '())
        (t (append (ll-minus (first q1) l2) (ql-minus (rest q1) l2)))))

(defun ll-minus (l1 l2)
  (demorgify l1 (subst (first l1) (first l2) (rest (rest l2)))))

(defun demorgify (l components)
  (cond ((null components) '())
        (t
         (let ((new-l (add-component l (negateComponent (car components)))))
           (if (or (not *simplify*) (l-sat? new-l))
               (cons new-l (demorgify l (rest components)))
               (demorgify l (rest components)))))))

```

### A.1.3 tuple-query-utils.lisp

```
;; *** 'Fast' Type Checkers ***
```

```
(defun q-query? (query)
  (equal (first query) 'OR))
```

```
(defun l-query? (query)
  (not (q-query? query)))
```

```

(defun complexCond? (aCond)
  (if (and (equal 3 (length aCond))
          (or (equal '+ (first aCond))
              (equal '- (first aCond)))) t nil))

(defun signed-query-comp? (onecond)
  (or (equal '- (first onecond)) (equal '+ (first onecond))))

(defun simple-condition? (onecond)
  (and (= (length onecond) 4) (not (listp (fourth onecond)))))

(defun set-condition? (onecond)
  (and (= (length onecond) 4) (listp (fourth onecond))))

(defun sqc? (oneCond)
  (or (sqc?-fol oneCond) (signed-query-comp? oneCond)))

(defun sqc?-fol (onecond)
  (or (equal 'not (first onecond)) (equal 'exists (first onecond))))

(defun join-cond?-dc (one-cond)
  (and (equal 3 (length one-cond))
       (equal '= (first one-cond))
       (dc-var? (second one-cond))
       (dc-var? (third one-cond))))

(defun negateComponent (aCond)
  (cond ((complexCond? aCond) (negateComplex aCond))
        (t (negateSimple aCond))))

(defun add-component (l component)
  (append l (list component)))

(defun negateSimple (aCond)
  (let ((op (first aCond))
        (condtail (rest aCond)))
    (cond ((equal op '>) (cons '<= condtail))
          ((equal op '<') (cons '>= condtail))
          ((equal op '>=') (cons '< condtail))
          ((equal op '<=') (cons '> condtail))
          ((equal op '=) (cons '<> condtail))
          ((equal op '<>') (cons '= condtail))
          ((equal op 'IN) (cons 'NOT_IN condtail))
          ((equal op 'NOT_IN) (cons 'IN condtail))
          (t (cons 'faulty_input_to_negatesimple (list aCond)))))

(defun negateComplex (aCond)
  (if (equal '+ (first aCond))
      (cons '- (rest aCond))
      (cons '+ (rest aCond))))

```

### A.1.4 utilities.lisp

```

;;; -*- Mode: Lisp; Syntax: Common-Lisp; -*- File: utilities.lisp

;;; Basic utility functions and macros, used throughout the code.

;;; The utilities are divided into control flow macros, list
;;; utilities, functions for 2-dimensional points, numeric utilities,
;;; some trivial functions, utilities for strings, symbolsand

```

```

;;; printing, a debugging tool, and a testing tool."

(eval-when (eval compile load)
  ;; Make it ok to place a function definition on a built-in LISP symbol.
  #+(or Allegro EXCL)
  (dolist (pkg '(excl common-lisp common-lisp-user))
    (setf (excl:package-definition-lock (find-package pkg)) nil)))

;;; Control Flow Macros

;;; We define iteration macros to match the book's pseudo-code.
;;; This could all be done with LOOP, but some users don't have
;;; the LOOP from the 2nd edition of 'Common Lisp: the Language'.

(defmacro while (test do &body body)
  "Execute body while the test is true."
  (assert (eq do 'do))
  `(do () ((not ,test) nil) ,@body))

(defmacro for-each (var in list do &body body)
  "Execute body for each element of list. VAR can be a list or tree
of variables, in which case the elements are destructured."
  (assert (eq in 'in)) (assert (eq do 'do))
  (typecase var
    (symbol `(dolist (,var ,list) ,@body))
    (cons (let ((list-var (gensym)))
            `(dolist (,list-var ,list)
              (destructuring-bind ,var ,list-var ,@body))))
    (t (error "~V is an illegal variable in (for each ~V in ~A ...)"
              var list))))

(defmacro for (var = start to end do &body body)
  "Execute body with var bound to successive integers."
  (cond ((eq var 'each) ; Allow (for each ...) instead of (for-each ...)
        `(for-each ,= ,start ,to ,end ,do ,@body))
    (t (assert (eq = '=)) (assert (eq to 'to)) (assert (eq do 'do))
       (let ((end-var (gensym "END")))
         `(do ((,var ,start (+ 1 ,var)) (,end-var ,end))
             (> ,var ,end-var) nil)
           ,@body))))

(defmacro deletf (item sequence &rest keys &environment env)
  "Destructively delete item from sequence, which must be SETF-able."
  (multiple-value-bind (temps vals stores store-form access-form)
    (get-setf-expansion sequence env)
    (assert (= (length stores) 1))
    (let ((item-var (gensym "ITEM")))
      `(let* ((,item-var ,item)
              ,@(mapcar #'list temps vals)
              (,first stores) (delete ,item-var ,access-form ,@keys))
         ,store-form)))

(defmacro define-if-undefined (&rest definitions)
  "Use this to conditionally define functions, variables, or macros that
may or may not be pre-defined in this Lisp. This can be used to provide
CLtL2 compatibility for older Lisps."
  `(progn
    ,@(mapcar #'(lambda (def)
                  (let ((name (second def)))
                    `(when (not (or (boundp ',name) (fboundp ',name)
                                     (special-form-p ',name)
                                     (macro-function ',name)))
                      ,def)))
              definitions)))

```

```

                (macro-function ',name)))
            ,def)))
    definitions)))

;;; List Utilities

(defun length>1 (list)
  "Is this a list of 2 or more elements?"
  (and (consp list) (cdr list)))

(defun length=1 (list)
  "Is this a list of exactly one element?"
  (and (consp list) (null (cdr list))))

(defun random-element (list)
  "Return some element of the list, chosen at random."
  (nth (random (length list)) list))

(defun mappend (fn &rest lists)
  "Apply fn to respective elements of list(s), and append results."
  (reduce #'append (apply #'mapcar fn lists) :from-end t))

(defun starts-with (list element)
  "Is this a list that starts with the given element?"
  (and (consp list) (eq (first list) element)))

(defun last1 (list)
  "Return the last element of a list."
  (first (last list)))

(defun left-rotate (list)
  "Move the first element to the end of the list."
  (append (rest list) (list (first list))))

(defun right-rotate (list)
  "Move the last element to the front of the list."
  (append (last list) (butlast list)))

(defun transpose (list-of-lists)
  "Transpose a matrix represented as a list of lists.
  Example: (transpose '((a b c) (d e f))) => ((a d) (b e) (c f))."
  (apply #'mapcar #'list list-of-lists))

(defun reuse-cons (x y x-y)
  "Return (cons x y), or reuse x-y if it is equal to (cons x y)"
  (if (and (eql x (car x-y)) (eql y (cdr x-y)))
      x-y
      (cons x y)))

"An expression is a list consisting of a prefix operator followed by args, Or it can be a symbol,
denoting an operator with no arguments. Expressions are used in Logic, and as actions for agents."

(defun make-exp (op &rest args) (cons op args))

(defun op (exp)
  "Operator of an expression"
  (if (listp exp) (first exp) exp))

(defun args (exp) "Arguments of an expression" (if (listp exp) (rest
exp) nil))

```

```

(defun arg1 (exp)
  "First argument"
  (first (args exp)))

(defun arg2 (exp)
  "Second argument"
  (second (args exp)))

(defsetf args (exp) (new-value)
  '(setf (cdr ,exp) ,new-value))

(defun prefix->infix (exp)
  "Convert a fully parenthesized prefix expression into infix notation."
  (cond ((atom exp) exp)
        ((length=1 (args exp)) exp)
        (t (insert-between (op exp) (mapcar #'prefix->infix (args exp))))))

(defun insert-between (item list)
  "Insert item between every element of list."
  (if (or (null list) (length=1 list))
      list
      (list* (first list) item (insert-between item (rest list)))))

;;; Functions for manipulating 2-dimensional points

(defstruct (xy (:type list)) "A two-dimensional (i.e. x and y) point." x y)

(defun xy-p (arg)
  "Is the argument a 2-D point?"
  (and (consp arg) (= (length arg) 2) (every #'numberp arg)))

(defun @ (x y) "Create a 2-D point" (make-xy :x x :y y))

(defun xy-equal (p q) (equal p q))

(defun xy-add (p q)
  "Add two points, component-wise."
  (@ (+ (xy-x p) (xy-x q)) (+ (xy-y p) (xy-y q))))

(defun xy-distance (p q)
  "The distance between two points."
  (sqrt (+ (square (- (xy-x p) (xy-x q)))
           (square (- (xy-y p) (xy-y q))))))

(defun x+y-distance (p q)
  "The 'city block distance' between two points."
  (+ (abs (- (xy-x p) (xy-x q)))
     (abs (- (xy-y p) (xy-y q)))))

(defun xy-between (xy1 xy2 xy3)
  "Predicate; return t iff xy1 is between xy2 and xy3. Points are collinear."
  (and (between (xy-x xy1) (xy-x xy2) (xy-x xy3))
       (between (xy-y xy1) (xy-y xy2) (xy-y xy3))))

(defun rotate (o a b c d)
  (let ((x (xy-x o))
        (y (xy-y o)))
    (@ (+ (* a x) (* b y)) (+ (* c x) (* d y)))))

(defun inside (l xmax ymax)
  "Is the point l inside a rectangle from 0,0 to xmax,ymax?"

```



```

(let ((x (xy-x 1)) (y (xy-y 1)))
  (and (>= x 0) (>= y 0) (< x xmax) (< y ymax)))

;;; Numeric Utilities

(defconstant infinity most-positive-single-float) (defconstant minus-infinity
most-negative-single-float)

(defun average (numbers)
  "Numerical average (mean) of a list of numbers."
  (/ (sum numbers) (length numbers)))

(defun running-average (avg new n)
  "Calculate new average given previous average over n data points"
  (/ (+ new (* avg n)) (1+ n)))

(defun square (x) (* x x))

(defun sum (numbers &optional (key #'identity))
  "Add up all the numbers; if KEY is given, apply it to each number first."
  (if (null numbers)
      0
      (+ (funcall key (first numbers)) (sum (rest numbers) key))))

(defun between (x y z)
  "Predicate; return t iff number x is between numbers y and z."
  (or (<= y x z) (>= y x z)))

(defun rms-error (predicted target)
  "Compute root mean square error between predicted list and target list"
  (sqrt (ms-error predicted target)))

(defun ms-error (predicted target &aux (sum 0))
  "Compute mean square error between predicted list and target list"
  (mapc #'(lambda (x y) (incf sum (square (- x y)))) predicted target)
  (/ sum (length predicted)))

(defun boolean-error (predicted target)
  (if (equal predicted target) 0 1))

(defun dot-product (l1 l2 &aux (sum 0)) ;;; dot product of two lists
  (mapc #'(lambda (x1 x2) (incf sum (* x1 x2))) l1 l2)
  sum)

(defun iota (n &optional (start-at 0))
  "Return a list of n consecutive integers, by default starting at 0."
  (if (<= n 0) nil (cons start-at (iota (- n 1) (+ start-at 1)))))

(defun random-integer (from to)
  "Return an integer chosen at random from the given interval."
  (+ from (random (+ 1 (- to from)))))

(defun normal (x mu sigma)
  (/ (exp (/ (- (square (- x mu))) (* 2 (square sigma))))
     (* (sqrt (* 2 pi)) sigma)))

(defun sample-with-replacement (n population)
  (let ((result nil))
    (dotimes (i n) (push (random-element population) result))
    result))

```

```

(defun sample-without-replacement (n population &optional
                                   (m (length population)))
  ;; Assumes that m = (length population)
  (cond ((<= n 0) nil)
        ((>= n m) population)
        ((>= (/ n m) (random 1.0))
         (cons (first population) (sample-without-replacement
                                   (- n 1) (rest population) (- m 1))))
        (t (sample-without-replacement n (rest population) (- m 1)))))

(defun fuzz (quantity &optional (proportion .1) (round-off .01))
  "Add and also subtract a random fuzz-factor to a quantity."
  (round-off (+ quantity
                (* quantity (- (random (float proportion))
                                (random (float proportion))))))
             round-off))

(defun round-off (number precision)
  "Round off the number to specified precision. E.g. (round-off 1.23 .1) = 1.2"
  (* precision (round number precision)))

;;; Trivial Functions

(defun nothing (&rest args)
  "Don't do anything, and return nil."
  (declare (ignore args))
  nil)

(defun declare-ignore (&rest args)
  "Ignore the arguments."
  ;; This is used to avoid compiler warnings in defmethod.
  ;; Some compilers warn "Variable unused" if it is bound by a method
  ;; but does not appear in the body. However, if you put in a
  ;; (declare (ignore var)), then other compilers warn "var declared
  ;; ignored, but is actually used", on the grounds that it is implicitly
  ;; used to do method dispatch. So its safest to use declare-ignore.
  ;; If you like, you can redefine declare-ignore to be a macro that
  ;; expands to either (declare (ignore args)), or to nothing, depending
  ;; on the implementation.
  (declare (ignore args))
  nil)

#-(or MCL Lispworks) ;; MCL, Lispworks already define this function (defun true (&rest args) "Always
return true." (declare (ignore args)) t)

#-(or MCL Lispworks) ;; MCL, Lispworks already define this function (defun false (&rest args) "Always
return false." (declare (ignore args)) nil)

(defun required (&optional (msg "A required argument is missing.") &rest args)
  "If this ever gets called, it means something that was required was not
supplied. Use as default value for &key args or defstruct slots."
  (apply #'error msg args))

;;; Utilities for strings and symbols and printing

(defun stringify (exp)
  "Coerce argument to a string."
  (cond ((stringp exp) exp)
        ((symbolp exp) (symbol-name exp))
        (t (format nil "~A" exp))))

```

```

(defun concat-symbol (&rest args)
  "Concatenate the args into one string, and turn that into a symbol."
  (intern (format nil "~~a~" args)))

(defun print-grid (array &key (stream t) (key #'identity) (width 3))
  "Print the contents of a 2-D array, numbering the edges."
  (let ((max-x (- (array-dimension array 0) 1))
        (max-y (- (array-dimension array 1) 1)))
    ;; Print the header
    (format stream "~&") (print-repeated " " width stream)
    (for x = 0 to max-x do
      (format stream "|") (print-dashes width stream)
      (format stream "|~%")
    ;; Print each row
    (for y1 = 0 to max-y do
      (let ((y (- max-y y1)))
        (print-centered y width stream)
        ;; Print each location
        (for x = 0 to max-x do
          (format stream "|")
          (print-centered (funcall key (aref array x y)) width stream)
          (format stream "|~%")
        ;; Print a dashed line
        (print-repeated " " width stream)
        (for x = 0 to max-x do
          (format stream "|") (print-dashes width stream)))
        (format stream "|~%"))
    ;; Print the X-coordinates along the bottom
    (print-repeated " " width stream)
    (for x = 0 to max-x do
      (format stream " ") (print-centered x width stream))
    array))

(defun print-centered (string width &optional (stream t))
  "Print STRING centered in a field WIDTH wide."
  (let ((blanks (- width (length (stringify string))))))
    (print-repeated " " (floor blanks 2) stream)
    (format stream "~A" string)
    (print-repeated " " (ceiling blanks 2) stream)))

(defun print-repeated (string n &optional (stream t))
  "Print the string n times."
  (dotimes (i n)
    (format stream "~A" string)))

(defun print-dashes (width &optional (stream t) separate-line)
  "Print a line of dashes WIDTH wide."
  (when separate-line (format stream "~&"))
  (print-repeated "-" width stream)
  (when separate-line (format stream "~%")))

;;; Assorted conversion utilities and predicates

(defun copy-array (a &aux (dim (array-dimensions a))
                   (b (make-array dim)))
  "Make a copy of an array."
  (copy-subarray a b nil dim)
  b)

(defun copy-subarray (a b indices dim)
  (if dim

```

```

(dotimes (i (first dim))
  (copy-subarray a b (append indices (list i)) (rest dim)))
(setf (apply #'aref (cons b indices))
      (apply #'aref (cons a indices))))

(defun array->vector (array)
  "Convert a multi-dimensional array to a vector with the same elements."
  (make-array (array-total-size array) :displaced-to array))

(defun plot-alist (alist file)
  (with-open-file (stream file :direction :output :if-does-not-exist :create
                      :if-exists :supersede)
    (dolist (xy alist)
      (format stream "~&~A ~A~%" (car xy) (cdr xy)))))

(defun copy-hash-table (H1 &optional (copy-fn #'identity))
  (let ((H2 (make-hash-table :test #'equal)))
    (maphash #'(lambda (key val) (setf (gethash key H2) (funcall copy-fn val)))
             H1)
    H2))

(defun hash-table->list (table)
  "Convert a hash table into a list of (key . val) pairs."
  (maphash #'cons table))

(defun hprint (h &optional (stream t))
  "prints a hash table line by line"
  (maphash #'(lambda (key val) (format stream "~&~A:~10T ~A" key val)) h)
  h)

(defun compose (f g)
  "Return a function h such that (h x) = (f (g x))."
  #'(lambda (x) (funcall f (funcall g x))))

(defun the-biggest (fn l)
  (let ((biggest (first l))
        (best-val (funcall fn (first l))))
    (dolist (x (rest l))
      (let ((val (funcall fn x)))
        (when (> val best-val)
          (setq best-val val)
          (setq biggest x))))
    biggest))

(defun the-biggest-random-tie (fn l)
  (random-element
   (let ((biggest (list (first l)))
         (best-val (funcall fn (first l))))
     (dolist (x (rest l))
       (let ((val (funcall fn x)))
         (cond ((> val best-val)
                (setq best-val val)
                (setq biggest (list x)))
              ((= val best-val)
                (push x biggest))))
       biggest)))

(defun the-biggest-that (fn p l)
  (let ((biggest (first l))
        (best-val (funcall fn (first l))))
    (best-val (funcall fn (first l))))

```

```

(dolist (x (rest l))
  (when (funcall p x)
    (let ((val (funcall fn x)))
      (when (> val best-val)
        (setq best-val val)
        (setq biggest x))))))
biggest))

(defun the-smallest (fn l)
  (the-biggest (compose #'- fn) l))

(defun the-smallest-random-tie (fn l)
  (the-biggest-random-tie (compose #'- fn) l))

(defun the-smallest-that (fn p l)
  (the-biggest-that (compose #'- fn) p l))

;;; Debugging tool

(defvar *debugging* nil)

(defun dprint (&rest args)
  "Echo all the args when *debugging* is true. Return the first one."
  (when *debugging* (format t "&~S ~%" args))
  (first args))

;;; Testing Tool: deftest and test

(defmacro deftest (name &rest examples)
  "Define a set of test examples. Each example is of the form (exp test)
  or (exp). Evaluate exp and see if the result passes the test. Within the
  test, the result is bound to *. The example ((f 2)) has no test to
  fail, so it always passes the test. But ((+ 2 2) (= * 3)) has the test
  (= * 3), which fails because * will be bound to the result 4, so the test
  fails. Call (TEST name) to count how many tests are failed within the
  named test. NAME is the name of an aima-system."
  `(add-test ',name ',examples))

(defun add-test (name examples)
  "The functional interface for deftest: adds test examples to a system."
  (let ((system (or (get-aima-system name)
                    (add-aima-system :name name :examples examples))))
    (setf (aima-system-examples system) examples))
  name)

(defun test (&optional (name 'all) (print? 't))
  "Run a test suite and sum the number of errors. If all is well, this
  should return 0. The second argument says what to print: nil for
  nothing, t for everything, or FAIL for just those examples that fail.
  If there are no test examples in the named system, put the system has
  other systems as parts, run the tests for all those and sum the result."
  (let ((*print-pretty* t)
        (*standard-output* (if print? *standard-output*
                                (make-broadcast-stream))))
    (system (aima-load-if-unloaded name)))
    (cond ((null system) (warn "No such system as ~A." name))
          ((and (null (aima-system-examples system))
                (every #'symbolp (aima-system-parts system))))
          (sum (aima-system-parts system)
                #'(lambda (part) (test part print?))))
          (t (when print? (format t "Testing System ~A%" name))
              (sum (aima-system-parts system)
                    #'(lambda (part) (test part print?))))))

```

```

      (let ((errors (count-if-not #'(lambda (example)
                                   (test-example example print?))
                                (aima-system-examples system))))
        (format *debug-io* "~%~2D error~P on system ~A~%"
                errors errors name)
        errors))))))

(defun test-example (example &optional (print? t))
  "Does the EXP part of this example pass the TEST?"
  (if (stringp example)
      (progn
        (when (eq print? t)
          (format t "~&;; ~A~%" example))
        t)
      (let* ((exp (first example))
             (* nil)
             (test (cond ((null (second example)) t)
                         ((constantp (second example))
                          '(equal * ,(second example)))
                         (t (second example))))
             test-result)
        (when (eq print? t)
          (format t "~&> ~S~%" exp))
        (setf * (eval exp))
        (when (eq print? t)
          (format t "~&~S~%" *))
        (setf test-result (eval test))
        (when (null test-result)
          (case print?
            ((FAIL) (format t "~&;; FAILURE on ~S; expected ~S, got:~%;;; ~S~%"
                            exp test *))
            ((T) (format t "~&;; FAILURE: expected ~S" test))
            (otherwise)))
          test-result)))

(defun mklist (x)
  "If x is a list, return it; otherwise return a singleton list, (x)."  

  (if (listp x) x (list x)))

```

## A.2 /sqd/lsat/

### A.2.1 load.lisp

```

(defvar *sat-debug* nil)

(compile-file "lsat/fol-utils") (load "lsat/fol-utils.lisp")
(compile-file "lsat/l2fol.lisp") (load "lsat/l2fol.lisp")
(compile-file "lsat/constants2fol.lisp") (load "lsat/constants2fol.lisp")
(compile-file "lsat/fds2fol.lisp") (load "lsat/fds2fol.lisp")
(compile-file "lsat/unify") (load "lsat/unify")
(compile-file "lsat/normal") (load "lsat/normal")
(compile-file "lsat/fol-sat") (load "lsat/fol-sat")

```

### A.2.2 constants2fol.lisp

```

(defun get-constants (l)
  (let* ((numbers '())
         (strings '()))

```

```

      (others '() )
      (components (cddr l)))
(labels ((pushTypedConst (c)
  (typecase c
    (string (pushnew c strings :test #'string-equal))
    (number (pushnew c numbers))
    (t (pushnew c others :test #'equalp))))
  (getConstantsFromComponents (components)
    (dolist (comp components (values numbers strings others))
      (cond ((sqc? comp)
        (getConstantsFromComponents (caddr comp)))
        ((setCond? comp)
        (mapcar #'pushTypedConst (fourth comp)))
        ((simple-condition? comp)
        (pushTypedConst (fourth comp)))
        (t nil))))))
  (getConstantsFromComponents components))))

(defun unas (sets)
  (reduce #'append (mapcar #'una sets)))

(defun una (set)
  (cond ((null set) '())
    (t (append (una-single (first set) (rest set))
      (una (rest set))))))

(defun una-single (val set)
  (cond ((null set) '())
    (t (append (una-rule val (first set)) (una-single val (rest set))))))

(defun una-rule (val1 val2)
  (list '(forall ($V) (=> (= $V ,val1) (not (= $V ,val2))))
    '(forall ($V) (=> (= $V ,val2) (not (= $V ,val1))))))



### A.2.3 fds2fol.lisp



```

(defun *fol-fds* (make-hash-table))

(defun getRelevantFDs (l)
  "Retrieves the FDs that are relevant to the relations in a query."
  (let ((relations-in-l (getTypes l))
    (relevant-fol-fds '()))
    (dolist (type relations-in-l relevant-fol-fds)
      (setf relevant-fol-fds (append (gethash type *fol-fds*) relevant-fol-fds))))))

(defun fds->fol (fds)
  (setGlobalFDs fds)
  (let ((res '()))
    (maphash #'(lambda (key val) (declare (ignore key)) (setf res (append val res))) *fol-fds*)
    res))

(defun setGlobalFDs (fds)
  "Transforms the functional dependencies in <fds> to FOL (Horn-clauses
really), and puts them in a hash table, with type (relation name) as key."
  (clrhash *fol-fds*)
  (dolist (fd fds t)
    (let ((type (first fd)))
      (setf (gethash type *fol-fds*) (append (fd2fol fd) (gethash type *fol-fds*)))))

(defun fd2fol (fd)

```


```

```

"Transforms a functional dependency in <fd> to a corresponding set of Horn-clauses."
(let* ((type (first fd))
      (LHS (second fd))
      (RHS (third fd))
      (attrs (cadr (assoc type *schema*)))
      (vars '())
      (impl-pairs '())
      (pred1 '())
      (pred2 '())
      (forall '())
      (new-fd-set '()))
  (dolist (attr attrs (progn (setf forall (reverse forall))
                             (setf pred1 (cons type (reverse pred1)))
                             (setf pred2 (cons type (reverse pred2)))))
    (cond ((member attr LHS)
           (setf vars (list (genVarName 'x))))
          ((member attr RHS)
           (setf vars (list (genVarName 'y) (genVarName 'y)))
           (push (cons '= vars) impl-pairs))
          (t ; neither in LHS nor RHS
           (setf vars (list (genVarName 'z) (genVarName 'z)))))
    (setf forall (append vars forall))
    (let* ((var1 (first vars))
          (var2 (if (length=1 vars) var1 (second vars))))
      (push var1 pred1)
      (push var2 pred2)))
  (dolist (impl-pair impl-pairs new-fd-set)
    (push (list 'forall forall (list '=> (list 'and pred1 pred2) impl-pair)) new-fd-set)))

```

## A.2.4 fol-sat.lisp

```

;; The 'expr' is the query and the 'constraints' are the functional
;; dependencies and constant constraints

```

```

(defun fol-sat? (expr constraints)
  (let ((kb (make-fol-kb)))
    (tell kb constraints)
    (if (unsatisfiable? kb (->stripped-cnf expr)) nil t)))

```

```

(defstruct fol-kb
  (positive-clauses (make-hash-table :test #'eq))
  (negative-clauses (make-hash-table :test #'eq)))

```

```

(defun tell (kb sentence)
  "Add a sentence to a FOL knowledge base."
  (let* ((cnf1 (->stripped-cnf sentence))
        (mods (get-modulators cnf1))
        (cnf (if mods (subst-bindings mods cnf1) cnf1)))
    (for each clause in cnf do
      (tell-stripped-clause kb clause))))

```

```

(defun unsatisfiable? (kb sos)
  "See if the kb is unsatisfiable."
  (prove-by-refutation kb
    sos
    #'(lambda ()
        (RETURN-FROM unsatisfiable? t))))

```

```

(defun get-cnf (kb)
  "Obtain the complete CNF, in stripped form, of the kb."

```



```

(let ((cnf '()))
  (maphash #'(lambda (key val)
    (declare (ignore key))
    (setf cnf (append val cnf)))
    (fol-kb-positive-clauses kb))
  (maphash #'(lambda (key val)
    (declare (ignore key))
    (setf cnf (append val cnf)))
    (fol-kb-negative-clauses kb))
  (remove-duplicates cnf :test #'equal)))

(defun possible-resolvers (kb literal)
  "Find clauses that might resolve with a clause containing literal."
  (if (eq (op literal) 'not)
      (gethash (op (arg1 literal)) (fol-kb-positive-clauses kb))
      (gethash (op literal) (fol-kb-negative-clauses kb))))

(defun clause-in-kb? (kb clause)
  (declare (optimize speed)) ; NEW
  (let ((literal (first clause)))
    (if (eq (op literal) 'not)
        (member clause
          (gethash (op (arg1 literal)) (fol-kb-negative-clauses kb))
          :test #'equal)
        (member clause
          (gethash (op literal) (fol-kb-positive-clauses kb))
          :test #'equal))))))

(defun tell-stripped-clause (kb clause)
  (unless (or (clause-in-kb? kb clause) (tautology? clause))
    (for each literal in clause do
      (if (eq (op literal) 'not)
          (push clause (gethash (op (arg1 literal))
                                (fol-kb-negative-clauses kb)))
          (push clause (gethash (op literal)
                                (fol-kb-positive-clauses kb)))))))

(defun prove-by-refutation (kb sos fn)
  "Try to prove that SOS is true (given KB) by resolution refutation."
  (declare (optimize speed)) ; NEW
  (if (get-modulators sos)
      (setf sos (subst-bindings (get-modulators sos) sos))
      (setf sos (sort (filter-tautologies sos) #'< :key #'length)))

  (let ((clause nil)
        (result nil)
        (step 0))
    (loop
      (setf step (+ 1 step))

      ;; Exhausted set of support (sos) - IS SATISFIABLE.
      (when (null sos) (RETURN nil))

      (if *sat-debug*
          (progn
            (format t "~%~% STEP ~d:" step)
            (format t "~% KB: ~s " (get-cnf kb))
            (format t "~% SET OF SUPPORT: ~s " sos)
            ))

      ;; Move clause from SOS to the usable KB

```

```

(setf clause (pop sos))
(tell-stripped-clause kb clause)

;; Process everything that resolves with CLAUSE
(for each literal in clause do
  (for each resolver in (possible-resolvers kb literal) do
    (for each candidate-literal in resolver do
      (let ((b (unify (negate-literal candidate-literal) literal)))
        (when b
          (setf result
            (factor
              (subst-bindings
                b
                (append
                  (remove-literal literal clause)
                  (remove-literal candidate-literal resolver)))
              )
            )
          ;; We generated the empty clause! - IS NOT SATISFIABLE
          (if (= (length result) 0)
              (funcall fn)

              (unless (or
                (member result sos :test #'equal)
                (tautology? result))
                (if (= (length result) 1)
                    (setf sos (insert result sos #'< :key #'length))
                    (tell-stripped-clause kb result))
                ))))))))

```

### A.2.5 fol-utils.lisp

```

(defun ->stripped-cnf (sentence)
  "Strips the 'and' and 'or's."
  (mapcar #'disjuncts (conjuncts (->cnf sentence))))

(defun remove-literal (literal clause)
  (cond ((null clause) '())
        ((equal literal (first clause))
         (remove-literal literal (rest clause)))
        (t (cons (first clause) (remove-literal literal (rest clause))))))

(defun negate-literal (literal)
  (if (eql (first literal) 'not) (second literal) '(not ,literal)))

(defun positive-literal (literal)
  (if (eql (first literal) 'not)
      (negate-literal literal)
      literal))

(defun insert (item list pred &key (key #'identity))
  (merge 'list (list item) list pred :key key))

(defun remove-equality-tautologies (clause)
  (cond ((null clause) '())
        ((equality-tautology? (first clause))
         (remove-equality-tautologies (rest clause)))
        (t (cons (first clause)
                  (remove-equality-tautologies (rest clause))))))

```

```

(defun filter-tautologies (cset)
  (cond ((null cset) '())
        (t (remove 'nil (cons (remove-equality-tautologies
                               (first cset)) (filter-tautologies (rest cset)))))))

(defun equality-tautology? (lit)
  (or
   (and (equal (op lit) '=)
        (equal (arg1 lit)
                (arg2 lit)))
   (and (equal (op lit) 'NOT)
        (equality-contradiction? (arg1 lit))))))

(defun equality-contradiction? (lit)
  (and (equal (op lit) '=)
       (not (variable? (arg1 lit)))
       (not (variable? (arg2 lit)))
       (not (skolem? (arg1 lit)))
       (not (skolem? (arg2 lit)))
       (not (equal (arg1 lit) (arg2 lit)))))

(defun tautology? (clause)
  "Is clause a tautology (something that is always true)?"
  (some #'(lambda (literal)
            (or (and (eq (op literal) 'not)
                    (member (arg1 literal) clause :test #'equal))
                (equality-tautology? literal)))
        clause))

(defun modulator? (clause)
  "Is it a unit clause that assigns a value to a Skolem constant? (e.g.(= +x 1)) "
  (if (and (= 1 (length clause))
           (equal (first (first clause)) '=)
           (skolem? (second (first clause))))
      t nil))

(defun get-modulators (cnf)
  (cond ((null cnf) '())
        ((modulator? (first cnf))
         (cons (make-binding (second (first (first cnf)))
                             (third (first (first cnf))))
               (get-modulators (rest cnf))))
        (t (get-modulators (rest cnf)))))

(defun remove-contra (clause)
  (cond ((null clause) '())
        ((and (equal (op (first clause)) 'not)
              (equality-tautology? (arg1 (first clause))))
         (remove-contra (rest clause)))
        (t (cons (first clause) (remove-contra (rest clause))))))

(defun factor (clause)
  (remove-duplicates
   (remove-contra clause)
   :test #'equal))

```



```

      (cons 'and (reverse innerConds))))))
    (when (eq '- (first outerCond))
      (setf sqc (list 'not sqc)))
    (push sqc allConds)))
  (push (changeSimpleCond innerCond) innerConds)
  (push (changeSimpleCond outerCond) allConds))))))

(defun getTypes (l)
  (let ((types (list (second l))))
    (oldvars (list (first l)))
    (dolist (outerCond (caddr l))
      (when (sqc? outerCond)
        (dolist (innerCond (third outerCond))
          (when (= 2 (length innerCond))
            (pushnew (first innerCond) types)
            (pushnew (second innerCond) oldvars))))))
      (values (reverse types) (reverse oldvars))))

(defun genVarName (&optional (prefix 'y))
  (intern (format nil "$A%D" prefix (incf *variable-counter*))))

(defun genVarMappings (types-oldvars)
  (let ((oldvar-attr-newvars '() )
        (type-newVars-list '() )
        (oldvar-newVars-list '() ))
    (flet ((doRelation (type-oldvar prefix)
      (let ((attr-newvar '() )
            (type-newVars '() ))
        (dolist (attr
                  (second (assoc (first type-oldvar) *schema*)))
          (progn
            (push (list (first type-oldvar) (reverse type-newVars)) type-newVars-list)
            (push (list (rest type-oldvar) (reverse type-newVars)) oldvar-newVars-list)
            (push (cons (rest type-oldvar) (reverse attr-newvar)) oldvar-attr-newvars)))
          (let ((newvar (genVarName prefix)))
            (push newvar type-newVars)
            (push (cons attr newvar) attr-newvar))))))
      (doRelation (first types-oldvars) 'x)
      (dolist (type-oldvar2 (rest types-oldvars) (values oldvar-attr-newvars type-newVars-list oldvar-newVars-list))
        (doRelation type-oldvar2 'y))))))

(defun setCond? (oneCond)
  (let ((op (first oneCond)))
    (or (equal op 'in)
        (equal op 'not_in))))

```

```

(defun setCmd2fol (sc)
  (let ((op (first sc))
        (var (second sc))
        (constants (third sc)))
    (if (equal 'in op)
        (cons 'or (mapcar #'(lambda (constant) (cons '= (cons var (list constant)))) constants))
        (cons 'and (mapcar #'(lambda (constant) (cons '<> (cons var (list constant)))) constants))))))

(defun remove-equiJoins-fol (fol-sentence)
  (cons (first fol-sentence)
        (cons (second fol-sentence)
              (remove-equiJoins-fol-sqcs (rest (rest fol-sentence))))))

(defun remove-equiJoins-fol-sqcs (fol-sqcs)
  (let ((fol-sqc (first fol-sqcs)))
    (cond ((null fol-sqcs) nil)
          (t
           (if (equal 'not (first fol-sqc))
               (cons
                (list 'not (remove-equiJoins-positive-fol-sqc (first (rest fol-sqc))))
                (remove-equiJoins-fol-sqcs (rest fol-sqcs)))
               (cons
                (remove-equiJoins-positive-fol-sqc fol-sqc)
                (remove-equiJoins-fol-sqcs (rest fol-sqcs)))))))

(defun remove-equiJoins-positive-fol-sqc (positive-fol-sqc)
  (cond ((sqc?-fol positive-fol-sqc)
         (let* ((joinConds-otherConds (splice-joinConds-otherConds-fol positive-fol-sqc))
                (joinConds (first joinConds-otherConds))
                (otherConds (second joinConds-otherConds))
                (exiVars (second positive-fol-sqc))
                (newExiVars-newOtherConds (remove-equiJoins-positive-fol-sqc-1
                                             exiVars joinConds otherConds)))
           (list 'exists (first newExiVars-newOtherConds)
                 (cons 'and (second newExiVars-newOtherConds))))
        (t positive-fol-sqc)))

(defun splice-joinConds-otherConds-fol (positive-fol-sqc)
  (splice-joinConds-otherConds-fol-1 (rest (third positive-fol-sqc)) nil nil))

(defun splice-joinConds-otherConds-fol-1 (unsplicedConds joinConds otherConds)
  (let ((unsplicedCond (first unsplicedConds)))
    (cond ((null unsplicedConds) (list joinConds otherConds))
          (t
           (let ((unsplicedConds (rest unsplicedConds))
                 (unsplicedCond (first unsplicedConds)))
             (splice-joinConds-otherConds-fol-1 (rest (third positive-fol-sqc)) nil nil))))))

```

```

((join-cond?-dc unsplicedCond)
 (splice-joinConds-otherConds-fo1-1
  (rest unsplicedConds) (cons unsplicedCond joinConds) otherConds))
(t
 (splice-joinConds-otherConds-fo1-1
  (rest unsplicedConds) joinConds (cons unsplicedCond otherConds))))))
(defun remove-equi-joins-positive-fo1-sqc-1 (exiVars joinConds otherConds)
 (cond ((null joinConds) (list exiVars otherConds))
       (t
        (let* ((joinCond (order-joinCond (first joinConds)))
               (newJoinConds (subst (second joinCond) (third joinCond) (rest joinConds)))
               (newOtherConds (subst (second joinCond) (third joinCond) otherConds))
               (newExiVars (remove (third joinCond) exiVars)))
              (remove-equi-joins-positive-fo1-sqc-1 newJoinConds newOtherConds))))))
(defun order-joinCond (joinCond)
 (cond ((free-dc-var? (second joinCond)) joinCond)
       ((free-dc-var? (third joinCond)) (list (first joinCond) (third joinCond)
                                               (second joinCond)))
       (t joinCond)))
(defun dc-var? (x)
 "Is x a variable (a symbol starting with $)?"
 (and (symbolp x) (eql (char (symbol-name x) 0) #)))
(defun free-dc-var? (x)
 (and (dc-var? x)
      (and (symbolp x) (eql (char (symbol-name x) 1)
                             #\X))))

```

## A.2.7 normal.lisp

```

;;; Convert Expressions to CNF

;;; Top-Level Functions

(defun ->cnf (p &optional vars)
  (setf p (eliminate-implications p))
  (case (op p)
    (NOT (let ((p2 (move-not-inwards (arg1 p))))
            (if (literal-clause? p2) p2 (->cnf p2 vars))))
    (AND (conjunction (mappend #'(lambda (q) (conjunction (->cnf q vars)))
                               (args p))))
    (OR (merge-disjuncts (mapcar #'(lambda (q) (->cnf q vars))
                                (args p))))
    (FORALL (let ((new-vars (mapcar #'new-variable (mklist (arg1 p))))
                  (->cnf (sublis (mapcar #'cons (mklist (arg1 p)) new-vars)
                              (arg2 p))
                          (append new-vars vars))))
            (EXISTIS (->cnf (skolemize (arg2 p) (arg1 p) vars) vars))
    (t p) ; p is atomic
  ))

;;; Auxiliary Functions

(defun eliminate-implications (p)
  (if (literal-clause? p)
      p
      (case (op p)
        (=> '(or ,(arg2 p) (not ,(arg1 p))))
        (<=> '(and (or ,(arg1 p) (not ,(arg2 p)))
                  (or (not ,(arg1 p)) ,(arg2 p))))
        (t (cons (op p) (mapcar #'eliminate-implications (args p))))))

(defun move-not-inwards (p)
  "Given P, return ~P, but with the negation moved as far in as possible."
  (case (op p)
    (TRUE 'false)
    (FALSE 'true)
    (NOT (arg1 p))
    (AND (disjunction (mapcar #'move-not-inwards (args p))))
    (OR (conjunction (mapcar #'move-not-inwards (args p))))
    (FORALL (make-exp 'EXISTS (arg1 p) (move-not-inwards (arg2 p))))
    (EXISTS (make-exp 'FORALL (arg1 p) (move-not-inwards (arg2 p))))
    (t (make-exp 'not p))))

(defun merge-disjuncts (disjuncts)
  "Return a CNF expression for the disjunction."
  ;; The argument is a list of disjuncts, each in CNF.
  ;; The second argument is a list of conjuncts built so far.
  (case (length disjuncts)
    (0 'false)
    (1 (first disjuncts))
    (t (conjunction
        (let ((result nil))
          (for each y in (conjuncts (merge-disjuncts (rest disjuncts))) do
            (for each x in (conjuncts (first disjuncts)) do
              (push (disjunction (append (disjuncts x) (disjuncts y))
                                result)))
            (nreverse result))))))

```



```

(defun skolemize (p vars outside-vars)
  "Within the proposition P, replace each of VARS with a skolem constant,
  or if OUTSIDE-VARS is non-null, a skolem function of them."
  (sublis (mapcar #'(lambda (var)
                    (cons var (if (null outside-vars)
                                (skolem-constant var)
                                (cons (skolem-constant var) outside-vars))))
          (mklist vars))
    p))

(defun skolem-constant (name)
  "Return a unique skolem constant, a symbol starting with '+'."
  (intern (format nil "+^A_~D" name (incf *new-variable-counter*))))

(defun renaming? (p q &optional (bindings +no-bindings+))
  "Are p and q renamings of each other? (That is, expressions that differ
  only in variable names?)"
  (cond ((eq bindings +fail+) +fail+)
        ((equal p q) bindings)
        ((and (consp p) (consp q))
         (renaming? (rest p) (rest q)
                    (renaming? (first p) (first q) bindings)))
        ((not (and (variable? p) (variable? q)))
         +fail+)
        ;; P and Q are both variables from here on
        ((and (not (get-binding p bindings)) (not (get-binding q bindings)))
         (extend-bindings p q bindings))
        ((or (eq (lookup p bindings) q) (eq p (lookup q bindings)))
         bindings)
        (t +fail+)))

;;; Utility Predicates and Accessors

(defconstant +logical-connectives+ '(and or not => <=>)) (defconstant +logical-quantifiers+ '(forall
exists))

(defun atomic-clause? (sentence)
  "An atomic clause has no connectives or quantifiers."
  (not (or (member (op sentence) +logical-connectives+)
          (member (op sentence) +logical-quantifiers+))))

(defun literal-clause? (sentence)
  "A literal is an atomic clause or a negated atomic clause."
  (or (atomic-clause? sentence)
      (and (negative-clause? sentence) (atomic-clause? (arg1 sentence)))))

(defun negative-clause? (sentence)
  "A negative clause has NOT as the operator."
  (eq (op sentence) 'not))

(defun conjuncts (sentence)
  "Return a list of the conjuncts in this sentence."
  (cond ((eq (op sentence) 'and) (args sentence))
        ((eq sentence 'true) nil)
        (t (list sentence))))

(defun disjuncts (sentence)
  "Return a list of the disjuncts in this sentence."
  (cond ((eq (op sentence) 'or) (args sentence))
        ((eq sentence 'false) nil)

```

```

      (t (list sentence))))

(defun conjunction (args)
  "Form a conjunction with these args."
  (case (length args)
    (0 'true)
    (1 (first args))
    (t (cons 'and args))))

(defun disjunction (args)
  "Form a disjunction with these args."
  (case (length args)
    (0 'false)
    (1 (first args))
    (t (cons 'or args))))

```

## A.2.8 unify.lisp

```

(defconstant +fail+ nil "Indicates unification failure")

(defvar *new-variable-counter* 0)

(defun new-variable (var)
  "Create a new variable. Assumes user never types variables of form $X.9"
  (concat-symbol (if (variable? var) "" "$")
    var ".") (incf *new-variable-counter*))

(defconstant +no-bindings+ '((nil))
  "Indicates unification success, with no variables.")

;;; Top Level Functions

(defun unify (x y &optional (bindings +no-bindings+))
  "See if x and y match with given bindings. If they do,
  return a binding list that would make them equal [p 303]."
  (cond ((eq bindings +fail+) +fail+)
        ((equal x y) bindings)
        ((variable? x) (unify-var x y bindings))
        ((variable? y) (unify-var y x bindings))
        ((and (consp x) (consp y))
         (unify (rest x) (rest y)
                (unify (first x) (first y) bindings)))
        (t +fail+)))

(defun rename-variables (x)
  "Replace all variables in x with new ones."
  (sublis (mapcar #'(lambda (var) (make-binding var (new-variable var)))
    (variables-in x))
    x))

;;; Auxiliary Functions

(defun unify-var (var x bindings)
  "Unify var with x, using (and maybe extending) bindings [p 303]."
  (cond ((get-binding var bindings)
         (unify (lookup var bindings) x bindings))
        ((and (variable? x) (get-binding x bindings))
         (unify var (lookup x bindings) bindings))
        ((occurs-in? var x bindings)
         (t +fail+))))

```

```

        +fail+)
      (t (extend-bindings var x bindings))))

(defun variable? (x)
  "Is x a variable (a symbol starting with $)?"
  (and (symbolp x) (eql (char (symbol-name x) 0) #)))

(defun skolem? (x)
  "Is x a variable (a symbol starting with +)?"
  (and (symbolp x) (eql (char (symbol-name x) 0) #\+)))

)))

(defun get-binding (var bindings)
  "Find a (variable . value) pair in a binding list."
  (assoc var bindings))

(defun binding-var (binding)
  "Get the variable part of a single binding."
  (car binding))

(defun binding-val (binding)
  "Get the value part of a single binding."
  (cdr binding))

(defun make-binding (var val) (cons var val))

(defun lookup (var bindings)
  "Get the value part (for var) from a binding list."
  (binding-val (get-binding var bindings)))

(defun extend-bindings (var val bindings)
  "Add a (var . value) pair to a binding list."
  (cons (make-binding var val)
        ; Once we add a "real" binding,
        ; we can get rid of the dummy +no-bindings+
        (if (eql bindings +no-bindings+)
            nil
            bindings)))

(defun occurs-in? (var x bindings)
  "Does var occur anywhere inside x?"
  (cond ((eql var x) t)
        ((and (variable? x) (get-binding x bindings)
              (occurs-in? var (lookup x bindings) bindings))
         ((consp x) (or (occurs-in? var (first x) bindings)
                       (occurs-in? var (rest x) bindings)))
         (t nil)))

(defun subst-bindings (bindings x)
  "Substitute the value of variables in bindings into x,
taking recursively bound variables into account."
  (cond ((eql bindings +fail+) +fail+)
        ((eql bindings +no-bindings+) x)
        ((and (or (variable? x) (skolem? x))
              (get-binding x bindings)
              (subst-bindings (lookup x bindings)))
         ((atom x) x)

```

```
(t (reuse-cons (subst-bindings bindings (car x))
              (subst-bindings bindings (cdr x))
              x)))

(defun unifier (x y)
  "Return something that unifies with both x and y (or fail)."  
  (subst-bindings (unify x y) x))

(defun variables-in (exp)
  "Return a list of all the variables in EXP."  
  (unique-find-anywhere-if #'variable? exp))

(defun unique-find-anywhere-if (predicate tree &optional found-so-far)
  "Return a list of leaves of tree satisfying predicate,  
  with duplicates removed."  
  (if (atom tree)
      (if (funcall predicate tree)
          (pushnew tree found-so-far)
          found-so-far)
      (unique-find-anywhere-if  
        predicate  
        (first tree)  
        (unique-find-anywhere-if predicate (rest tree)  
                                  found-so-far))))))

(defun find-anywhere-if (predicate tree)
  "Does predicate apply to any atom in the tree?"  
  (if (atom tree)
      (funcall predicate tree)
      (or (find-anywhere-if predicate (first tree))
          (find-anywhere-if predicate (rest tree)))))
```