

Towards an Algebraic Theory of Information Integration

Gösta Grahne and Victoria Kiricenko
Department of Computer Science, Concordia University
Montreal, Quebec, Canada H3G 1M8
{grahne,kiricen}@cs.concordia.ca

Abstract

Information integration systems provide uniform interfaces to varieties of heterogeneous information sources. For query answering in such systems, the current generation of query answering algorithms in local-as-view (source-centric) information integration systems all produce what has been thought of as "the best obtainable" answer, given the circumstances that the source-centric approach introduces incomplete information into the virtual global relations. However, this "best obtainable" answer does not include all information that can be extracted from the sources because it does not allow partial information. Neither does the "best obtainable" answer allow for composition of queries, meaning that querying a result of a previous query will not be equivalent to the composition of the two queries.

In this paper we provide a foundation for information integration, based on the algebraic theory of incomplete information. Our framework allows us to define the semantics of partial facts and introduce the notion of the exact answer—that is the answer that includes partial facts. We show that querying under the exact answer semantics is compositional. We also present two methods for actually computing the exact answer. The first method is tableau-based, and it is a generalization of the "inverse-rules" approach. The second, much more efficient method, is a generalization of the rewriting approach, and it is based on partial containment mappings introduced in the paper.

1 Introduction

In the mid eighties the first author, together with Serge Abiteboul, interested Paris Kanellakis in incomplete information in relational databases. Paris immediately realized the importance of the problem, and the profoundness of the algebraic foundation laid by Imielinski and Lipski in their landmark paper [IL84]. Paris also saw that an understanding of the complexity theoretic aspects, in particular data complexity, was lacking. This resulted in the paper [AKG91], which has become the standard reference for the computational complexity of incomplete information.

In this paper we show that the problem of incomplete information is still highly relevant, due to its application to information integration. By extending the classical framework, we are able to provide a solid algebraic foundation for reasoning about information integration. As an additional benefit, our approach allows an important extension of the capabilities of these applications, namely to provide partial answers to user queries. These partial answers are important, not only because they are more informative than the types of answers considered before, but also because they enable query evaluation to be compositional. Without compositional query semantics, a user would not

be able to “search within the result,” *i.e.* pose a second query on the result of a first query, nor would it be possible to have a view mechanism.

To illustrate the problem and to make this discussion more concrete let us consider a simple example. Suppose we want to integrate three sources, let us call them S_1 , S_2 , and S_3 . The sources S_1 and S_3 are Internet movie databases, and source S_2 is a Hollywood gossip website that provides information about the cohabitation of movie stars. A small, but illustrative example of the sources is given below.

- Source S_1 , an Internet movie database about movie stars.

| Actor | Origin | Domicile |
|-------------------|---------|-------------------------|
| Gloria Swanson | U.S. | Sunset Blvd., Hollywood |
| Eric von Stroheim | Germany | Sunset Blvd., Hollywood |
| Greta Garbo | Sweden | Manhattan, New York |

- Source S_2 , a Hollywood gossip database about cohabitation of movie stars.

| Actor1 | Type1 | Actor2 | Type2 |
|------------------|-------------|----------------|-------------|
| Britt Ekland | Comedienne | Peter Sellers | Comedian |
| Elisabeth Taylor | Tragedienne | Richard Burton | Tragic hero |

- Source S_3 , an Internet movie database about affiliations of movie stars.

| Actor | Affiliation |
|-------------------|-------------|
| Gloria Swanson | UA |
| Eric von Stroheim | UA |
| Britt Ekland | UA |
| Peter Sellers | UA |
| Charlie Chaplin | UA |
| Greta Garbo | MGM |

A plausible interpretation of this source collection is to stipulate the global schema as

- **Star**(Name, Origin, Type, Domicile), that is actor’s name, role type, and address; and
- **Affiliation**(Actor, Studio), that is information of the studio of the star.

The relationship between the global schema and the schemas of S_1 , S_2 , and S_3 could then be expressed by the following definitions:¹

$$\begin{aligned}
 S_1(\text{Actor}, \text{Origin}, \text{Domicile}) &\leftarrow \text{Star}(\text{Actor}, \text{Origin}, \text{Type}, \text{Domicile}). \\
 S_2(\text{Actor}_1, \text{Type}_1, \text{Actor}_2, \text{Type}_2) &\leftarrow \text{Star}(\text{Actor}_1, \text{Origin}_1, \text{Type}_1, \text{Domicile}), \\
 &\quad \text{Star}(\text{Actor}_2, \text{Origin}_2, \text{Type}_2, \text{Domicile}). \\
 S_3(\text{Actor}, \text{Studio}) &\leftarrow \text{Affiliation}(\text{Actor}, \text{Studio}).
 \end{aligned}$$

¹Note that we are using the positional version of the relational model in our queries.

Suppose now the user issues the query

$$Q_1(\textit{Actor}, \textit{Origin}, \textit{Type}, \textit{Domicile}) \leftarrow \textit{Star}(\textit{Actor}, \textit{Origin}, \textit{Type}, \textit{Domicile}), \\ \textit{Affiliation}(\textit{Actor}, 'UA').$$

That is, the user is interested in obtaining all available information about stars affiliated with the United Artists studio. In all current generation information integration systems, such as [LRO96, LMSS95], the user would be returned an empty answer. This is because these systems do not have a way to get tuples for the subgoal *Star*, and would produce an empty rewriting.

Everybody who has ever queried integrated information, especially in the case of the World Wide Web, knows this situation very well. The next logical action of the user, who wants to get at least partial information, is to try to make his/her query less restrictive. In our example, to get at least some information about the stars affiliated with United Artists, the user has to modify the original query by projecting out some of the attributes of the relation *Star*. Since the user is unaware of the internals or the system, he/she has to try all possible combinations of the attributes and then manually assemble the final answer. From the point of view of the user, the system should take on this burden and compute the answer containing this partial information. This is feasible, since the query could be rewritten as the union of the following unsafe conjunctive queries.

$$Q_1(\textit{Actor}, \textit{Origin}, \textit{Type}, X) \leftarrow S_1(\textit{Actor}, \textit{Origin}, \textit{Type}), S_3(\textit{Actor}, 'UA'). \\ Q_1(\textit{Actor}_1, X, \textit{Type}_1, Y) \leftarrow S_2(\textit{Actor}_1, \textit{Type}_1, \textit{Actor}_2, \textit{Type}_2), S_3(\textit{Actor}_1, 'UA'). \\ Q_1(\textit{Actor}_2, X, \textit{Type}_2, Y) \leftarrow S_2(\textit{Actor}_1, \textit{Type}_1, \textit{Actor}_2, \textit{Type}_2), S_3(\textit{Actor}_2, 'UA').$$

The unrestricted variables *X* and *Y* represent unknown values. The answer can then be presented to the user as a table with some values missing, for our example as the table below.

| Actor | Origin | Type | Domicile |
|-------------------|---------|------------|-------------------------|
| Gloria Swanson | U.S. | ⊥ | Sunset Blvd., Hollywood |
| Eric von Stroheim | Germany | ⊥ | Sunset Blvd., Hollywood |
| Britt Ekland | ⊥ | Comedienne | ⊥ |
| Peter Sellers | ⊥ | Comedian | ⊥ |

Producing such partial answers is not without its intricacies. Since the missing values obviously represent nulls of the type “unknown,” there is a connection to incomplete information. This connection was first utilized by Abiteboul and Duschka [AD98], who used the conditional tables of Imielinski and Lipski [IL84], as extended in [Gra84], as a tool to obtain an effective query evaluation mechanism and complexity theoretic results for information integration. Notably, the complexity results for non-recursive queries in [AD98] are closely related to those in [AKG91]. Here we see again an example of the importance of incomplete information, and Abiteboul subsequently argues in his PODS 1999 invited talk [Abi99] that the problem of information integration “should be approached with an incomplete information model,” and he recalls Imielinski’s and Lipski’s contribution as “*the* model for incomplete information.”

However, in [AD98], Abiteboul and Duschka only use the *certain answer* aspect of incomplete information and conditional tables. The framework of Imielinski and Lipski is much richer than

that, and we shall see below that a notion called the *exact answer* will play a crucial role in extending the theory of information integration to take partial information into account.

First of all, the answer above containing nulls is essentially the exact answer, an alternative semantic to the certain answer. Note that all previous approaches, including [AD98], explicitly or implicitly use the certain answer. Second, the exact answer allows for composition of queries. One of the important benefits of compositionality of queries is that it allows the user to “search within the results.” For example, in a web based setting, the number of tuples in an answer typically is much larger than the user wants. In such a case, to reduce the size of the answer, the user would most likely want to make his/her query more restrictive. This means that an information integration system, in order to be of any practical use, has to provide the user with the ability to search within the answers.

Lets assume that in our example the user is interested in knowing about the cohabitation of United Artists stars. In current generation information integration systems the user would have to write the query as

$$Q_2(Actor_1, Actor_2) \leftarrow Star(Actor_1, Origin_1, Type_1, Domicile), Affiliation(Actor_1, 'UA'), \\ Star(Actor_2, Origin_2, Type_2, Domicile), Affiliation(Actor_2, 'UA').$$

Since we just saw that we were able to provide the user with partial information about the join of the two first (and two last) subgoals in terms of the relation Q_1 , containing nulls, returned before. It is easy to see that Q_2 is equivalent to

$$Q'_2(Actor_1, Actor_2) \leftarrow Q_1(Actor_1, Origin_1, Type_1, Domicile), \\ Q_1(Actor_2, Origin_2, Type_2, Domicile).$$

From the users point of view it is obvious that he/she should “search within the results,” in other words, issue query Q'_2 .

Can we use the results of the previous query in the computation of the answer? The answer is yes but only if the system remembers that some of the null values in this result represent the same unknown values. In particular, the unknown domiciles of Britt Ekland and Peter Sellers are the *same* domicile. Intuitively, this is because only source S_2 mentions these actors, and according to the definition of S_2 they would live at the same unknown address. Therefore, the information integration system will now have to account for this fact in some way, for instance by presenting the answer with marked nulls instead.

| Actor | Origin | Type | Domicile |
|-------------------|-----------|------------|-------------------------|
| Gloria Swanson | U.S. | \perp_1 | Sunset Blvd., Hollywood |
| Eric von Stroheim | Germany | \perp_2 | Sunset Blvd., Hollywood |
| Britt Ekland | \perp_3 | Comedienne | \perp_4 |
| Peter Sellers | \perp_5 | Comedian | \perp_4 |

Then we could evaluate Q'_2 of the result of Q_1 above, and obtain the table below. Note that Britt Ekland and Peter Sellers occur in the output since they join on the value \perp_4 in the fourth column. For an exposition on marked nulls, see [Ull89].

| Actor1 | Actor2 |
|----------------|-------------------|
| Gloria Swanson | Eric von Stroheim |
| Britt Ekland | Peter Sellers |

The question that now arises is how these marked nulls should be accounted for.

The rest of this paper is organized as follows. In Section 2 we review the algebraic approach to incomplete information. In Section 3 we describe how querying in information integration systems amounts querying incomplete databases, producing partial facts. In Section 4 we extend the notion of rewriting to take into account the partial facts, and in Section 5 we give a unification based method for computing the extended rewritings. Section 6 contains the conclusions.

2 Basic Definitions and Results

Relational Databases.

The Relational Model structures information in the form of tables. There are various ways to formalize this notion. The first is a direct adaptation of the mathematical concept of a relation: Let the atomic values come from a countably infinite set **dom**.² Let k be a positive integer. A k -ary relation is then a finite subset of $\mathbf{dom} \times \mathbf{dom} \times \dots \times \mathbf{dom}$, k -times, sometimes denoted \mathbf{dom}^k . In order to have names for relations, such as R , R_1 , R_2 , etc, we assume a function **db** that when applied to a relation name R gives the *instance* $\mathbf{db}(R)$ of the actual tuples from \mathbf{dom}^k . For a relational “schema” $\{R_1, R_2, \dots, R_n\}$, the database instance consists of an instance $\mathbf{db}(R_i)$ for each relation name R_i , $i \in \{1, \dots, n\}$. If the schema is understood, we sometimes write simply **db**. Note that there are no column names in this formalization, only positions (first column, second column, etc).

In this paper we adopt a variation of this modeling influenced by the logic programming approach to databases. From the domain **dom** and the relation names we build up a (Herbrand) universe consisting of all expressions of the form $R(a_1, a_2, \dots, a_k)$, where R is a relation name and the a_i ’s are values in **dom**. Such expressions are called *facts*. A database instance is then simply a finite set of facts, i.e. a finite subset of the universe, such as $\{R_1(a_1, a_3), R_1(a_2, a_3), R_2(a_3, a_4)\}$. Since relation names are part of the facts, we do not need to list facts from different relations separately.

For more information on the various ways to formalize the relational model, see [AHV95].

A *query* is an expression that defines a function from all databases to (usually) a single relation. In this paper we focus on a simple but the most practical class of queries, namely the conjunctive queries. For this we need the concept of an *atom*. An atom is like a fact, except that we allow *variables*, taken from an infinite set **var**, as well as constants. For instance, $R(a, X)$ is an atom, whereas $R(a, b)$ is a fact.³ The variables in the atoms are placeholders, and they stand for any domain value. Variables will be denoted by uppercase letters, such as X and Y .

A *conjunctive query* φ is an expression of the form

$$\mathit{head}(\varphi) \leftarrow \mathit{body}(\varphi),$$

²We assume for simplicity and clarity of exposition that there is only one domain from which all values are drawn.

³ $R(a, b)$ is of course also an atom, since a fact is a special case of an atom.

where $body(\varphi)$ is a set of atoms over relation names in the database schema, and $head(\varphi)$ is an atom over an answer relation name not used in the database. We assume that all variables occurring in $head(\varphi)$ also occur in $body(\varphi)$, i. e. that the query φ is *safe*.

A conjunctive query φ can be applied to a database \mathbf{db} resulting in a set of facts

$$\varphi(\mathbf{db}) = \{\sigma(head(\varphi)) : \sigma(body(\varphi)) \subseteq \mathbf{db} \text{ for some valuation } \sigma\}.$$

A *valuation* σ , is formally a finite partial mapping from $\mathbf{var} \cup \mathbf{dom}$ to \mathbf{dom} that is the identity on \mathbf{dom} .

Example 1 If the database \mathbf{db} is $\{R_1(a, b), R_1(c, b), R_2(b, d)\}$, and φ is $Q(X) \leftarrow R_1(X, Y), R_2(Y, Z)$, then $\varphi(\mathbf{db}) = \{Q(a), Q(c)\}$. On the other hand, if φ is $Q(X) \leftarrow R_1(a, Y), R_2(Y, X)$, then $\varphi(\mathbf{db}) = \{Q(d)\}$. ■

Incomplete Databases and Tableaux.

Here we shall briefly review Imielinski's and Lipski's algebraic approach to incomplete information [IL84]. We shall present the approach in our uniform framework.

Usually a database is a complete description of the state of the world modeled by it. This is actually true only if we adopt the *Closed World Assumption*, CWA [Rei78], according to which facts not explicitly stored in the database are false. For example, in the database $\{R_1(a_1, a_3), R_1(a_2, a_3), R_2(a_3, a_4)\}$, the fact $R_1(a_3, a_1)$ is false. The closed world assumption is convenient since there in general is an infinitude of false facts. The CWA is however not appropriate for modeling all situations. We might for instance have a database $\mathbf{db} = \{Affiliation(Greta Garbo, MGM), Affiliation(Charlie Chaplin, UA)\}$. If we consider the fact $Affiliation(Elisabeth Taylor, UA)$ we might not want to draw the conclusion that Elisabeth Taylor is not affiliated with UA, we just don't have evidence recorded about it. We then adopt the *Open World Assumption*, OWA, which regards the database as an *incomplete* description of the world: All facts stored in the database are true, the truth value of any other fact is *unknown*. Semantically this means that the stored database \mathbf{db} actually is a finite description of a *set of possible worlds*, defined as

$$\{\mathbf{db}' : \mathbf{db} \subseteq \mathbf{db}'\}.$$

Each \mathbf{db}' represents one possible complete state of affairs (or world). Thus for instance the fact $Affiliation(Elisabeth Taylor, UA)$ will be true in some possible worlds, and false in others. Facts that are true in *some* possible worlds are called *possible* facts, and facts true in *all* possible worlds are called *certain* facts.

A database under the OWA is the simplest example of an incomplete database. To go a step further, suppose we know that Elisabeth Taylor is affiliated with some studio, but we don't know which one. This could be represented by the atom $Affiliation(Elisabeth Taylor, X)$. Here again, we have an incomplete description of the world. The world could correspond to $Affiliation(Elisabeth Taylor, MGM)$, or to $Affiliation(Elisabeth Taylor, UA)$, and so on.

A very simple (and efficient) way to represent possible worlds is to allow the database to consist of atoms, as opposed to pure facts only. A set of atoms is called a *tableau* [Men84], also known as *naïve tables* [IL84], and *equality tables* [AKG91]. A tableau is denoted T , as opposed to \mathbf{db} which stands for a finite set of facts (not atoms).

Example 2 Let $T = \{\text{Affiliation}(\text{Greta Garbo}, \text{MGM}), \text{Affiliation}(\text{Elisabeth Taylor}, X)\}$. This tableau contains two atoms, the first of which actually is a fact (a special case of an atom). ■

A tableau T represents a set of databases. This set is denoted $\text{rep}(T)$, and it is defined by

$$\text{rep}(T) = \{\mathbf{db} : \text{there is a valuation } \sigma \text{ such that } \sigma(T) \subseteq \mathbf{db}\}.$$

The definition says that a database \mathbf{db} is represented by a tableau T , if there is a valuation σ such that when all variables in T are replaced by their image under σ , the set of facts thus obtained is a subset of \mathbf{db} .

Example 3 In our tableau $T = \{\text{Affiliation}(\text{Greta Garbo}, \text{MGM}), \text{Affiliation}(\text{Elisabeth Taylor}, X)\}$, we will for instance have $\{\text{Affiliation}(\text{Greta Garbo}, \text{MGM}), \text{Affiliation}(\text{Elisabeth Taylor}, \text{MGM})\} \in \text{rep}(T)$, $\{\text{Affiliation}(\text{Greta Garbo}, \text{MGM}), \text{Affiliation}(\text{Elisabeth Taylor}, \text{UA})\} \in \text{rep}(T)$, etc. According to $\text{rep}(T)$, $\text{Affiliation}(\text{Greta Garbo}, \text{MGM})$ is a certain fact since it is true in all possible worlds in $\text{rep}(T)$, and $\text{Affiliation}(\text{Elisabeth Taylor}, \text{MGM})$ is a possible fact since it is true in some possible worlds represented by T . ■

There are many other ways to represent sets of possible worlds, see [Gra84, IL84, AKG91], and [Gra91, AHV95] for an overview.

Querying incomplete databases and tableaux.

Now what does it mean to query an incomplete database? Since an incomplete database is a set, each element representing a possible database, the natural answer to a query is the *set* of answers obtained by applying the query on each possible database. If φ is a query, and \mathcal{X} is a set of possible databases, then

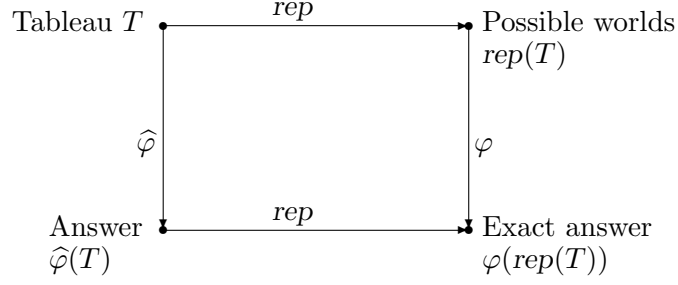
$$\varphi(\mathcal{X}) = \{\varphi(\mathbf{db}) : \mathbf{db} \in \mathcal{X}\}$$

This set of answers is called the *exact answer*.

Suppose now we have a query φ on an incomplete database represented by a tableau T . The exact answer would be $\varphi(\text{rep}(T))$. But just as the original database was represented by a tableau T we would like to represent the exact answer by another tableau U , such that $\text{rep}(U) = \varphi(\text{rep}(T))$. More practically thinking, given a query φ and a tableau T , we would like to have an algorithm, or evaluation mechanism, call it $\hat{\varphi}$, applicable on tableau, such that

$$\text{rep}(\hat{\varphi}(T)) = \varphi(\text{rep}(T)).$$

This requirement can be illustrated by the following commutative diagram.



Is the requirement in $rep(\widehat{\varphi}(T)) = \varphi(rep(T))$ achievable in general? The answer is yes, but only if we use a representation mechanism more complicated than tableaux (see [IL84]). If we wish to use tableaux we have to give up the strict equality requirement in the commutative diagram.

Since we cannot represent the exact answer to a query, it is natural to focus on the *certain answer* instead. The certain answer to a query φ consists of those tuples that are in the answer to φ , no matter which possible database we consider. The certain answer is denoted φ_* , and is defined formally as

$$\varphi_*(\mathcal{X}) = \bigcap_{\mathbf{db} \in \mathcal{X}} \varphi(\mathbf{db}).$$

Note that $\varphi_*(\mathcal{X})$ is a relation (a set of facts, no atoms). Then we could require instead that given a query φ and a tableaux T , there is an effectively constructible relation, call it $\varphi_*(T)$, such that

$$\varphi_*(T) = \bigcap_{\mathbf{db} \in rep(T)} \varphi(\mathbf{db}).$$

However, as illustrated in the introduction, if we return the relation $\varphi_*(T)$, then we cannot use the result for subsequent querying. In other words, there are (conjunctive) queries φ and ψ , and tableaux T , such that

$$(\varphi \circ \psi)_*(T) \neq \varphi_*(\psi_*(T)).$$

This also means that we could not have a query evaluation mechanism, found in all real query engines, that would operate in a uniform recursive fashion, i.e. $\varphi \circ \psi$ could be evaluated by first evaluating ψ , and then φ .

The reason that the equality $(\varphi \circ \psi)_*(T) = \varphi_*(\psi_*(T))$ fails to hold is that the intermediate result contains some *partial facts* that are discarded by the certain answer of ψ , and that nevertheless contributes to some certain facts found in $\varphi \circ \psi$.

If we like to achieve uniform evaluation and composability of queries we, thus, have to preserve the partial facts in the answer. We shall, therefore, require that given a query φ and a tableau T , there is a tableaux $\widehat{\varphi}(T)$, such that

$$\text{Truth Preservation} \quad \bigcap rep(\widehat{\varphi}(T)) = \bigcap \varphi(rep(T))$$

$$\text{Query Compositionality} \quad \widehat{\psi}(\widehat{\varphi}(T)) = \widehat{\psi \circ \varphi}(T)$$

for all queries ψ applicable at a result of φ .

Lipski, in an all but forgotten paper [Lip84], proved that for monotone query languages, a sufficient condition for truth preservation and compositionality is that

$$\text{Coinitiality} \quad \text{rep}(\widehat{\varphi}(T)) \approx \varphi(\text{rep}(T)),$$

where \approx -sign means that the two sets have the same \subseteq -minimal elements. In the sequel we will use this sufficient condition.

For tableaux, it turns out that if we, for evaluation purposes, treat the variables as pairwise distinct constants, and distinct from all “true” constants, we can apply standard evaluation and obtain a representation of the result that is coinitial with the exact answer. To formally explain the evaluation, we need the concept of a substitution. A *substitution* is a valuation, except that we allow variables to be mapped into variables, not only constants. Thus, a substitution θ is a partial function from $\mathbf{dom} \cup \mathbf{var}$ to $\mathbf{dom} \cup \mathbf{var}$, keeping in mind that constants have to be mapped to themselves. Then

$$\widehat{\varphi}(T) = \{\theta(\text{head}(\varphi)) : \theta(\text{body}(\varphi)) \subseteq T \text{ for some substitution } \theta\}.$$

Theorem 1 $\text{rep}(\widehat{\varphi}(T)) \approx \varphi(\text{rep}(T))$.

Proof. Let \mathbf{db} be a \subseteq -minimal element in $\text{rep}(\widehat{\varphi}(T))$. Then there exist a valuation σ such that $\sigma(\widehat{\varphi}(T)) = \mathbf{db}$. Let t be an arbitrary fact in \mathbf{db} . Then there is a fact $u \in \widehat{\varphi}(T)$ such that $\sigma(u) = t$, and there is a substitution θ such that $\theta(\text{body}(\varphi)) \subseteq T$ and $u = \theta(\text{head}(\varphi))$.

Let σ' be an extension of σ that maps every variable that is in T but not in $\widehat{\varphi}(T)$ to a distinct new constant. Since $\theta(\text{body}(\varphi)) \subseteq T$, we have $\sigma'\theta(\text{body}(\varphi)) \subseteq \sigma'(T)$. It now follows that $t = \sigma'(\theta(\text{head}(\varphi))) \in \varphi(\sigma'(T))$. Note that $\sigma'(T)$ is a \subseteq -minimal element in $\text{rep}(T)$. From the monotonicity of φ it follows that $\varphi(\sigma'(T))$ is a \subseteq -minimal element in $\varphi(\text{rep}(T))$. We have established that $\mathbf{db} \subseteq \varphi(\sigma'(T))$

That concludes the proof that any \subseteq -minimal element \mathbf{db} in $\text{rep}(\widehat{\varphi}(T))$ is also in $\varphi(\text{rep}(T))$.

For inclusion in the other direction let \mathbf{db} be a \subseteq -minimal element in $\varphi(\text{rep}(T))$. Then there is a valuation σ , such that $\mathbf{db} = \varphi(\sigma(T))$. Let t be an arbitrary tuple in \mathbf{db} . Then there is a valuation ρ , such that $t = \rho(\text{head}(\varphi))$ and all facts in $\rho(\text{body}(\varphi))$ are in $\sigma(T)$. Now we have two cases to consider.

Case 1: The valuation σ is one-to-one. Then there is an inverse σ^{-1} , and hence $\sigma^{-1}(\rho(\text{body}(\varphi))) \subseteq \sigma^{-1}(\sigma(T)) = T$, and consequently $\sigma^{-1}(t) = \sigma^{-1}(\rho(\text{head}(\varphi)))$ is in $\widehat{\varphi}(T)$. Since $\sigma(\widehat{\varphi}(T)) \in \text{rep}(\widehat{\varphi}(T))$, it follows that $t \in \sigma(\widehat{\varphi}(T)) \in \text{rep}(\widehat{\varphi}(T))$. Likewise, if t' is any other tuple in $\mathbf{db} = \varphi(\sigma(T))$, it is generated by some valuation ρ' , and we have $\sigma^{-1}(t') = \sigma^{-1}(\rho'(\text{head}(\varphi))) \in \widehat{\varphi}(T)$. Therefore $\mathbf{db} \subseteq \sigma(\widehat{\varphi}(T))$.

Case 2: There is (at least one) pair of distinct variables X and Y in T , such that $\sigma(X) = \sigma(Y)$. If $\sigma(X) = \rho(U)$, and $\sigma(Y) = \rho(W)$, for $U \neq W$, then the valuation ω , that is like $\sigma^{-1} \circ \rho$, except $\omega(U) = X$, and $\omega(V) = Y$, gives us $\omega(\text{body}(\varphi)) \subseteq T$, and $\sigma^{-1}(t) = \omega(\text{head}(\varphi)) \in \widehat{\varphi}(T)$. Consequently $t = \sigma(\sigma^{-1}(t)) \in \sigma(\widehat{\varphi}(T))$.

Suppose then that $\sigma(X) = \sigma(Y) = \rho(W)$, and that there are (at least) two occurrences of W in $\text{body}(\varphi)$. Consider now the valuation σ' , that is exactly like σ , except it maps Y to a fresh constant, say a . Clearly $t \notin \varphi(\sigma'(T))$, and any fact in $\varphi(\sigma'(T))$ is also in $\varphi(\sigma(T))$ (because there is an embedding of $\sigma'(T)$ into $\sigma(T)$.) Therefore we have a contradiction to the assumption that t belonged to a \subseteq -minimal element of $\varphi(\text{rep}(T))$. \blacksquare

Theorem 1 means that the diagram on page 8 commutes with respect to \approx , *i.e.* cointiality. Furthermore, the certain facts of $\varphi(\text{rep}(T))$, that is $\cap(\varphi(\text{rep}(T)))$, can be obtained from $\widehat{\varphi}(T)$ by retaining only the pure variable-free atoms. Moreover, the resulting tableau $\widehat{\varphi}(T)$ contains all the facts necessary for subsequent query evaluation.

Example 4 Let $T = \{R_1(a, b), R_1(d, X), R_2(b, c), R_2(X, e), R_2(Y, f)\}$, and $\varphi = Q(X, Y, Z) \leftarrow R_1(X, Y), R_2(Y, Z)$. Then $\widehat{\varphi}(T) = \{Q(a, b, c), Q(d, X, e)\}$. The only certain fact in the answer is $Q(a, b, c)$. If $\varphi = Q(X, b) \leftarrow R_1(X, b)$ we have $\widehat{\varphi}(T) = \{Q(a, b)\}$. This fact is also certain. ■

Theorem 1 was first discovered by Imielinski and Lipski in [IL84], and independently by Vardi in [Var86].

3 Global databases: The meaning of integrated information

Information Integration systems [Ull97, Hal00, Hal01, Len02] aim to provide a uniform query interface to multiple heterogeneous sources. One particular and useful way of viewing these systems, first proposed within the Information Manifold project [LRO96], is to postulate a *global schema* (called a world view) that provides a unifying data model for all the information sources. A query processor is in charge of accepting queries written in terms of this global schema, translating them to queries on the appropriate sources, and assembling the answers into a global answer. Each source is modeled as a *materialized view* defined in terms of the global relations, which are virtual. Note the reversal of the classical model: instead of thinking of views as virtual artifacts defined on stored relations, we think of the views as stored, and the relations that the views are defined on as virtual.⁴

A question of semantics now arises: what is the meaning of a query? Since a query is expressed in terms of the global schema, and the sources implicitly represent an instance of this global schema, it would be natural—at least conceptually—to reconstruct the global database represented by the views and apply the query to this global database.

Well, it turns out that the global database actually is incomplete. In other words, there might be *several* (usually infinitely many) global databases that are consistent with the definition of, and the data in the sources.

Example 5 For a simple example, suppose we have a global relations $Affiliation(Actor, Studio)$, and two sources: Source S_1 has definition $V_1(X, Y) \leftarrow Affiliation(X, Y)$ and extension $\mathbf{s}_1 = \{V_1(Greta Garbo, MGM)\}$. Source S_2 has definition $V_2(X) \leftarrow Affiliation(X, Y)$ and extension $\mathbf{s}_2 = \{V_2(Elisabeth Taylor)\}$. Then it is natural to think that any database that contains at least the facts $Affiliation(Greta Garbo, MGM)$, and $Affiliation(Elisabeth Taylor, a)$, for some $a \in \mathbf{dom}$, is a possible global database for $\mathcal{S} = \{S_1, S_2\}$. ■

Formally, for a source collection $\mathcal{S} = \{S_1, S_2, \dots, S_m\}$, where each S_i has as definition a conjunctive query ψ_i , and a finite set of facts \mathbf{s}_i over the relation name in $head(\psi_i)$ as extension, we call its set

⁴Whether the views are actually stored at each source or materialized by other means, and whether they consist of relations or semi-structured objects or files are issues irrelevant to our discussion and hidden from the query processor by appropriate wrappers.

of possible databases $\text{poss}(\mathcal{S})$, and define it as

$$\text{poss}(\mathcal{S}) = \{\mathbf{db} : \mathbf{s}_i \subseteq \psi_i(\mathbf{db}), i \in \{1, \dots, m\}\}.$$

Now the set $\text{poss}(\mathcal{S})$ can be conveniently represented by a tableau, denoted $T(\mathcal{S})$, over the relation names in the bodies of the definitions ψ_i , such that $\text{rep}(T) = \text{poss}(\mathcal{S})$. To construct T we shall follow the approach in [GM99] (see also [GK02, GK03]). We define a function, which, by abuse of notation, we also denote T , from source extensions \mathbf{s} , with definition ψ , to tableaux over the relation names in the body of ψ . We also need an auxiliary function refresh , that, when applied to a set of atoms, replaces each variable with a distinct unused (“fresh”) variable. Given a source $S = (\psi, \mathbf{s})$, we set

$$T(S) = \bigcup_{t \in \mathbf{s}} \{\text{refresh}(\sigma(\text{body}(\psi))) : \sigma(\text{head}(\psi)) = t\},$$

for some valuation σ .

Example 6 Let S have definition $\psi = V(X, Z) \leftarrow R_1(X, Y), R_2(Y, Z)$ and extension $\mathbf{s} = \{V(a, b), V(c, d)\}$. Then $T(S) = \{R_1(a, Y_1), R_2(Y_1, b), R_1(c, Y_2), R_2(Y_2, d)\}$, where Y_1 and Y_2 are fresh variables. ■

When there are several sources in \mathcal{S} we set

$$T(\mathcal{S}) = \bigcup_{S \in \mathcal{S}} T(S).$$

Example 7 For $\mathcal{S} = \{S_1, S_2\}$ in Example 5, we have $T(\mathcal{S}) = \{\text{Affiliation}(\text{Greta Garbo}, \text{MGM}), \text{Affiliation}(\text{Elisabeth Taylor}, X)\}$. ■

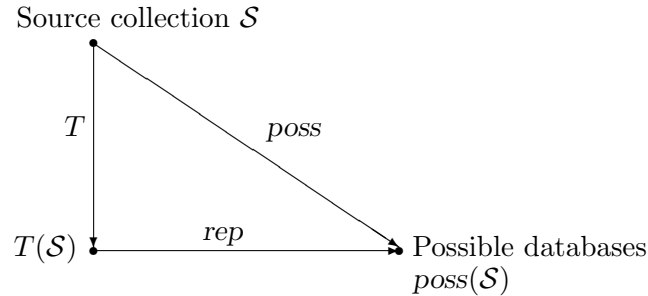
The tableau constructed by the function T has the following desirable property.

Theorem 2 $\text{rep}(T(\mathcal{S})) = \text{poss}(\mathcal{S})$.

Proof. Let $\mathbf{db} \in \text{rep}(T(\mathcal{S}))$. To prove that $\mathbf{db} \in \text{poss}(\mathcal{S})$ we need to show that for all sources $S_i = (\psi_i, \mathbf{s}_i)$ in \mathcal{S} , we have $\mathbf{s}_i \subseteq \psi_i(\mathbf{db})$. Since $\mathbf{db} \in \text{rep}(T(\mathcal{S}))$ there is a valuation σ such that $\sigma(T(\mathcal{S})) \subseteq \mathbf{db}$. Let $S_i = (\psi_i, \mathbf{s}_i)$ be an arbitrary source in \mathcal{S} , and let t be an arbitrary fact in \mathbf{s}_i . Then there must be a substitution θ , such that $t = \theta(\text{head}(\psi_i))$ and all facts in $\theta(\text{body}(\psi_i))$ are in $T(\mathcal{S})$. It follows that $\theta(\sigma(\text{body}(\psi_i))) \subseteq \mathbf{db}$ and, consequently, $\theta(\sigma(\text{head}(\psi_i))) \in \mathbf{db}$. Since $\theta(\sigma(\text{head}(\psi_i))) = t$, we have $\mathbf{s}_i \subseteq \psi_i(\mathbf{db})$ as desired.

For inclusion in the other direction, let $\mathbf{db} \in \text{poss}(\mathcal{S})$. From construction of $T(\mathcal{S})$ it immediately follows that there is a valuation σ such that $\sigma(T(\mathcal{S})) \subseteq \mathbf{db}$ and, thus, that $\mathbf{db} \in \text{rep}(T(\mathcal{S}))$. ■

Theorem 2 is illustrated in the following diagram.

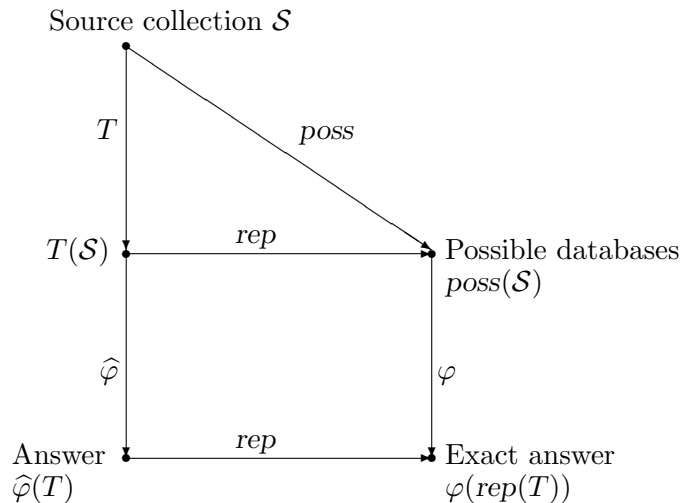


In definition of the set of possible databases we made a version of the OWA. The facts in the sources are considered to be *sound*, in the sense that they all are generated by facts in the global database. An alternative view would be to state that a source can contain (spurious) facts that do not stem from facts in the global database. In this case the source would be considered *complete*, but not necessarily sound. If the sources are required to be both sound and complete, they are called *exact*.⁵ The situations where sources are allowed to be both sound and complete are treated in [AD98, GM99]. The presence of complete sources however usually increase the computational complexity of decision problems related to integration, see [AD98, GM99]. In [MM01] the authors consider a generalization of the sound/complete view assumptions. Each source is considered to be both *partially sound*, and *partially complete*. A soundness of say 80% would say that the source contains 80% facts that are derived from facts in the global database, and 20% of “spurious” facts. A completeness of 80% says that 80% of the facts derivable from the global database are in the source.

Now that we have seen that a source collection \mathcal{S} actually defines an incomplete database $poss(\mathcal{S})$ that we can represent by a tableau $T(\mathcal{S})$, and that an incomplete database represented as a tableau can be queried using the $\widehat{\varphi}$ -evaluation of a query φ , we have a method for computing a representation of the exact answer to a query φ posed on the global database: First compute the tableau $T(\mathcal{S})$, then apply $\widehat{\varphi}$ to obtain $\widehat{\varphi}(T(\mathcal{S}))$. Theorem (1) gave us $rep(\widehat{\varphi}(T)) \approx \varphi(rep(T))$, for any tableau T , and since Theorem (2) gives $rep(T(\mathcal{S})) = poss(\mathcal{S})$, we get

Theorem 3 $rep(\widehat{\varphi}(T(\mathcal{S}))) \approx \varphi(poss(\mathcal{S}))$.

In other words, we combine the two previous commutative diagrams, and get the following illustration of Theorem 3.



Example 8 In Examples 5 and 7, if the query φ were the identity, i.e. it asked for the facts in the *Affiliation* -relation, a representation of the exact answer would be the tableau $\{Affiliation(Greta Garbo, MGM), Affiliation(Elisabeth Taylor, X)\}$. The certain fact in this answer is $Affiliation(Greta Garbo, MGM)$. ■

⁵If all sources are complete, the definition of $poss$ would be $poss(\mathcal{S}) = \{\mathbf{db} : \mathbf{s}_i \supseteq \psi_i(\mathbf{db}), i \in \{1, \dots, m\}\}$, and for exact sources we would have $poss(\mathcal{S}) = \{\mathbf{db} : \mathbf{s}_i = \psi_i(\mathbf{db}), i \in \{1, \dots, m\}\}$.

4 Using rewritings to answer queries on global databases

In the previous section we saw how to answer a query φ by first constructing the tableau representing all possible global databases. However, computing $T(\mathcal{S})$ might involve a lot of redundant work, since it amounts to constructing the tableau corresponding to the entire source collection \mathcal{S} , whereas the global relations that are in the body of φ might be mentioned in only few source definitions. Furthermore, the query might have selections and joins that could be computed directly at the sources.

Let $\{\psi_1, \dots, \psi_n\}$ be a set of source definitions, and φ a query whose body mentions relation names in $\{body(\psi_i)\}_i$. A *rewriting* of a query φ , with respect to $\{\psi_1, \dots, \psi_n\}$, is a query χ with the same head-relation name as φ , and body-relation names in $\{head(\psi_i)\}_i$.

The idea of a rewriting is that it can be evaluated on source extensions. It appears that rewritings were first proposed for a restricted setting in [YL87]. They were generalized in [LMSS95] and [LRO96]. The correctness criteria for rewritings were defined in terms of a certain containment between the rewriting and the original query.

Example 9 Let $\mathcal{S} = \{S_1, S_2\}$, with definitions $V_1(X, Z) \leftarrow R_1(X, Y), R_2(Y, Z)$, and $V_2(X, Z) \leftarrow R_1(X, Y), R_3(Y, Z)$. Let the query be $\varphi = Q(X, Y) \leftarrow R_1(X, Y)$. According to the correctness criteria in [LRO96], the desired rewriting φ' would be the union of $Q(X, Y) \leftarrow V_1(X, Y)$, and $Q(X, Y) \leftarrow V_2(X, Y)$. ■

The correctness criteria did not however state what the answer produced by the rewriting meant. It was subsequently shown in [DG98, GM99] that the rewritings actually produced the *certain answer* $\cap(\varphi(poss(\mathcal{S})))$. This is easy to see in Example 9. Source S_1 contains facts of R_1 that join with a fact in R_2 , and S_2 contains facts of R_1 that join with a fact in R_3 . Since the sources provide us no information about R_2 and R_3 , there is a possible database where R_2 and R_3 are empty. Therefore, the facts of R_1 that are in every possible database is the union of the facts in S_1 and S_2 , which is exactly what the “correct” rewriting would produce in Example 9.

The first algorithms for constructing rewritings essentially explored the whole (finite) solution space in order to find all χ 's. Later more efficient algorithms were developed. These include the set-covering based algorithm in [GM99], its more detailed implementation MiniCon [PL00], the CoreCover algorithm [ALU01], and the algorithm [ALM02] taking arithmetic comparisons into account. Other extensions of the basic technique are surveyed in [Hal00, Hal01]. The cases where the queries defining the sources, as well as the user query, are allowed to be more general than conjunctive ones are analyzed in [AD98].

However, perhaps since the relationship between information integration and incomplete information had not been clearly articulated, there were no algorithms for computing the exact answer. To see the difference of the certain and exact answers in information integration applications, consider the following example.

Example 10 Suppose that the global schema contains two relations, $Affiliation(Actor, Studio)$, and $Star(Name, Type)$. A fact $Affiliation(Greta Garbo, MGM)$ means that Greta Garbo is affiliated with the MGM studio, and a fact such as $Star(Elisabeth Taylor, Tragedienne)$ means that Elisabeth Taylor is a Tragedienne. We have two sources, S_1 and S_2 , with definitions

$V_1(X, Y) \leftarrow Affiliation(X, Y), Star(X, Z)$, and $V_2(X, Z) \leftarrow Affiliation(X, Y), Star(X, Z)$. The user issues the query $\varphi = Q(X, Y, Z) \leftarrow Affiliation(X, Y), Star(X, Z)$.

Using the correctness criteria in [LRO96], we get an empty rewriting, and thus an empty query result for the user. However, suppose the extension $\mathbf{s}_1 = \{V_1(\text{Greta Garbo, MGM})\}$, and $\mathbf{s}_2 = \{V_2(\text{Elisabeth Taylor, Tragedienne})\}$. Then $T(\{S_1, S_2\}) = \{\text{Affiliation}(\text{Greta Garbo, MGM}), \text{Affiliation}(\text{Elisabeth Taylor, } X_1), \text{Star}(\text{Elisabeth Taylor, Tragedienne}), \text{Star}(\text{Greta Garbo, } X_2)\}$, and $\widehat{\varphi}(T(\{S_1, S_2\})) = \{Q(\text{Greta Garbo, MGM, } Y_1), Q(\text{Elisabeth Taylor, } Y_2, \text{Tragedienne})\}$. This exact answer tells the user that Greta Garbo is affiliated with MGM and is an actress of an unknown type, and that Elisabeth Taylor is a Tragedienne and is affiliated with an unknown studio. This is all the information we can deduce based on the incomplete information about the global relations provided by the sources. ■

We now proceed as follows. In the remainder of this section we first generalize the notion of query containment, which enables comparisons between different reformulations of queries, to *p-containment*. Next we define *p-rewritings* based on this more general containment. Then we extend the standard query evaluation so that given a p-rewriting, it computes the exact answer. In Section 5 we provide a unification-based method for actually producing the p-rewritings.

P-containment and p-containment mappings

We can broaden the classical notion of query containment as follows.

Let φ_1 and φ_2 be conjunctive queries. A query φ_1 is said to be *p-contained* in φ_2 , denoted $\varphi_1 \subseteq_p \varphi_2$, if and only if there exists a conjunctive query ϕ , where φ_1 is equivalent to $\pi_L(\phi)$ (π is relational projection), for some list L of columns in $\text{head}(\phi)$ taken in the original order, such that for all databases \mathbf{db} , $\phi(\mathbf{db}) \subseteq \varphi_2(\mathbf{db})$. Note that p-containment is a generalization of query containment since L can be the list of all columns in φ_2 .

In order to facilitate testing of p-containment, the classical notion of a containment mapping [CM77], can be generalized to define p-containment

mappings. A *p-containment mapping* from a conjunctive query φ_2 to a conjunctive query φ_1 is a mapping μ , from variables of φ_2 to variables and constants of φ_1 , such that

1. $\mu(\text{body}(\varphi_2)) \subseteq \text{body}(\varphi_1)$, and
2. for every variable X in $\text{head}(\varphi_1)$ there is a variable Y in $\text{head}(\varphi_2)$, such that $\mu(Y) = X$.

Example 11 Consider the following queries

$$\varphi_1 = Q_1(X) \leftarrow R_1(X, Y), R_2(Y, Y), R_3(Y, Z)$$

$$\varphi_2 = Q_2(X', Y') \leftarrow R_1(X', Y'), R_2(Y', Z')$$

There is a p-containment mapping $\mu = \{X'/X, Y'/Y, Z'/Y\}$ from φ_2 to φ_1 . ■

As consequence, we can now use p-containment mappings to test p-containment of conjunctive queries.

Theorem 4 *A query φ_1 is p-contained in a query φ_2 if and only if there is a p-containment mapping from φ_2 to φ_1 .*

Proof. Let μ be a p-containment mapping from φ_2 to φ_1 , and let \mathbf{db} be an arbitrary database. A fact t_1 in $\varphi_1(\mathbf{db})$ is generated by some valuation σ . Then $\sigma \circ \mu$ is a valuation that generates the corresponding fact t_2 in $\varphi_2(\mathbf{db})$. To see that this is indeed so, let $b \in \text{body}(\varphi_2)$. Then $\sigma \circ \mu(b) = \sigma(c) \in \mathbf{db}$, for some $c \in \text{body}(\varphi_1)$. Therefore, $\sigma \circ \mu(b) \in \mathbf{db}$. From requirement 2 of a p-containment mapping it follows that $\sigma \circ \mu(\text{head}(\varphi_2)) = \pi_L(\sigma \circ \mu(\text{head}(\varphi_1)))$, where L is a list of variables in $\text{head}(\varphi_2)$ in the original order. Thus, $\varphi_1 \subseteq_p \varphi_2$.

Let $\varphi_1 \subseteq_p \varphi_2$. Let \mathbf{db} be the canonical database that is the “frozen” $\text{body}(\varphi_1)$, that is, $\mathbf{db} = \rho(\text{body}(\varphi_1))$, where ρ is an injective valuation (the “freezing mapping”). By the definition of p-containment there exist a conjunctive query ϕ , such that $\phi(\mathbf{db}) \subseteq \varphi_2(\mathbf{db})$ and $\varphi_1 = \pi_L(\phi)$ for some ordered list L of columns in $\text{head}(\phi)$. Obviously, $\varphi_1(\mathbf{db})$ contains a fact t_1 , which is the “frozen” $\text{head}(\varphi_1)$. Since $\varphi_1 = \pi_L(\phi)$ there must be a fact t_2 in $\phi(\mathbf{db})$, such that $\pi_L(t_2) = t_1$. Since $\phi(\mathbf{db}) \subseteq \varphi_2(\mathbf{db})$, we have $t_2 \in \varphi_2(\mathbf{db})$.

Let σ be a valuation that generates the fact t_2 in $\varphi_2(\mathbf{db})$. Let ρ be the the “freezing” mapping, which also is a valuation that generates the fact t_1 in $\varphi_1(\mathbf{db})$. Then $\rho^{-1} \circ \sigma$ is a p-containment mapping from φ_2 to φ_1 .

To see that this is indeed so note two things. First, that each atom $b \in \text{body}(\varphi_2)$ is mapped by σ to some fact in \mathbf{db} , which is a frozen version of some atom $c \in \text{body}(\varphi_1)$, so $\rho^{-1} \circ \sigma$ maps b to the unfrozen fact, that is to c itself.

Second, note that those variables in $\text{head}(\varphi_2)$ that are also in $\text{head}(\varphi_1)$ are mapped by σ to constants in the fact t_1 , which is the frozen $\text{head}(\varphi_1)$, so that all of the head variables in φ_1 are covered. Thus $\rho^{-1} \circ \sigma$ maps corresponding variables in $\text{head}(\varphi_2)$ to the unfrozen $\text{head}(\varphi_1)$. Thus, $\rho^{-1} \circ \sigma$ is a p-containment mapping from φ_2 to φ_1 . ■

P-rewritings

We will now extend the classical notion of rewriting [LMSS95, Ull97] to *p-rewriting*.

Think of a set of source definitions $\{\psi_1, \dots, \psi_n\}$ for a class of source collections.

Let χ be a query over relation names V_i in $\{\text{head}(\psi_i)\}_i$. Then the *expansion* of χ , denoted χ^{exp} , is defined only if, for each $V_i(\cdot)$ in $\text{body}(\chi)$, there is a containment mapping μ from $\text{head}(\psi_i)$ to $V_i(\cdot)$. In this case χ^{exp} is obtained from χ by replacing all the $V_i(\cdot)$'s in $\text{body}(\chi)$ with $\mu(\text{body}(\psi_i))$. Existential variables in $\mu(\text{body}(\psi_i))$ are replaced by fresh variables when constructing χ^{exp} .

A query χ is a *contained rewriting* of φ if $\chi^{exp} \subseteq \varphi$. The query χ is a *p-contained rewriting* of φ if $\chi^{exp} \subseteq_p \varphi$.

Let χ be a p-contained rewriting of φ , and \mathcal{S} a source collection with definitions $\{\psi_1, \dots, \psi_n\}$ and extension $\mathbf{s} = \cup\{\mathbf{s}_1, \dots, \mathbf{s}_n\}$. We define the φ -evaluation of χ on \mathcal{S} , denoted $\chi_\varphi(\mathcal{S})$, as follows.

$$\chi_\varphi(\mathcal{S}) = \{\sigma_\mu(\text{head}(\varphi)) : \sigma(\text{body}(\chi)) \subseteq \mathbf{s}\},$$

where μ is a p-containment mapping from φ to χ^{exp} and σ is a valuation. Suppose a head variable X of φ is mapped by μ to a variable in the expansion of the atom $V_i(X_1, \dots, X_k)$, and $\mu(X)$ is the j^{th} existential variable in the definition ψ_i . Then the extension σ_μ of σ , is defined by setting

$$\sigma_\mu(X) = \begin{cases} \sigma(\mu(X)), & \text{if } \mu(X) \text{ occurs in } \text{head}(\chi) \\ f_{i,j}(\sigma(X_1), \dots, \sigma(X_k)), & \text{otherwise.} \end{cases}$$

In order to define $\tilde{\varphi}(\mathcal{S})$, the evaluation of conjunctive queries with function symbols in the head, we need an auxiliary function *replace* that when applied to a set of atoms with function terms, replaces each unique term with a unique variable. Now we set

$$\tilde{\varphi}(\mathcal{S}) = \text{replace} \left(\bigcup \{ \chi_\varphi(\mathbf{s}) : \chi^{exp} \subseteq_p \varphi \} \right),$$

and state the following important result:

Theorem 5 *For all source collections \mathcal{S} with definitions $\{\psi_1, \dots, \psi_n\}$,*

$$\tilde{\varphi}(\mathcal{S}) = \hat{\varphi}(T(\mathcal{S})),$$

up to renaming of the variables.

Proof. Let t be an arbitrary atom in $\tilde{\varphi}(\mathcal{S})$. Then there was a conjunctive query χ in the union of all p-rewritings of φ and a φ -evaluation of χ , using a valuation σ_μ such that $t = \sigma_\mu(\text{head}(\varphi))$, and $\sigma(\text{body}(\chi)) \subseteq \mathbf{s}$, where μ is a containment mapping from φ to χ^{exp} .

Let the extension of \mathcal{S} be \mathbf{s} . Since $\sigma(\text{body}(\chi)) \subseteq \mathbf{s}$, it means that all atoms in $\sigma(\text{body}(\chi^{exp}))$ are in $T(\mathcal{S})$ (with fresh existential variables). Since μ is a containment mapping from φ to χ^{exp} , we have that $\sigma(\mu(\text{body}(\varphi))) \subseteq T(\mathcal{S})$. Thus $\sigma(\mu(\text{head}(\varphi))) \in \hat{\varphi}(T(\mathcal{S}))$. Now $\sigma(\mu(\text{head}(\varphi)))$ is equal to $\sigma_\mu(\text{head}(\varphi))$, except for positions that don't occur in $\text{head}(\chi)$, these have been replaced by function terms in $\sigma_\mu(\text{head}(\varphi))$. Now let us consider the case that a given term appears somewhere else in $\tilde{\varphi}(\mathcal{S})$. There are two cases, either the repeating term appears in the same atom or it appears in a different atom.

Case 1: (See Example 12 further.) The repeating term appears in the same atom. This means that either there were repeating variables in $\text{head}(\varphi)$ and it is obvious that in $\hat{\varphi}(T(\mathcal{S}))$ there is two occurrences of the same variable. Another possibility is that in χ^{exp} there were two occurrences of the same existential variable. In this situation in $T(\mathcal{S})$ there would be two occurrences of the same variables and since $\hat{\varphi}$ always substitutes variables of the query for the variables of the tableau they would also appear in $\hat{\varphi}(T(\mathcal{S}))$.

Case 2: (See Example 13 further.) The repeating term appears in some atom t' in $\tilde{\varphi}(\mathcal{S})$. In this case there are two possibilities, either t' was produced by the φ -evaluation of χ or by the φ -evaluation of some other rewriting χ' . In either situation both terms originated from the same fact in \mathbf{s} and, moreover, from the same existential variable because function terms are subscripted with variable index. It is obvious that there would be two occurrences of the same variables in $T(\mathcal{S})$.

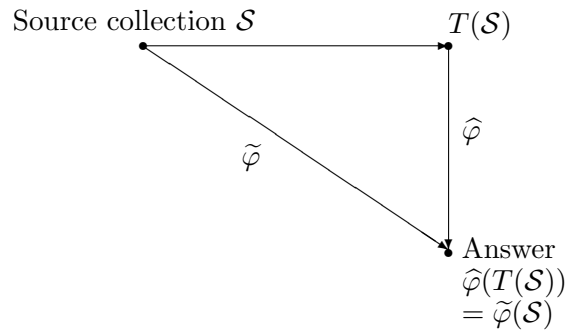
For the proof of inclusion in the other direction, let t be an arbitrary atom in $\hat{\varphi}(T(\mathcal{S}))$. Suppose $\text{body}(\varphi)$ consists of atoms b_1, b_2, \dots, b_n . Then there is a substitution θ , such that $\theta(b_i) \in T(\mathcal{S})$, for all $i \in \{1, \dots, n\}$. But each atom $\theta(b_i)$ is in $T(\mathcal{S})$, which follows from the fact that there is a source S_{i_j} , with definition ψ_{i_j} and extension \mathbf{s}_{i_j} , a valuation σ_{i_j} , and a fact $t_{i_j} \in \mathbf{s}_{i_j}$, such that $t_{i_j} = \sigma_{i_j}(\text{head}(\psi_{i_j}))$ and $\theta(b_i) \in \text{refresh}(\sigma_{i_j}(\text{body}(\psi_{i_j})))$. If we take a query ϕ such that $\text{body}(\phi) = \sigma_{1_j}(\text{head}(\psi_{1_j})), \sigma_{2_j}(\text{head}(\psi_{2_j})), \dots, \sigma_{n_j}(\text{head}(\psi_{n_j}))$, and $\text{head}(\phi) = (\sigma_{1_j} \cup \sigma_{2_j} \cup \dots \cup \sigma_{n_j})(\text{head}(\varphi))$, we have a p-containment mapping (namely θ), from φ to χ^{exp} . Consequently χ will be an element in the union of queries that $\tilde{\varphi}$ considers, and obviously χ generates the fact t when applied to \mathbf{s} . Now let us consider the case that a given variable appears somewhere else in $\hat{\varphi}(T(\mathcal{S}))$. Again there are two cases to consider, either the repeating variable appears in the same atom or it appears in

a different atom. In both cases, the same variables are in $\widehat{\varphi}(T(\mathcal{S}))$ because there were the same variables in $T(\mathcal{S})$ —remember that $\widehat{\varphi}$ always substitutes variables from the query for the variables from the tableau. The same variables could appear in the tableau only if they come from inversion of the same fact and, moreover, from the same existential variable in the query defining the view. And since this is the case it is obvious that in the corresponding atoms in the $\widehat{\varphi}$ -evaluation on \mathcal{S} there will be the same terms in the corresponding positions. ■

Example 12 For example, let the source collection be $\mathcal{S} = \{S_1, S_2\}$, where $\psi_1 = V_1(X_1, X_3) \leftarrow R_1(X_1, X_2), R_2(X_2, X_3)$, $\mathbf{s}_1 = \{V_1(a, c)\}$, $\psi_2 = V_2(X_1) \leftarrow R_3(X_1, X_2)$, and $\mathbf{s}_2 = \{V_2(c)\}$. Let the query φ be $Q(X, Y, Y, Z, W) \leftarrow R_1(X, Y), R_2(Y, Z), R_3(Z, W)$. In this case there is one (minimal) p-contained rewriting of φ , namely $\chi = Q(X, Y) \leftarrow V_1(X, Y), V_2(Y)$. Applying χ_φ gives us the atom $Q(a, f_{1,2}(a), f_{1,2}(a), c, f_{2,2}(c))$, and applying the *replace* function yields the answer $\{Q(a, Y, Y, c, W)\}$. Had we used the tableau method instead, we would have gotten $T(\mathcal{S}) = \{R_1(a, Y), R_2(Y, c), R_3(c, W)\}$. Then applying $\widehat{\varphi}$ to $T(\mathcal{S})$ would have given $\{Q(a, Y, Y, c, W)\}$. ■

Example 13 For an additional example, consider the following source collection $\mathcal{S} = \{S_1, S_2, S_3\}$, where $\psi_1 = V_1(X_1, X_3) \leftarrow R_1(X_1, X_2, X_3)$, $\mathbf{s}_1 = \{V_1(a, b)\}$, $\psi_2 = V_2(X_1, X_2) \leftarrow R_2(X_1, X_2)$, $\mathbf{s}_2 = \{V_2(c, a)\}$, $\psi_3 = V_3(X_2) \leftarrow R_2(X_1, X_2)$, and $\mathbf{s}_3 = \{V_3(a)\}$. Let φ be $Q(X, Y, Z, W) \leftarrow R_2(X, Y), R_1(Y, Z, W)$. The two (minimal) p-contained rewritings of φ are $\chi_1 = Q(X, Y, W) \leftarrow V_2(X, Y), V_1(Y, W)$ and $\chi_2 = Q(Y, W) \leftarrow V_3(Y), V_1(Y, W)$. Applying $(\chi_1)_\varphi$ and $(\chi_2)_\varphi$ gives us two atoms $Q(c, a, f_{1,2}(a, b), b)$ and $Q(f_{3,1}(a), a, f_{1,2}(a, b), b)$, and applying the *replace* function yields the answer $\{Q(c, a, Z, b), Q(X, a, Z, b)\}$. Had we constructed the tableau for this source collection first, we would have gotten $T(\mathcal{S}) = \{R_1(a, Z, b), R_2(c, a), R_2(X, a)\}$. Then applying $\widehat{\varphi}$ to $T(\mathcal{S})$ would have given $\{Q(c, a, Z, b), Q(X, a, Z, b)\}$. ■

Theorem 5 can be illustrated by the following diagram.



Note that since $\widetilde{\varphi}(\mathcal{S})$ computes a tableau that is equivalent to $\widehat{\varphi}(T(\mathcal{S}))$ the result of this computation can be used for subsequent querying.

5 Unification-based Rewriting.

The definition of a p-rewriting in the previous section is purely declarative. We now provide a constructive method for producing a p-rewriting given a query φ , and the description $\{\psi_1, \dots, \psi_n\}$ for a class of source collections.

A *unifier* for two atoms b_1 and b_2 is a substitution θ such that $\theta(b_1) = \theta(b_2)$. A substitution θ is *more general* than a substitution ζ if for some substitution ζ' , $\zeta = \theta \circ \zeta'$. A *most general unifier* for a and b is a unifier θ such that, for each unifier ζ of a and b , θ is more general than ζ .

Let \mathcal{A} be a subset of the atoms in $\{body(\psi_i)\}_i$, such that \mathcal{A} unifies with the set of atoms in $body(\varphi)$.

Let the *mgu* that achieves this unification be θ . If θ does not equate any view variable X to any other view variable or constant, unless X appears in the head of its view then do the following. Choose a subset $\{\psi_{i_1}, \dots, \psi_{i_k}\}$ of the view definitions, such that every atom in \mathcal{A} occurs in some $body(\psi_{i_j})$ and every $body(\psi_{i_j})$ contains at least one atom in \mathcal{A} . Now, define χ as

$$\pi_L(\theta(head(\varphi)) \leftarrow \theta(head(\psi_{i_1})), \dots, \theta(head(\psi_{i_k}))),$$

where L is a list of all variables in $\theta(head(\varphi))$ that also occur in

$$\theta(head(\psi_{i_1})), \dots, \theta(head(\psi_{i_k})).$$

The output, $rew(\varphi)$, of the algorithm is the union of all possible χ 's thus constructed. Note that there is a finite number (modulo renaming) of such χ 's since each one is based on a subset \mathcal{A} of the atoms in the bodies of the view queries, and there is a unique (up to renaming) *mgu* for each pair of atom sets. Furthermore, there is a finite number of ways of choosing the “covering” subset \mathcal{A} of the view atoms.

Theorem 6 *rew(φ) is equivalent to the union of all p-contained rewritings of φ .*

Proof. To prove that $rew(\varphi)$ contains *only* semantically correct rewritings we need to show that for each $\chi \in rew(\varphi)$, there is a p-containment mapping from φ to χ^{exp} . Let χ be an arbitrary element in $rew(\varphi)$ and let θ be the *mgu* that was used to produce χ . Consider the expansion of χ

$$\chi^{exp} = \pi_L(head(\theta(\varphi)) \leftarrow \theta(body(\psi_{i_1})), \dots, \theta(body(\psi_{i_k}))).$$

It is obvious that θ also gives a p-containment mapping from φ to χ^{exp} . Thus χ is p-contained in φ .

To prove that the method computes *all* semantically correct rewritings we will show that for each p-rewriting χ of a given conjunctive query φ , there will be a query χ' in $rew(\varphi)$ such that $\chi \subseteq \chi'$.

Suppose χ is of the form

$$head(\chi) \leftarrow V_{i_1}(\cdot), \dots, V_{i_k}(\cdot),$$

where each $V_{i_j}(\cdot)$ comes from $head(\psi_{i_j})$. Without loss of generality we assume that no variable is unnecessarily projected out from $head(\chi)$.

Since χ is a semantically correct rewriting, we know that there is a p-containment mapping μ from φ to χ^{exp} . It now follows that there is a subset \mathcal{A} of the atoms in the bodies of $\{\psi_{i_1}, \dots, \psi_{i_k}\}$, such that μ is an unifier for \mathcal{A} and $body(\varphi)$. Suppose that μ is actually the most general unifier. Then the query χ' obtained as $\pi_L(head(\varphi)) \leftarrow \theta(head(\psi_{i_1})), \dots, \theta(head(\psi_{i_k}))$, where L is a list of variables $X \in head(\theta(\varphi))$, such that $X \in \theta(head(\psi_{i_1})), \dots, \theta(head(\psi_{i_k}))$, will be in $rew(\varphi)$. Obviously the identity function will be a containment mapping from χ' to χ .

Suppose then that θ is not the most general unifier for \mathcal{A} and $body(\varphi)$. Then there will in $rew(\varphi)$ be a query χ' , such that the body of χ'^{exp} is a more specific instance of the body of the expansion of χ' . This guarantees that there will be a containment mapping from χ' to χ .

The claim of the theorem now follows from the characterization of equivalence of unions of conjunctive queries by Sagiv and Yannakakis [SY80]. ■

Example 14 In Example 10 the unification method would produce as rewriting the union of $Q(X, Y) \leftarrow V_1(X, Y)$ and $Q(X, Z) \leftarrow V_2(X, Z)$. When using the $\tilde{\varphi}$ -evaluation on this rewriting to the sources in Example 10 would produce the result $\{Q(\textit{Greta Garbo}, \textit{MGM}, Y_1), Q(\textit{Elisabeth Taylor}, Y_2, \textit{Tragedienne})\}$. ■

Obviously it would be desirable to compute a rewriting that encodes the containment mappings μ needed in the $\tilde{\varphi}$ evaluation in the rewriting directly, rather than recompute the μ 's at evaluation time.

To do that we can extend our unification based method as follows. Everything up to producing a rewriting χ remains the same but for every rewriting χ , where

$$body(\chi) = \theta(head(\psi_{i_1})), \dots, \theta(head(\psi_{i_k})),$$

we set $head(\chi) = \eta(\theta(head(\varphi)))$, where η is defined as follows. Suppose a variable X of φ is unified by θ with a variable in an atom in $body(\psi_{i_1})$, $\theta(X)$ is the j^{th} distinguished variable in the $body(\psi_{i_1})$, and existential variables in ψ_{i_1} are X_1, \dots, X_k . Then the

$$\eta(X) = \begin{cases} \theta(X), & \text{if } \theta(X) \text{ occurs in } head(\psi_{i_j}) \\ f_{i,j}(\theta(X_1), \dots, \theta(X_k)), & \text{otherwise.} \end{cases}$$

The union of all rewritings produced by our extended method can then be evaluated on the source collection in the usual manner and *replace* function can be applied to the resulting set.

The following lemma directly follows from Theorem 6 and the construction of skolemized rewriting by our extended unification based method.

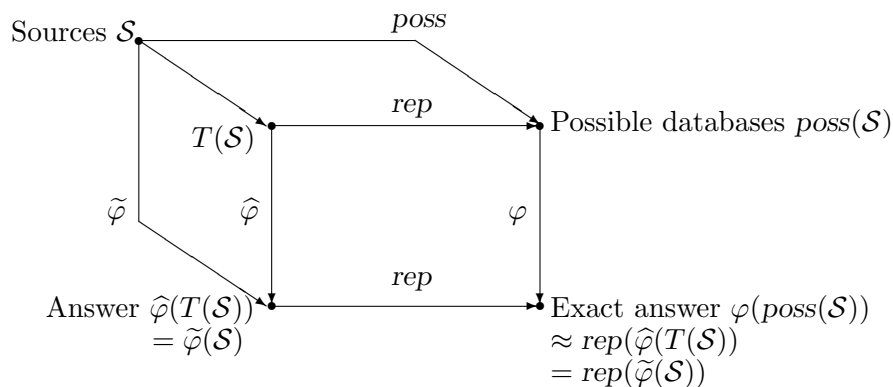
Lemma 1 *The union of all rewritings produced by the extended unification based method relative to a query φ is equivalent to $rew(\varphi)$.*

Proof. It is obvious that the union of all rewritings produced by our extended unification based method contains one skolemized rewriting for each p-contained rewriting produced by the original method. The claim follows from the Theorem 6 and the fact that function terms are inserted by our extended method in the same way as by $\tilde{\varphi}$. ■

Example 15 In Example 14 the extended unification method would produce as rewriting the union of $Q(X, Y, f_{1,1}(X, Y)) \leftarrow V_1(X, Y)$ and $Q(X, f_{2,1}(X, Z), Z) \leftarrow V_2(X, Z)$. Evaluation of the rewriting on the sources in Example 10 will indeed produce $\{Q(\textit{Greta Garbo}, \textit{MGM}, Y_1), Q(\textit{Elisabeth Taylor}, Y_2, \textit{Tragedienne})\}$. ■

6 Conclusion

Query answering in Information Integration systems can now be summarized by combining the previous commutative diagrams into the following.



We note that the computational complexity of query evaluation does not change when moving from certain to exact answers. This means that the basic complexity results, such as those in [AKG91, AD98] carry over to our more general theory. This does not mean that the theory does not open up any new complexity questions, they will be the topic of future papers.

References

- [Abi99] Abiteboul S. On Views and XML. In *Proc. 18th Annual ACM Symp. Principles of Databases (PODS '98)*, Philadelphia, Pennsylvania 1999, pp. 1–9.
- [AD98] Abiteboul S. and O. M. Duschka. Complexity of Answering Queries Using Materialized Views. In *Proc. 17th Annual ACM Symp. Principles of Databases (PODS '98)*, Seattle, Washington 1998, pp. 254–263.
- [AHV95] Abiteboul S., R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, Reading Ma. 1995.
- [AKG91] Abiteboul S., P.C. Kanellakis, and G. Grahne. On the Representation and Querying of Sets of Possible Worlds. In *J. TCS* **78**:1, 1991, pp. 158–187.
- [ALM02] Afrati F., C. Li and P. Mitra. Answering Queries Using Views with Arithmetic Comparisons. In *Proc. 21st ACM Symp. on Principles of Database Systems (PODS '02)*, Madison, Wisconsin 2002, pp. 209–220.
- [ALU01] Afrati F., C. Li and J. D. Ullman. Generating Efficient Plans for Queries Using Views. In *Proc. of the ACM SIGMOD International Conference on Management of Data (SIGMOD '01)*, Dallas, Texas 2001.

- [CM77] A. K. Chandra, P. M. Merlin. Optimal implementation of conjunctive queries. In *Proc. ACM SIGACT Symp. on the Theory of Computing (STOC '77)*, pp. 77–90, 1977.
- [DG98] Dushka O. M. and M. R. Genesereth. Query Planning with Disjunctive Sources. In *Proceedings of the AAAI Workshop on AI and Information Integration*, Madison, Wisconsin 1998.
- [Gra84] Grahne G. Dependency satisfaction in databases with incomplete information. *Proc. of the 10th International Conference on Very Large Databases (VLDB '84)*, Singapore 1984, pp. 37–45.
- [Gra91] Grahne G. *The Problem of Incomplete Information in Relational Databases*. Lecture Notes in Computer Science, vol. 554. Springer-Verlag, Berlin 1991.
- [GM99] Grahne G. and A. O. Mendelzon. Tableau Techniques for Querying Information Sources through Global Schemas. In *Proc. 7th International Conference on Database Theory (ICDT '99)*, Jerusalem, Israel 1999, pp. 332–347.
- [GK02] Grahne G. and V. Kiricenko. Obtaining more answers from Information Integration Systems. *Proc. Fifth International Workshop on the Web and Databases (WebDB '02)*, Madison, Wisconsin 2002, pp. 67–76.
- [GK03] Grahne G. and V. Kiricenko. Partial answers in Information Integration Systems. *Proc. 5th ACM CIKM International Workshop on Web Information and Data Management (WIDM '03)*, New Orleans, Louisiana 2003, pp. 98–101.
- [Hal00] Halevy A. Y. Theory of Answering Queries Using Views. *SIGMOD Record* 29(4): 40-47 (2000).
- [Hal01] Halevy A. Y. Answering queries using views: A survey. *VLDB Journal* 10(4): 270-294 (2001).
- [IL84] Imielinski T. and W. Lipski Jr. Incomplete Information in Relational Databases. *J. ACM* 31:4, 1984, pp. 761–791.
- [Len02] Lenzerini M. Data Integration: A Theoretical Perspective. Invited tutorial in *Proc. 21st ACM Symp. on Principles of Database Systems (PODS '02)*, Madison, Wisconsin 2002, pp. 233–246.
- [LMSS95] Levy A. Y., A. O. Mendelzon, Y. Sagiv and D. Srivastava. Answering Queries Using Views. In *Proc. 14th ACM Symp. on Principles of Database Systems (PODS '95)*, San Jose, California 1995, pp. 95–104.
- [LRO96] Levy A. Y, A. Rajaraman, J. J. Ordille. Querying Heterogeneous Information Sources Using Source Descriptions. *Proc. 22nd Int'l. Conf. on Very Large Databases (VLDB '96)*, Mumbai (Bombay), India 1996, pp. 251–262.
- [Lip84] W. Lipski Jr. On Relational Algebra with Marked Nulls. In *Proc. 3rd ACM Symp. on Principles of Database Systems (PODS '84)*, Waterloo, Ontario 1984, pp. 201–203.

- [Men84] Mendelzon A. O. Database States and Their Tableaux. In *ACM Trans. on Databases Systems* **9**:2, 1984, pp. 264–282.
- [MM01] Mendelzon A. O. and G. Mihaila. Querying Partially Sound and Complete Data Sources. In *Proc. 20th ACM Symp. on Principles of Database Systems (PODS '01)*, Santa Barbara, California 2001, pp. 162–170
- [PL00] Pottinger R. and A. Y. Levy. A Scalable Algorithm for Answering Queries Using Views. In *Proc. 26th Int'l. Conf. on Very Large Databases (VLDB '00)*, Cairo, Egypt 2000, pp.484–495.
- [Rei78] Reiter R. On Closed World Databases. In *Logic and Databases*, H. Gallaire and J. Minker (Eds.), Plenum Press, New York, 1978, pp. 56–76.
- [SY80] Y. Sagiv, M. Yannakakis. Equivalence among relational expressions with the union and difference operators. In *J. ACM* **27**:4, 1980, pp. 633–655.
- [Ull89] Ullman J. D. *Principles of Database and Knowledge-Base Systems, Volume II*. Computer Science Press, Rockville, Maryland 1989.
- [Ull97] Ullman J. D. Information Integration Using Logical Views. In *Proc. 6th International Conference on Database Theory (ICDT '97)*. Delphi, Greece 1997, pp. 19–40.
- [Var86] Vardi M. Y. Querying Logical Databases. *Journal of Computer and System Sciences* **33**, 1986, pp. 142–160.
- [YL87] Yang H. Z. and P. A. Larson. Query Transformation for PSJ Queries. In *Proc. 13th Int'l. Conf. on Very Large Data Bases (VLDB '87)* Brighton, England 1987, pp. 245–254.