



## Optimizing Recursive Information Gathering Plans in EMERAC\*

SUBBARAO KAMBHAMPATI<sup>†</sup>

rao@asu.edu

ERIC LAMBRECHT

ULLAS NAMBIAR

ZAIQING NIE

GNANAPRAKASAM SENTHIL

*Department of Computer Science and Engineering, Arizona State University, Tempe, AZ 85287, USA*

*Received July 24, 2002; Revised April 3, 2003; Accepted April 4, 2003*

**Abstract.** In this paper we describe two optimization techniques that are specially tailored for information gathering. The first is a greedy minimization algorithm that minimizes an information gathering plan by removing redundant and overlapping information sources without loss of completeness. We then discuss a set of heuristics that guide the greedy minimization algorithm so as to remove costlier information sources first. In contrast to previous work, our approach can handle recursive query plans that arise commonly in the presence of constrained sources. Second, we present a method for ordering the access to sources to reduce the execution cost. This problem differs significantly from the traditional database query optimization problem as sources on the Internet have a variety of access limitations and the execution cost in information gathering is affected both by network traffic and by the connection setup costs. Furthermore, because of the autonomous and decentralized nature of the Web, very little cost statistics about the sources may be available. In this paper, we propose a heuristic algorithm for ordering source calls that takes these constraints into account. Specifically, our algorithm takes both access costs and traffic costs into account, and is able to operate with very coarse statistics about sources (i.e., without depending on full source statistics). Finally, we will discuss implementation and empirical evaluation of these methods in *Emerac*, our prototype information gathering system.

**Keywords:** data integration, information gathering, Web and databases, query optimization

### 1. Introduction

The explosive growth and popularity of the world-wide web have resulted in thousands of structured queryable information sources on the Internet, and the promise of unprecedented information-gathering capabilities to lay users. Unfortunately, the promise has not yet been transformed into reality. While there are sources relevant to virtually any user-queries, the morass of sources presents a formidable hurdle to effectively accessing the information. One

\*This research is supported in part by NSF young investigator award (NYI) IRI-9457634, Army AASERT grant DAAH04-96-1-0247, and NSF grant IRI-9801676. We thank Selçuk Candan for many helpful comments. Preliminary versions of parts of this work have been presented at IJCAI (Lambrech et al., 1999), and workshops on Intelligent Information Integration (Kambhampati and Gnanaprakasam, 1999; Lambrecht and Kambhampati, 1998).

<sup>†</sup>Author to whom all correspondence should be addressed.

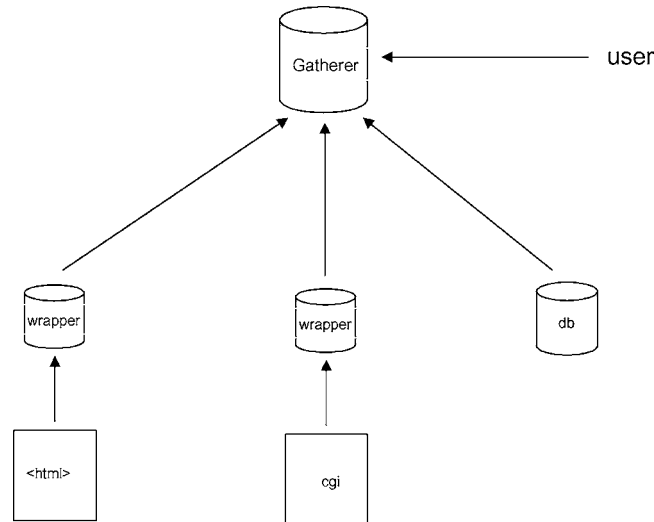


Figure 1. The information gatherer acts as an intermediary between the user and information sources on the Internet.

way of alleviating this problem is to develop *information gatherers* (also called mediators) which take the user's query, and develop and execute an effective *information gathering plan*, that accesses the relevant sources to answer the user's query efficiently.<sup>1</sup> Figure 1 illustrates the typical architecture of such a system for integrating diverse information sources on the internet. Several first steps have recently been taken towards the development of a theory of such gatherers in both database and artificial intelligence communities.

The information gathering problem is typically modeled by building a virtual global schema for the information that the user is interested in, and describing the accessible information sources as materialized views on the global schema.<sup>2</sup> The user query is posed in terms of the relations of the global schema. Since the global schema is virtual (in that its extensions are not stored explicitly anywhere), computing the query requires rewriting (or "folding" (Qian, 1996)) the query such that all the extensional (EDB) predicates in the rewrite correspond to the materialized view predicates that represent information sources. Several researchers (Levy et al., 1996; Qian, 1996; Kwok and Weld, 1996) have addressed this rewriting problem. Recent research by Duschka and his co-workers (Duschka and Genesereth, 1997; Duschka and Levy, 1997) subsumes most of this work, and provides a clean methodology for constructing information gathering plans for user queries posed in terms of a global schema. The plans produced by this methodology are "maximally contained" in that any other plan for answering the given query is contained in them.<sup>3</sup>

Generating source complete plans however is only a first step towards efficient information gathering. A crucial next step, which we focus on in this paper, is that of query plan optimization. Maximally contained plans produced by Duschka's methodology are conservative in that they in essence wind up calling any information source that may be remotely relevant to the query. Given the autonomous and decentralized nature of the

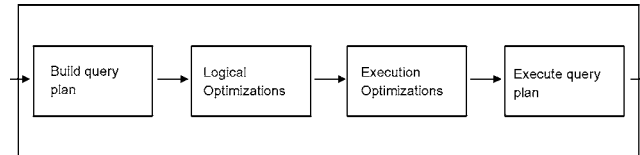


Figure 2. The full process of query planning.

Internet, sources tend to have significantly overlapping contents (e.g. mirror sources), as well as varying access costs (premium vs. non-premium sources, high-traffic vs. low-traffic sources). Naive execution of maximally contained plans will access all potentially relevant sources and be prohibitively costly, in terms of the network traffic, response time, and access costs (in the case of “premium” sources that charge for access).

At first blush, it would seem that we should be able to directly apply the rich body of work on query optimization in databases (Chaudhuri, 1998) to solve this problem. Unfortunately, this does not work because many of the assumptions made in the traditional database query optimization do not hold in information gathering scenarios. To begin with, in traditional databases, redundancy and overlap among different sources is not a major issue, while it is a very crucial issue in information gathering. Similarly, traditional query optimization methods depend on elaborate statistical models (histograms, selectivity indices etc.) of the underlying databases. Such statistical models may not be easily available for sources on the Internet.<sup>4</sup> Finally, even the work on optimizing queries in the presence of materialized views (c.f. Chaudhuri and Krishnamurthy, 1995) is not directly applicable as in such work materialized views are assumed to be available *in addition* to the main database. In contrast, the global database in information gathering is “virtual” and the only accessible information resides in materialized views whose statistical models are not easily available. For all these reasons, it is now generally recognized (c.f. Florescu et al., 1998) that query optimization for information gathering is a very important open problem.

In this paper we describe the query optimization techniques that we have developed in the context of *Emerac*, a prototype information gathering system under development. Figure 2 provides a schematic illustration of the query planning and optimization process in *Emerac*. We start by generating a query plan using the source inversion techniques described by Duschka and Genesereth (1997) and Duschka and Levy (1997). This polynomial time process gives us a “maximally contained” query plan which serves as the input for the optimization methods. As in traditional databases, our optimization phase involves two steps: logical optimization and execution optimization. In traditional databases, logical optimization involves rewriting a query plan, using relational algebra equivalences, to make it more efficient; while execution optimization involves steps such as ordering the access to the base relations to make computations of joins cheaper. For *Emerac*, the logical optimization step involves minimizing the maximally contained query plan such that access to redundant sources is removed. Execution optimization involves ordering the access to the information sources in the minimized plan so as to reduce execution cost.

*Our contributions.* For logical optimization, we present a technique that operates on the recursive plans generated by Duschka’s algorithm and greedily minimizes it so as to remove

access to costly and redundant information sources, without affecting the completeness of the plan. For this purpose, we use the so-called localized closed world (LCW) statements that characterize the completeness of the contents of a source relative to either the global (virtual) database schema or the other sources. Our techniques are based on an adaptation of Sagiv's (Sagiv, 1988) method for minimizing datalog programs under uniform equivalence. Although there exists some previous research on minimizing information gathering plans using LCW statements (Duschka, 1997; Friedman and Weld, 1997), none of it is applicable to minimization of information gathering plans containing recursion. Our ability to handle recursion is significant because recursion appears in virtually all information gathering plans either due to functional dependencies, binding constraints on information sources, or recursive user queries (Duschka and Genesereth, 1997). Additionally, in contrast to existing methods, which do pairwise redundancy checks on source accesses, our approach is capable of exploiting cases where access to one information source is rendered redundant by access to a combination of sources together. Large performance improvements in our prototype information gatherer, *Emerac*, attest to the cost-effectiveness of our minimization approach.

Ultimately plan execution in our context boils down to doing joins between the sources efficiently. When gathering information on the Internet, we typically cannot instruct two sources to join with each other. It is thus necessary to order the access to the sources. The existing methods for subgoal ordering assume that the plan is operating on a single "fully relational" (i.e., no binding restrictions) database, and that the plan execution cost is dominated by the number of tuples transferred. In contrast, sources on the Internet have a variety of access limitations and the execution cost in information gathering is affected significantly by the connection setup costs. We describe a way of representing the access capabilities of sources, and provide a greedy algorithm for ordering source calls that respects source limitations, and takes both access costs and traffic costs into account, without requiring full source statistics.

Although there exist other research efforts that address source redundancy elimination and optimization in the presence of sources with limited capabilities, *Emerac* is the first to consider end-to-end issues of redundancy elimination and optimization in recursive information gathering plans. It is also the first system to consider the source access costs as well as traffic costs together in doing optimization.

*Organization.* The rest of the paper is organized as follows. Section 2 provides the background on information integration and describes the motivation for our work. Section 3, we review the work on integrating diverse information sources by modeling them as materialized views on a virtual database. We pay special attention to the work of Duschka and Genesereth (1997) and Duschka and Levy (1997), which forms the basis for our own work. Section 4 briefly reviews the use of LCW statements and Sagiv's algorithm for datalog program minimization under uniform equivalence. Section 5 presents our greedy minimization algorithm that adapts Sagiv's algorithm to check for source redundancy in the context of the given LCW statements. We also explain how the inter-source subsumption relations can be exploited in addition to LCW statements (Section 5.1). We then discuss the complexity of the minimization and present heuristics for biasing the greedy minimization strategy. Section 6.1 describes our algorithm for ordering source accesses during execution. Section 7 describes the architecture of *Emerac*, our prototype information gatherer, and presents an empirical

evaluation of the effectiveness of our optimization techniques. Section 8 discusses related work. Section 9 presents our conclusions, and outlines our current research directions.

## 2. Background and motivation

### 2.1. Background

In order to place the work on *Emerac* in the proper context, it is important to provide a broad classification of the prior work on information integration. Information integration (aka data integration) has received a significant amount of attention in the recent years, and several systems have been developed. These include InfoMaster (Geddis et al., 1995), Information Manifold (Levy et al., 1996), Garlic (Haas et al., 1997), TSIMMIS (Garcia-Molina et al., 1997), HERMES (Adalı et al., 1996), and DISCO (Raschid et al., 1998). The similarities and differences among these systems can be broadly understood in terms of (1) the approach used to relate the mediator and source schemas and (2) the type of application scenario considered by the system.

The application scenarios considered to date can be usefully classified into two categories:

*Authorized integration of databases.* Integrating a set of heterogeneous database systems owned by a given corporation/enterprise.

*Information Gathering.* Integrating a set of information sources that export information related to some specific application area (e.g. comparison shopping of books, integration of multiple bibliography sources etc.).

In the first case, we would expect that the set of data sources are relatively stable, and that the mediation is “authorized” (in that the data sources are aware that they are being integrated). In the second, “information gathering” scenario, the set of data sources may be changing, and more often than not, the mediation may not have been explicitly authorized by the sources. Systems such as Garlic (Haas et al., 1997), TSIMMIS (Garcia-Molina et al., 1997), HERMES (Adalı et al., 1996), and DISCO (Raschid et al., 1998) can be characterized as aiming at authorized database integration, while InfoMaster (Geddis et al., 1995), Information Manifold (Levy et al., 1996), Occam (Kwok and Weld, 1996), Razor (Friedman and Weld, 1997) as well as the *Emerac* system presented in this paper address the information gathering scenario.

Although related in many ways, authorized database integration and information gathering systems do differ in two important ways—the way source and mediator schemas are modeled, and the types of approaches used for query optimization.

There are two broad approaches for modeling source and mediator schemas (Halevy, 2001). The “global as view” (GAV) approach involves modeling the mediator schema as a view on the (union of) source schemas. The “local as view” (LAV) approach involves modeling the source schemas as views on the mediated schema. The GAV approaches make query planning relatively easy as a query posed on the mediated schema can directly be rewritten in terms of sources. The LAV approach, in contrast, would require a more complex rewriting phase to convert the query posed on the mediator schema to a query on the sources.

The advantage of LAV approach however is that adding new sources to a mediator involves just modeling them as views on the mediated schema. In contrast, in the GAV approach, the mediated schema has to be rewritten every time a new source is added. Systems that address authorized integration of known databases, such as Garlic, Disco, TSIMMIS and HERMES use the GAV approach as they can be certain of a relatively stable set of sources. In contrast, systems aimed at information gathering, such as the Information Manifold (Levy et al., 1996) and InfoMaster (Geddis et al., 1995) use the LAV approach.

The other main difference among information integration systems is the types of approaches used for query optimization. Systems addressing authorized integration of known databases can count on the availability of statistics on the databases (sources) being integrated. Thus Garlic (Haas et al., 1997), TSIMMIS (Chawathe et al., 1994), HERMES (Adalı et al., 1996) and DISCO (Tomasic et al., 1997) systems attempt to use cost-based optimization algorithms for query planning. Systems addressing information gathering scenarios, on the other hand, cannot count on having access to statistics about the information sources. Thus, either the mediator has to learn the statistics it needs, or will have to resort to optimization algorithms that are not dependent on complete statistics about sources. Both Infomaster and Information Manifold system use heuristic techniques for query optimization.

Within the above classification, *Emerac* is aimed at information gathering. Thus, it is most closely related to systems such as InfoMaster (Geddis et al., 1995), Information Manifold (Levy et al., 1996), Occam (Kwok and Weld, 1996) and Razor (Friedman and Weld, 1997). Like these other systems, *Emerac* too uses the LAV approach to model source and mediator schemas, and uses heuristic techniques for query optimization. The specific contributions of *Emerac* over the other systems are:

- A systematic approach for handling minimization of (recursive) datalog query plans using the LCW information.
- A heuristic optimization technique for query plans that takes into account both the access and transfer costs.

### 3. Building query plans: Background

Suppose our global schema contains the world relation  $advisor(S, A)$ , where  $A$  is the advisor of  $S$ . Furthermore, suppose we have an information source ADDB, such that for every tuple  $(S, A)$  returned by it,  $A$  is the advisor of  $S$ . This can be represented as a materialized view on the global schema as follows:

$$ADDB(S, A) \rightarrow advisor(S, A)$$

We make the “open world assumption,” (OWA) on the sources (Abiteboul and Duschka, 1998), meaning that the ADDB source has some but not necessarily all of the tuples satisfying the *advisor* relation.

Suppose we want to retrieve all the students advised by Weld. We can represent our goal by the query  $Q$ :

$$query(S, A) :- advisor(S, A) \wedge A = \text{“Weld”}$$

Duschka and Genesereth (1997) and Duschka and Levy (1997) show how we can generate an information gathering plan that is “maximally contained” in that it returns every query-satisfying tuple that is stored in any of the accessible information sources. This method works by *inverting* all source (materialized view) definitions, and adding them to the query. The inverse,  $v^{-1}$ , of the materialized view definition with head  $v(X_1, \dots, X_m)$  is a set of logic rules in which the body of each new rule is the head of the original view, and the head of each new rule is a relation from the body of the original view. When we invert our definition above, we get:

$$\text{advisor}(S, A) \text{ :- ADDB}(S, A)$$

When this rule is added to the original query  $Q$ , we effectively create a datalog<sup>5</sup> program whose execution produces all the tuples satisfying the query.

Note that we are modeling sources as “conjunctive views” on the mediated schema. The complexity of finding the maximally contained plan depends on the expressiveness of the language used to describe sources. Abiteboul and Duschka (1998) show that as long as sources are described as conjunctive views on the mediated schema, and we use the open world assumption on the sources, maximally contained plans can be found in polynomial time. The complexity becomes NP-hard when the sources are written as disjunctive, and undecidable when the sources are written as recursive views on the mediated schema.

*Constrained sources & Recursion.* The materialized view inversion algorithm can be modified in order to model databases that have binding pattern requirements. Suppose we have a second information source, CONDB that requires the student argument to be bound, and returns the advisor of that given student. We denote this in its view as follows:

$$\text{CONDB}(\$S, A) \rightarrow \text{advisor}(S, A)$$

The ‘\$’ notation denotes that  $S$  must be bound for any query sent to CONDB. A straightforward inversion of this source will get us a rule of the form:

$$\text{advisor}(S, A) \text{ :- CONDB}(\$S, A)$$

which is unexecutable as  $S$  is not bound. This is handled by making up a new relation called *dom* whose extension is made to correspond to all possible constants that can be substituted for  $S$ . In our example, assuming that we have both the ADDB source and the CONDB source, the complete plan for the query, which we shall refer to as  $\mathcal{P}$ , is:

$$\begin{aligned} r1: & \text{query}(S, A) \text{ :- advisor}(S, A) \wedge A = \text{“Weld”} \\ r2: & \text{advisor}(S, A) \text{ :- ADDB}(S, A) \\ r3: & \text{advisor}(S, A) \text{ :- dom}(S) \wedge \text{CONDB}(S, A) \\ r4: & \text{dom}(S) \text{ :- ADDB}(S, A) \\ r5: & \text{dom}(A) \text{ :- ADDB}(S, A) \\ r6: & \text{dom}(A) \text{ :- dom}(S) \wedge \text{CONDB}(S, A) \end{aligned}$$

Notice that all extensional (EDB) predicates in the program correspond to source predicates (materialized views). Notice also the presence of  $dom(S)$  relation in the rule  $r3$ . Rules  $r4$ ,  $r5$  and  $r6$  define the extension of  $dom$  by collecting all possible constants that can be derived from source calls. Finally, note that rule  $r6$  is recursive, which makes the overall plan recursive, *even though* the original query as well as the source description are non-recursive. Given the ubiquitousness of constrained sources on the Internet, it is thus important that we know how to handle recursive information gathering plans.

It is worth noting that the complexity of finding maximally contained plans remains polynomial when we have sources with access constraints. The only change is that the query plan itself is going to be a recursive datalog program. This change can in turn significantly increase the execution cost of the plans. Consequently, we focus on using any information about the source overlap to minimize the query plan and remove the recursion as much as possible.

#### 4. Plan minimization preliminaries

The plan  $\mathcal{P}$  above accesses two different advisor databases to answer the query. It would be useful to try and cut down redundant accesses, as this would improve the execution cost of the plan. To do this however, we need more information about the sources. While the materialized view characterizations of sources explicate the world relations that are respected by each tuple returned by the source, there is no guarantee that all tuples satisfying those properties are going to be returned by that source.

One way to support minimization is to augment the source descriptions with statements about their relative coverage, using the so-called localized closed world (LCW) statements (Etzioni et al., 1997). An LCW statement attempts to characterize what information (tuples) the source is *guaranteed* to contain in terms of the global schema. Suppose, we happen to know that the source ADDB is guaranteed to contain all the students advised by Weld and Hanks. We can represent this information by the statement (note the direction of the arrow):

$$\begin{aligned} \text{ADDB}(S, A) &\leftarrow \text{advisor}(S, A) \wedge A = \text{“Weld”} \\ \text{ADDB}(S, A) &\leftarrow \text{advisor}(S, A) \wedge A = \text{“Hanks”} \end{aligned}$$

*Pair-wise rule subsumption.* Given the LCW statement above, intuitively it is obvious that we can get all the tuples satisfying the query  $\mathcal{Q}$  by accessing just ADDB. We now need to provide an automated way of making these determinations. Suppose we have two datalog rules, each of which has one or more materialized view predicates in its body that also have LCW statements, and we wish to determine if one rule subsumes the other. The obvious way of checking the subsumption is to replace the source predicates from the first rule with the bodies of their view description statements, and the source predicates from the second rule with the bodies of the LCW statements corresponding to those predicates. We now have the transformed first rule providing a “liberal” bound on the tuples returned by that rule, while the transformed second rule gives a “conservative” bound. If the conservative bound subsumes the liberal bound, i.e., if the transformed second rule “contains” (entails) the transformed first rule, we know that second rule subsumes the first rule. Duschka (1997)



shows that this check, while sufficient, is not a necessary condition for subsumption. He proposes a modified version that involves replacing each source predicate  $s$  with  $s \wedge v$  in the first rule, and with  $s \vee l$  in the second rule, where  $v$  is the view description of  $s$ , and  $l$  is the conjunction of LCW statements of  $s$ . If after this transformation, the second rule contains the first, then the first rule is subsumed by it.<sup>6</sup>

*Minimization under uniform equivalence.* Pair-wise rule subsumption checks alone are enough to detect redundancy in non-recursive plans (Levy, 1996; Friedman and Weld, 1997), but are inadequate for minimizing recursive plans. Specifically, recursive plans correspond to infinite union of conjunctive queries and checking if a particular rule of the recursive plan is redundant will involve trying to see if that part is subsumed by any of these infinite conjuncts (Ullman, 1989, pp. 908). We instead base our minimization process on the notion of “uniform containment” for datalog programs, presented in Sagiv (1998). To minimize a datalog program, we might try removing one rule at a time, and checking if the new program is equivalent to the original program. Two datalog programs are equivalent if they produce the same result for all possible assignments of EDB predicates (Sagiv, 1988). Checking equivalence is known to be undecidable. Two datalog programs are uniformly equivalent if they produce the same result for all possible assignments of EDB *and* IDB predicates. Uniform equivalence is decidable, and implies equivalence. Sagiv (1998) offers a method for minimizing a datalog program under uniform equivalence that we illustrate by an example (and later adapt for our information gathering plan minimization). Suppose that we have the following datalog program:

$$\begin{aligned} r1: p(X) &:- p(Y) \wedge j(X, Y) \\ r2: p(X) &:- s(Y) \wedge j(X, Y) \\ r3: s(X) &:- p(X) \end{aligned}$$

We can check to see if  $r1$  is redundant by removing it from the program, then instantiating its body to see if the remaining rules can derive the instantiation of the head of this rule through a simple bottom-up evaluation (Ullman, 1989). Our initial assignment of relations is  $p(\text{“Y”})$ ,  $j(\text{“X”}, \text{“Y”})$ . If the remaining rules in the datalog program can derive  $p(\text{“X”})$  from the assignment above, then we can safely leave rule  $r1$  out of the datalog program. This is indeed the case. Given  $p(\text{“Y”})$  we can assert  $s(\text{“Y”})$  via rule  $r3$ . Then, given  $s(\text{“Y”})$  and  $j(\text{“X”}, \text{“Y”})$ , we can assert  $p(\text{“X”})$  from rule  $r2$ . Thus the above program will produce the same results without rule  $r1$  in it.

## 5. Greedy minimization of recursive plans

We now adapt the algorithm for minimizing datalog programs under uniform equivalence to remove redundant sources and unnecessary recursion from the information gathering plans. Our first step is to transform the query plan such that the query predicate is directly related to the source calls. This is done by removing global schema predicates, and replacing them with bodies of source inversion rules that define those predicates (see Ullman, 1989, Section 13.4).<sup>7</sup> Our example plan  $\mathcal{P}$ , from Section 3, after this transformation with the LCW

statements in Section 4 looks as follows:

$$\begin{aligned}
 r2 : \text{query}(S, A) &:- \text{ADDB}(S, A) \wedge A = \text{“Weld”} \\
 r3 : \text{query}(S, A) &:- \text{dom}(S) \wedge \text{CONDB}(S, A) \wedge A = \text{“Weld”} \\
 r4 : \text{dom}(S) &:- \text{ADDB}(S, A) \\
 r5 : \text{dom}(A) &:- \text{ADDB}(S, A) \\
 r6 : \text{dom}(A) &:- \text{dom}(S) \wedge \text{CONDB}(S, A)
 \end{aligned}$$

We are now ready to consider minimization. Our basic idea is to iteratively try to remove each rule from the information gathering plan. At each iteration, we use the method of replacing information source relations with their views or LCW's as in the rule subsumption check (see previous section) to transform the removed rule into a representation of what could possibly be gathered by the information sources in it, and transform the remaining rules into a representation of what is guaranteed to be gathered by the information sources in them. Then, we instantiate the body of the transformed removed rule and see if the transformed remaining rules can derive its head. If so, we can leave the extracted rule out of the information gathering plan, because the information sources in the remaining rules guarantee to gather at least as much information as the rule that was removed. The full algorithm is shown in figure 3.

For our example plan above, we will try to prove that rule  $r3$ , containing an access to the source CONDB, is unnecessary. First we remove  $r3$  from our plan, then transform it and the remaining rules so they represent the information gathered by the information sources in them. For the removed rule, we want to replace each information source in it with a representation of all the possible information that the information source could return.

Replace all global schema predicates in  $\mathcal{P}$   
with bodies of their inversion rules.

**repeat**

  let  $r$  be a rule in  $\mathcal{P}$  that has not yet been considered

  let  $\hat{\mathcal{P}}$  be the program obtained by deleting rule  $r$  from  $\mathcal{P}$   
and simplifying it by deleting any unreachable rules.

  let  $\hat{\mathcal{P}}'$  be  $\hat{\mathcal{P}}[s \mapsto s \vee l]$

  let  $r'$  be  $r[s \mapsto s \wedge v]$

**if** there is a rule,  $r_i$  in  $r'$ ,

  such that  $r_i$  is uniformly contained by  $\hat{\mathcal{P}}'$

**then** replace  $\mathcal{P}$  with  $\hat{\mathcal{P}}$

**until** each rule in  $\mathcal{P}$  has been considered once

Figure 3. The greedy plan minimization algorithm.

Specifically, we want to transform it to  $r[s \mapsto s \wedge v]$ . This produces:

$$\begin{aligned} \text{query}(S, A) &:- \text{dom}(S) \wedge \text{CONDB}(S, A) \\ &\quad \wedge \text{advisor}(S, A) \wedge A = \text{“Weld”} \end{aligned}$$

For the remaining rules,  $\mathcal{P} - r_3$ , we transform them into  $\mathcal{P}' = (\mathcal{P} - r_3)[s \mapsto s \vee l]$ , which represents the information guaranteed to be produced by the information sources in the rules. For our example, we produce:

$$\begin{aligned} r21 : \text{query}(S, A) &:- \text{ADDB}(S, A) \wedge A = \text{“Weld”} \\ r22 : \text{query}(S, A) &:- \text{advisor}(S, A) \wedge A = \text{“Weld”} \\ r23 : \text{query}(S, A) &:- \text{advisor}(S, A) \wedge A = \text{“Hanks”} \\ \text{dom}(S) &:- \text{ADDB}(S, A) \\ \text{dom}(S) &:- \text{advisor}(S, A) \\ \text{dom}(A) &:- \text{ADDB}(S, A) \\ \text{dom}(A) &:- \text{advisor}(S, A) \\ \text{dom}(A) &:- \text{dom}(S) \wedge \text{CONDB}(S, A) \\ \text{dom}(A) &:- \text{dom}(S) \wedge \text{advisor}(S, A) \end{aligned}$$

When we instantiate the body of the transformed removed rule  $r_3$ , we get the ground terms:  $\text{dom}(\text{“S”})$ ,  $\text{CONDB}(\text{“S”}, \text{“A”})$ ,  $A = \text{“Weld”}$ ,  $\text{advisor}(\text{“S”}, \text{“A”})$ . After evaluating  $\mathcal{P}'$  the remaining rules given with these constants, we find that we can derive  $\text{query}(\text{“S”}, \text{“A”})$ , using the rule  $r22$ , which means we can safely leave out the rule  $r_3$  that we've removed from our information gathering program.

If we continue with the algorithm on our example problem, we will not be able to remove any more rules. The remaining  $\text{dom}$  rules can be removed if we do a simple reachability test from the user's query, as they are not referenced by any rules reachable from the query.

### 5.1. Handling inter-source subsumption relations

The algorithm above only makes use of LCW statements that describe sources in terms of the global schema. It is possible to incorporate inter-source subsumption statements into the minimization algorithm. Specifically, suppose we are considering the removal of a rule  $r$  containing a source relation  $s$  from the plan  $P$ . Let  $U$  be the set of inter-source subsumption statements that have  $s$  in the tail, and  $U^{\leftarrow \mapsto}$  be the statements of  $U$  with the  $\leftarrow$  notation replaced by  $:-$  notation (so  $U$  is a set of datalog rules). We have to check if  $r[s \mapsto s \wedge v]$  is uniformly contained in  $(P - r + U^{\leftarrow \mapsto})[s \mapsto s \vee l]$ . If so, then we can remove  $r$ .

As an example, suppose we know that  $s_1$  and  $s_2$  are defined by the views:

$$\begin{aligned} s_1(x) &:- r(x) \\ s_2(x) &:- r(x) \end{aligned}$$

Suppose we know that  $s_1$  contains all tuples that  $s_2$  contains. This corresponds to the statement

$$s_1(x) :- s_2(x)$$

Suppose we have the query:

$$Q(x) :- r(x)$$

The corresponding maximally contained plan  $P$  will be:

$$Q(x) :- r(x)$$

$$r(x) :- s1(x)$$

$$r(x) :- s2(x)$$

To recognize that we can remove third rule from this plan because of the source subsumption statements, we check if that rule is uniformly contained in the program  $(P - r + "s1(x):s2(x)")$ , which is:

$$Q(x) :- r(x)$$

$$r(x) :- s1(x)$$

$$s1(x) :- s2(x)$$

The uniform containment holds here since if we add the tuple  $s2(A)$  to the program, bottom up evaluation allows us to derive  $s1(A)$  and subsequently  $r(A)$ , thus deriving the head of the removed rule.

## 5.2. Heuristics for ordering rules for removal

The final information gathering plan that we end up with after executing the minimization algorithm will depend on the order in which we remove the rules from the original plan. In the example given in Section 5.1, suppose we had another LCW statement:

$$\text{CONDB}(S, A) \leftarrow \text{advisor}(S, A)$$

In such a case, we could have removed  $r2$  from the original information gathering plan  $\mathcal{P}$ , instead of removing  $r3$ . Since both rules will lead to the generation of the same information, the removal would succeed. Once  $r2$  is removed however, we can no longer remove  $r3$ . This is significant, since in this case, a plan with rule  $r3$  in it is much costlier to execute than the one with rule  $r2$  in it. The presence of  $r3$  triggers the *dom* recursion through rules  $r4 \dots r6$ , which would have been eliminated otherwise. Recursion greatly increases the execution cost of the plan, as it can generate potentially boundless number of accesses to remote sources (see Section 7). We thus consider for elimination rules containing non-recursive predicates before those containing recursive predicates (such as *dom* terms). Beyond this, we also consider any gathered statistics about the access costs of the sources (such as contact time, response time, probability of access etc.) to break ties (Lambrech and Kambhampati, 1997).

*Complexity of Minimization.* The complexity of the minimization algorithm in figure 3 is dominated by the cost of uniform containment checks. As Sagiv (Sagiv, 1988) points out,

the running time of the uniform containment check is in the worst case exponential in the size of the query plan being minimized. However, things are brighter in practice since the exponential part of the complexity comes from the “evaluation” of the datalog program. The evaluation here is done with respect to a “small” database – consisting of the grounded literals of the tail of the rule being considered for removal. Nevertheless, the exponential complexity justifies our greedy approach for minimization, as finding a globally minimal plan would require considering all possible rule-removal orders.

## 6. Plan execution

Once the datalog query plan has been minimized to remove any redundant source accesses, *Emerac* attempts to execute the minimized plan. In this section, we describe how the techniques for datalog plan execution (c.f. Ullman, 1989) are adapted to the information gathering scenario to efficiently execute *Emerac*’s information gathering plans.

Two efficient approaches for executing datalog programs are (1) top-down relational evaluation and (2) bottom-up evaluation with magic sets transformation (Ullman, 1989). In *Emerac* we use the top-down relational evaluation. The top-down relational evaluation scheme attempts to avoid the inefficiencies of the top-down tuple-by-tuple evaluation scheme by directly manipulating relations (c.f. Ullman, 1989, Algorithm 12.17). The standard version of this scheme involves generating a rule/goal graph for the datalog program and evaluating the graph until fix point. To make this evaluation feasible as well as more efficient, a “conjunct ordering” algorithm is used to re-order the conjuncts (Morris, 1988).

In order to adapt the top-down relational evaluation to information gathering, we make the following extensions to it:

- We provide a framework for modeling the access restrictions on the source relations. These restrictions include attributes that must be bound in order to access the relation, as well as those attributes whose bindings will be ignored by the source.
- We describe a novel conjunct ordering approach that takes into consideration the access restrictions, and qualitative costs in reordering the rules (and thereby the rule/goal graphs).

In the rest of this section, we elaborate on these two extensions. We should mention here that in addition to these main changes, we also make another minor but important change to the evaluation of the rule goal graph. The plan is executed by traversing the relational operator graph. When ever a *union* node is encountered during traversal of the rule/goal graph, new threads of execution are created to traverse the children of the node in parallel. Use of separate threads allows us to reduce the response time as well as return answers to the user asynchronously.

### 6.1. Ordering source calls during execution

A crucial practical choice we have to make during the evaluation of datalog programs is the order in which predicates are evaluated. Our objective is to reduce the “cost” of execution, where cost is a function of the access cost (including connection time), traffic costs (the

number of tuples transferred), and processing cost (the time involved in processing the data). Typically, traffic and processing costs are closely correlated.

In our cost model, we assume that the access cost dominates the other terms. This is a reasonable assumption given the large connection setup delays involved in accessing sources on the Internet. While the traffic costs can also be significant, this is offset to some extent by the fact that many data sources on the Internet do tend to have smaller extractable tables.<sup>8</sup>

Although the source call ordering problem is similar to the “join ordering” phase in the traditional database optimization algorithms (Chaudhuri, 1998), there are several reasons why the traditional as well as distributed-database techniques are not suitable:

- Join ordering algorithms assume that all sources are relational databases. The sources on the Internet are rarely fully relational and tend to support limited types of queries. These limitations need to be represented and respected by the join ordering algorithm.
- Join ordering algorithms in distributed databases typically assume that the cost of query execution is dominated by the number of tuples transferred during execution. Thus, the so-called “bound-is-easier” assumption makes good sense. In the Internet information gathering scenario, the cost of accessing sources tends to dominate the execution cost. Consequently, we cannot rely solely on the bound-is-easier assumption and would need to consider the number of source calls.
- Typically, join ordering algorithms use statistics about the sizes of the various predicates to compute an optimal order of joining. These techniques are not applicable for us as our predicates correspond to source relations, about which we typically do not have complete statistics.
- The fact that source latencies make up a significant portion of the cost of execution argues for parallel (or “bushy”) join trees instead of the “left-linear” join trees considered by the conventional algorithms (Chaudhuri, 1998).

**6.1.1. Representing source access capabilities.** As we mentioned, in the information gathering scenarios, the assumption that information sources are fully relational databases is not valid. An information source may now be a wrapped web page, a form interfaced database, or a fully relational database. A wrapped web page is a WWW document interfaced through a wrapper program to make it appear as a relational database. The wrapper retrieves the web page, extracts the relational information from it, and then answers relational queries. A form-interfaced database refers to a database with an HTML form interface on the web which only answers selection queries over a subset of the attributes in the database. A WWW airline database that accepts two cities and two dates and returns flight listings is an example of a form interfaced database.

In our system, we use a simple way to inform the gatherer as to what types of queries an information source would accept.<sup>9</sup> We use the “\$” annotation to identify variables that must be bound, and “%” annotation to identify unselectable attributes (i.e., those attributes whose bindings cannot be pushed to the source to narrow down the selection). Thus a fully relational source would be adorned  $source(X, Y)$ , a form interfaced web-page that only accepts bindings for its first argument would be adorned  $source(X, \%Y)$ , while a

wrapped web-page source would have all its attributes marked unselectable, represented as  $source(\%X, \%Y)$ . Finally, a form interfaced web-page that requires bindings for its first argument, and is able to do selections only on the second argument would be adorned as  $source(\$X, Y, \%Z)$ .

Often times, a single source might support multiple binding patterns. These are supported by listing all the feasible binding patterns for that source. For example, if source  $S_1$  has two binding patterns:  $S_1(\$X, Y, \%Z)$ ,  $S_1(X, \$Y, \%Z)$ , it means that  $S_1$  can be accessed either with  $X$  bound or with  $Y$  bound. In either case, the attribute  $Z$  cannot be selected at the source (and must be filtered locally).

The “\$” and “%” annotations are used to identify feasible binding patterns for queries on a source, to establish generality relations between two binding patterns, and to ensure that soundness is preserved in pushing variable selection constraints (such as “ $Y = 7$ ”) into source calls. Given a source with annotations  $S_1(\$X, \%Y, Z)$ , only the binding patterns of the form  $S_1^{b--}$  are feasible (where “-” stands for either *bound* or *free* argument). Similarly, we are not allowed to push selection constraints on  $Y$  to the source  $S_1$  (they must be filtered locally). Thus the call  $S_1^{bbf}$  is modeled as  $S_1^{bff}$  filtered locally with the binding on  $Y$ .

A binding pattern  $S^p$  is more general than  $S^q$  (written  $S^p \succ_g S^q$ , if every selectable (non “%”-annotated) variable that is free in  $q$  is also free in  $p$ , but not *vice versa*). Thus, for the source  $S_1$  above, the binding pattern  $S_1^{bbf}$  is more general than the binding pattern  $S_1^{bfb}$  (while such a relation would not have held without “%” annotations). Intuitively, the more general a binding pattern, the higher the number of tuples returned by the source when called with that binding pattern. We ignore binding status of “%”-annotated variables since by definition they will not have any effect on the amount of data transferred. Finally, we define  $\#(\alpha)$  as the number of bound variables in  $\alpha$  that are not %-annotated. Notice that “ $\succ_g$ ” holds only between binding patterns of the same source while “ $\#(.)$ ” can be used to relate binding patterns of different sources.

**6.1.2. Plans and costs.** In *Emerac*, we support simple select/project/join queries on the information sources. Given the access restrictions on the sources, plans for such queries involve computing “dependent joins” (c.f. Chaudhuri and Shim, 1993; see below) over the source relations. Moreover, it may not always be feasible to push selections to the sources. As an example, consider two data sources  $S_1(\%X, Y)$ ,  $S_2(\$Y, Z)$  that export the relations  $R_1(X, Y)$  and  $R_2(Y, Z)$  respectively. Suppose we have a query

$$Q(X, Z) : -R_1(X, Y), R_2(Y, Z), X = “a”$$

The query rewriting phase in *Emerac* will convert this to

$$Q(X, Z) : -S_1(X, Y), S_2(Y, Z), X = “a”$$

It would seem that a good execution plan for the query would be to push the selection over  $X$  to  $S_1$  and do join between the two sources (in any order). However, since  $S_1$  has an unselectable attribute restriction on  $X$ , it is not possible to do a selection at the source. Further, since  $S_2$  requires  $Y$  to be bound, we need to do a dependent join between  $S_1$  and

$S_2$ . A feasible plan for this query would thus be:

$$\sigma_{X="a"}(S_1^{ff}(X, Y)) \stackrel{Y}{\bowtie} S_2^{bf}(Y, Z)$$

Here source  $S_1$  is being called with the binding pattern  $ff$ , and the results are processed (at the mediator) with a selection on  $X = "a"$ . Next, the source  $S_2$  is called with the binding pattern  $bf$ , where calls are issued once for each unique value of  $Y$  from the left sub-tree. (This form of passing bindings from the left subtree of a join to the right subtree is called a "dependent join," and has been studied in the literature in connection with optimization in the presence of foreign functions (Chaudhuri and Shim, 1993).)

In computing the plans, we are thus interested in deciding not only the order in which the joins are carried out, as is the case in traditional system-R style query optimization (Chaudhuri, 1998), but also about what specific binding patterns are used in source calls, and how far it is feasible to push selections into query plans. We note that the execution cost is a function of the access cost (including connection time), traffic costs (the number of tuples transferred), and processing cost (the time involved in processing the data). Thus, optimal plans will need to minimize:

$$\sum_s (C_a^s * n_s + C_t^s * D_s)$$

where  $n_s$  is the number of times a source  $s$  has been accessed during the plan and  $C_a^s$  is the cost per access, and  $C_t^s$  is the per tuple transfer cost for source  $s$ , and  $D_s$  is the number of tuples transferred by  $s$ . We note that this cost metric imposes a tension between the desire to reduce network traffic, and the desire to reduce access costs. To elaborate, reducing the network traffic involves accessing sources with less general binding patterns. This in turn typically increases the number of separate calls made to a source, and leads to increased access cost. To illustrate this further, consider the subgoals:

$$S_1(X, Y) \wedge S_2(Y, Z)$$

Suppose that the query provides bindings for  $X$ . How should we access the sources? The conventional wisdom says that we should access  $S_1$  first since it has more bound arguments. As a result of this access, we will have bindings for  $Y$  which can then be fed into calls to  $S_2$ . The motivation here is to reduce the costs due to network traffic. However, calling  $S_1$  and using its outputs to bind the arguments of  $S_2$  may also lead to a potentially large number of separate calls to  $S_2$  (one per each of the distinct  $Y$  values returned by  $S_1$ ),<sup>10</sup> and this can lead to a significant connection setup costs, thus worsening the overall cost. On the other hand, calling  $S_2$  without propagating bindings from  $S_1$  would reduce the source calls to two. We need to thus consider both the access costs and the traffic costs to optimize the ordering of the sources.

Since we often do not have the requisite statistics to do the full cost-based optimization, we propose an approach that is less dependent on full statistics. Our algorithm does not make use of a quantitative measure of access cost or transfer costs, but rather a qualitative measure of high and low cost of access to a source. We describe this approach in the next section.



**6.1.3. An algorithm for ordering source calls.** We make two important assumptions in designing our source-call ordering algorithm:

- Exact optimization of the execution cost requires access to source selectivity statistics. While such statistics may be available for intra-corporation information integration scenarios (c.f. GARLIC (Haas et al., 1997)), they are harder to get in the case of autonomous and decentralized sources on the Internet.
- We assume that by default source access costs (rather than network traffic) are the dominating cost of a query plan. This becomes reasonable given the large connection setup delays involved in accessing sources on the Internet. Many Internet sources tend to have small extractable tables which help offset the traffic costs that at times can be proportional to or greater than access cost.<sup>11</sup>

If our assumptions about the secondary importance of network traffic costs were always true, then we can issue calls to any source as soon as its binding constraints are met (i.e., all the variables requiring bindings have bindings available). Furthermore, we need only access the source with the most general feasible binding pattern (since this will reduce the number of accesses to the source). We do provide an escape clause for this assumption (see below), as sometimes sources can transfer arbitrarily large amounts of data for calls with sufficiently general binding patterns.

*High-traffic Binding Patterns.* To ensure that we don't get penalized excessively for focusing concentrating primarily on access costs, we also maintain a table, called HTBP, of least general (w.r.t. " $>_g$ ") source binding patterns that are still known to be high-traffic producing. The general idea is to postpone calling a source as long as all the feasible binding patterns for that source supported by the currently bound variables are equal to or more general than a binding pattern listed in HTBP.

An underlying assumption of this approach is that while full source statistics are rarely available, one can easily gain partial information on the types of binding patterns that cause excessive traffic. For example, given a source that exports the relation

*Book(Author, Title, ISBN, Subject, Price, Pages)*

we might know that calls that do not bind at least one of the first four attributes tend to generate high traffic. The information as to which binding patterns generate high traffic could come either from source modeling phase, or could be learned with rudimentary probing techniques. The latter approach involves probing the sources with a set of sample queries, and logging for each source the binding patterns and the cardinalities of generated result sets, identifying the HTBP patterns using a threshold on the result set cardinality.

There are two useful internal consistency conditions on HTBP. First, if  $S^\alpha$  is listed in HTBP (where  $\alpha$  is a binding pattern on  $S$ ), then every  $S^\beta$  where  $\beta >_g \alpha$  is also implicitly in HTBP. Similarly, if  $S^\alpha$  is in HTBP, then it cannot be the case that the only free variables in  $\alpha$  are all "%" -annotated variables.

A greedy algorithm to order source calls based on these ideas appears in figure 4. It is along the lines of "bound-is-easier" type ordering procedures (Morris, 1988; Levy et al.,

```

Inputs: FBP: table of forbidden binding patterns
HTBP: table of high traffic binding patterns
V := all variables bound by the head; V' :=  $\emptyset$ 
C[1 .. m] :=
Array where C[i] lists sources chosen at  $i^{th}$  stage;
P[1 .. m] :=
Array where P[i] lists sources postponed at  $i^{th}$  stage
for i := 1 to m (where m is the number of subgoals)
do begin
C[i] :=  $\emptyset$ ; P[i] :=  $\emptyset$ ; V := V  $\cup$  V'
for each unchosen subgoal S
do begin
B := All feasible binding patterns for S w.r.t. V
and FBP sorted using " $\succ_g$ " relation.
for each  $\beta \in B$ 
do begin
if  $\nexists \beta' \in HTBP$  s.t. ( $\beta = \beta'$ )  $\vee$  ( $\beta \succ_g \beta'$ )
then begin
Push S with binding pattern  $\beta$  into C[i];
Mark S as "chosen";
add to V' all variables appearing in S;
end
end
if B  $\neq \emptyset$  and S is not chosen
then Push  $S^\gamma$  into P[i], where
 $\gamma \in B$  has the maximum  $\#(.)$  value;
end
if C[i] =  $\emptyset$  and P[i]  $\neq \emptyset$ 
then begin
Take the source  $S^\beta \in P[i]$  with maximum  $\#(.)$ 
value and push it into C[i];
add to V all variables appearing in S;
else fail
end
end
Return the array C[1..i].

```

Figure 4. A greedy source call ordering algorithm that considers both access costs and traffic costs.

1996). By default, it attempts to access each source with the most general feasible binding pattern. This default is reasonable given our assumption that access costs dominate transfer costs. The default is overridden if a binding pattern is known to produce too much traffic—by being present in HTBP (or being more general than a pattern present in HTBP).

The procedure takes as input a rule with  $m$  subgoals and a given binding pattern for its head. The input also includes FBP, a table of forbidden binding patterns for each source (constructed from the “\$” annotations), and the table HTBP, which contains all source binding patterns that are known to be high-traffic producing. At each level  $i$ , the algorithm considers the feasible binding patterns of each unchosen source from most general to least general, until one is found that is not in HTBP. If such a binding pattern is found, that source, along with that binding pattern, is added to the set of selected sources at that level (maintained in the array of sets  $C[i]$ ). If not, the source, along with the least general feasible binding pattern (given the “\$” restrictions as well as the currently available bound variables), is temporarily postponed (by placing it in  $P[i]$ ). If at the end of considering all unchosen sources at level  $i$ , we have not chosen at least one source (i.e., all of them have only high-traffic inducing binding patterns), then one of the sources from the list of postponed sources ( $P[i]$ ) is chosen and placed in  $C[i]$ . This choice is made with the “bound-is-easier” assumption by selecting the source with the least general binding pattern (in terms of  $\#(\cdot)$ ).

Anytime a source is placed in  $C[i]$ , the set  $V$  of variables that currently have available bindings is updated.<sup>12</sup> This ensures that at the next stage more sources will have feasible, as well as less general, binding patterns available. This allows the algorithm to progress since less general binding patterns are also less likely to be present in the HTBP table. Specifically, when  $C[i]$  is empty and  $P[i]$  is non-empty, we push only one source from  $P[i]$  into  $C[i]$  since it is hoped that the updated  $V$  will then support non-high-traffic binding patterns at the later stages. (By consulting HTBP, and the set of unchosen sources, the selection of the source from  $P[i]$  can be done more intelligently to ensure that the likelihood of this occurrence is increased).

Notice that each element of  $C$  is a (possibly non-singleton) set of source calls with associated binding patterns (rather than a single source call)—thus supporting parallel source calls which reduce the time spent on connection delays. Thus,  $C[i]$  may be non-empty for only a prefix of values from 1 to  $m$ . The complexity of our ordering algorithm is  $O(n^2)$  where  $n$  is the length of the rule. Note that HTBP table determines the behavior of our algorithm. An empty HTBP table makes our algorithm to focus on reducing source accesses (similar to Yerneni and Li (1999)), whereas presence of all source binding patterns in HTBP table makes the algorithm focus on reducing network traffic by using a variant of bound-is-easier (Morris, 1988). When some source patterns are present in the HTBP table our algorithm attempts to reduce both access and transfer costs, as appropriate.

**6.1.4. Example.** Figure 5 shows an example illustrating the operation of the source-call ordering procedure. We have two sources in the query plan: DP(A,T,Y), which is a source of all database papers, and SM98(T,U), which is a source of all papers from SIGMOD-98. The query binds the year to 1998, and wants the author, title and URL tuples. We will illustrate the algorithm under three different scenarios.

*Case 1.* In the first case, shown on the left, HTBP contains  $DP^{bbb}$ , and  $SM98^{bb}$ . Notice that this means that every possible call to these sources is considered to be high-traffic. Given that the query binds one variable,  $Y$ , the only possible source call bindings we have are:  $DP^{fff}$ ,  $DP^{ffb}$ , and  $SM98^{ff}$ . Among these, the algorithm finds no possible feasible source call that is not in HTBP—all of them are more general than the source call patterns stored in HTBP. Thus, it winds up picking up  $DP^{ffb}$  since this is the one with most bound variables. At this point, the second iteration starts with one remaining source  $SM98$ , and bindings for three variables,  $A, T, Y$  (where  $A$  and  $T$  are supplied by the first call). The two possible source calls are  $SM98^{ff}$  and  $SM98^{bf}$ , both of which are again in the HTBP. So, the algorithm picks  $SM98^{bf}$ . The query plan thus involves doing a dependent join (Chaudhuri and Shim, 1993) between  $DP^{ffb}$  and  $SM98^{bf}$ , with the unique titles retrieved by the DP tuples being used to invoke  $SM98$  ( $DP^{ffb}(A, T, Y) \xrightarrow{T} SM98^{bf}(T, U)$ )

*Case 2.* In the second case, only the call  $DP^{ffb}$  is in the HTBP. Thus, in the first iteration, neither of the DP calls are feasible, but the  $SM98$  call is. So,  $SM98^{ff}$  is chosen. In the second iteration, we have four DP calls, of which three are in HTBP. The call  $DP^{fbf}$  is however not in HTBP, and is thus feasible. The query plan thus involves a dependent join between  $SM98$  and DP with the URL values from  $SM98$  call being passed to DP ( $SM98^{ff}(T, U) \xrightarrow{T} DP^{fbf}(A, T, Y)$ ).

*Case 3.* In the third case, HTBP is empty. Thus, we simply pick the most general feasible source calls in the very first iteration—leading to calls  $SM98^{ff}$  and  $DP^{fff}$ . The query plan is thus a (non-dependent) join between the sources  $SM98$  and DP ( $DP^{fff}(A, T, Y) \bowtie SM98^{ff}(T, U)$ ).

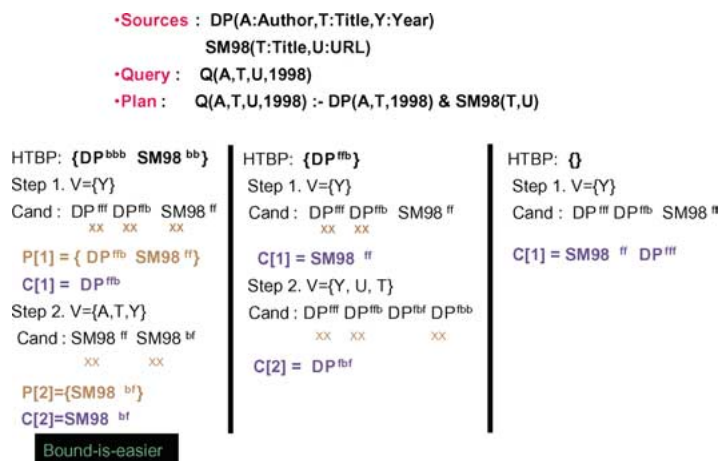


Figure 5. Example illustrating the operation of the source-call ordering procedure. Note that based on the contents of the HTBP table, the procedure can degenerate into bound-is-easier type ordering.

**6.1.5. Converting the call graph to a join tree.** When the algorithm in figure 4 terminates, the array  $C$  specifies which sources are to be called at each stage, and what binding patterns are to be used in those calls—specifically, all the source calls in  $C[i]$  are issued in *parallel* before those in  $C[i + 1]$ . There is still the matter of what is the exact *join tree* that is to be executed at the mediator to derive answers from these source calls.

Consider running the algorithm on the query with the given binding restrictions

$$Q(X, Y, W, Z) : \neg S_1^{ff}(X, Y) \wedge S_2^{bf}(Y, Z) \wedge S_3^{ff}(T, W) \wedge S_4^{bf}(W, Z)$$

We also assume that HTBP is empty.<sup>13</sup> Our algorithm will end with  $C[1] = \{S_1^{ff}, S_3^{ff}\}$  and  $C[2] = \{S_2^{bf}, S_4^{bf}\}$ . Analyzing the common variables among the sources, it is easy to represent this as a call graph as shown in figure 6—which has both directed and undirected arcs. Arcs correspond to shared variables between source calls. A directed arc from  $S_1$  to  $S_2$  states that a shared variable needs to be bound by  $S_1$  before reaching  $S_2$ .

The join tree can be greedily derived from the call graph by combining vertices in the graph till it becomes a single node. The vertex combination is done by traversing the graph in a topologically sorted order. All the vertices that have directed arcs between them are first combined (by doing a dependent join (Chaudhuri and Shim, 1993) between the corresponding sources). When the resulting graph doesn't have any directed arcs, then the undirected arcs are processed by converting them to joins. Finally, if we are left with a disconnected graph, the corresponding subtrees are joined by a cartesian product. Accordingly converting the example in figure 6 gives us the following plan:

$$\left( S_1^{ff}(X, Y) \overset{Y}{\bowtie} S_2^{bf}(Y, Z) \right) \bowtie \left( S_3^{ff}(T, W) \overset{Z}{\bowtie} S_4^{bf}(W, Z) \right)$$

As the example above shows, our approach supports bushy join trees (instead of sticking merely to left-linear join trees). As is evidenced by this example (and pointed out first in Florescu et al. (1999)), bushy join trees allow us to avoid cartesian products in more situations than left linear join trees alone would—when there are binding restrictions on sources.

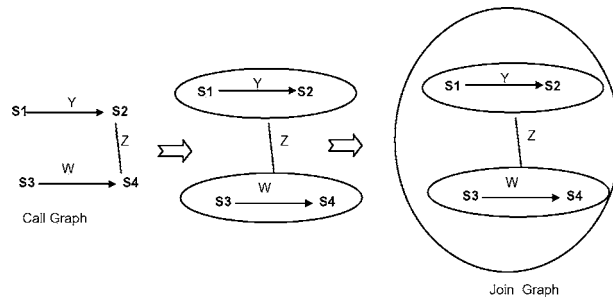


Figure 6. Converting a call graph into a join graph.

## 7. Implementation and evaluation

We will start by describing the architecture and implementation of *Emerac* (Section 7.1). Next, we will discuss a variety of experiments we conducted to evaluate its effectiveness. Section 7.2 evaluates the effectiveness of the plan minimization routines described in Section 5 in terms of their costs and benefits. Section 7.3 describes the experiments we conducted to evaluate the techniques for improving the plan execution that we presented in Section 6. This section starts with an empirical validation of our assumptions about the domination of source access costs (Section 7.3.1). The source call ordering scheme is evaluated over simulated sources in Section 7.3.2, and over sources on the Internet in Section 7.3.3.

### 7.1. Architecture of Emerac

*Emerac* is written in the Java programming language, and is intended to be a library used by applications that need a uniform interface to multiple information sources. Full details of *Emerac* system are available in Lambrecht (1998). *Emerac* presents a simple interface for posing queries and defining a global schema. *Emerac* is internally split into two parts: the query planner and the plan executor. The default planner uses algorithms discussed in this paper, but it can be replaced with alternate planners. The plan executor can likewise be replaced, and the current implementation attempts to execute an information gathering plan in parallel after transforming it into a relational operator graph.

The query planner accepts and parses datalog rules, materialized view definitions of sources, and LCW statements about sources. Given a query, the query planner builds a source complete information gathering plan (using the method from (Duschka and Genesereth, 1997)) and attempts to minimize it using the minimization algorithm presented in Section 5.

The optimized plan is passed to the plan executor, which transforms the plan into a relational operator graph. The plan executor makes use of “\$” and “%”-adornments to determine the order to access each information source in a join of multiple sources, as described in this paper. The plan is executed by traversing the relational operator graph. When a *union* node is encountered during traversal, new threads of execution are created to traverse the children of the node in parallel. Use of separate threads also allows us to return answers to the user asynchronously, facilitating return of first tuples faster.

*Handling Recursion during execution.* Since information gathering plans can contain recursion, handling recursive plans during execution becomes an important issue. Since each recursive call to a node in the *r/g* graph (Ullman, 1989) can potentially generate an access call to a remote source, evaluating a program until it reaches fix point can get prohibitively expensive. Currently, we take a practical solution to this problem involving depth-limited recursion. Specifically, we keep a counter on each node in the *r/g* graph to record how many times the node has been executed. When the counter reaches a pre-specified depth-limit, the node would not be executed, and an empty set will be returned to represent the result of executing the node. Since the recursion induced by the binding restrictions does not involve any negation in the tail of the rules, this strategy remains sound—i.e., will produce only correct answers.

*Wrapper Interface.* *Emerac* assumes that all information sources contain tuples of information with a fixed set of attributes, and can only answer simple *select* queries. To interface an information source with *Emerac*, a Java class needs to be developed that implements a simple standard interface for accessing it. The information source is able to identify itself so as to provide a mapping between references to it in materialized view and LCW definitions and its code.

In order to facilitate construction of wrappers for web pages, a tool was created to convert the finite state machine based wrappers created by SoftMealy (Hsu, 1998) into Java source code that can be compiled into information sources usable by *Emerac*. We have successfully adapted 28 computer science faculty listing web pages wrapped with SoftMealy into information sources usable by *Emerac*.

### 7.2. Evaluating the effectiveness of plan minimization

We used the prototype implementation of *Emerac* to evaluate the effectiveness of the optimization techniques proposed in this paper. We used two sets of experimental data. The first were a set of small artificial sources containing 5 tuples each. Our second data set was derived from the University of Trier's Database and Logic Programming (DBLP) online database, which contains bibliographical information on database-related publications. Individual sources used in the experiments corresponded to different subsets of DBLP data (ranging from 128 to 2048 tuples). In each case, some of the sources are unconstrained, while others have binding restrictions (leading to recursive plans). To normalize for differences caused by individual source implementations, we extracted the data into tables which we stored on disk as Java serialized data. All experiments were conducted using a simple wrapper (written in compiled Java) to return the contents of the serialized tables.

The sources delay answering each query for a set period of time in order to simulate actual latency on the Internet. In all our experiments, this delay was set to 2 seconds, which is quite reasonable in the context of current day Internet sources.

*Utility of minimization.* To see how the planner and executor performed with and without minimization, we varied the number of duplicate information sources available and relevant to the query, and compared the total time taken for optimization (if any) and execution. Given that the minimization step involves an exponential "uniform containment" check, it is important to ensure that the time spent in minimization is made up in improved execution cost. Notice that we are looking at only the execution time, and ignoring other costs (such as access cost for premium sources), which also can be reduced significantly with the minimization step. The naive method simply builds and executes source complete plans. The "LCW" method builds source complete plans, then applies the minimization algorithm described in Section 5 before executing the plans. For both methods, we support fully parallel execution at the union nodes in the *r/g* graph. Since in practice, recursive plans are handled with depth bounded recursion, we experimented with a variety of depth limits (i.e., the number of times a node is executed in the rule-goal graph), starting from 1 (which in essence prunes the recursion completely).

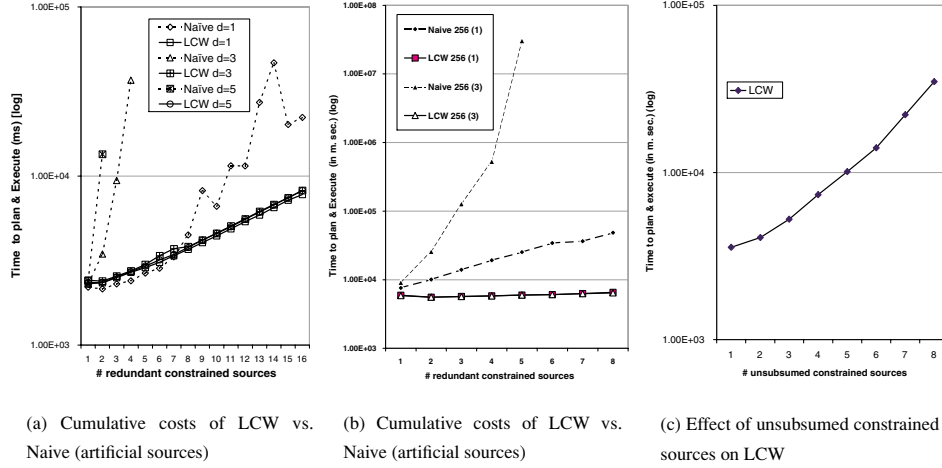


Figure 7. Results characterizing utility of minimization algorithm.

The plots in figure 7 show the results of our experiments. Plot *a* is for the artificial sources, and shows the relative time performances of LCW against the naive algorithm when the number of redundant constrained sources is increased. In this set of experiments, LCW statements allow us to prove all constrained sources to be redundant, and the minimization algorithm prunes them. The y-axis shows the cumulative time taken for minimization and execution. We note that the time taken by the LCW algorithm remains fairly independent of recursion depth as well as number of constrained sources. The naive algorithm, in contrast, worsens exponentially with increasing number of constrained sources. The degradation is more pronounced for higher recursion depths, with the LCW method outperforming the naive one when there are two or more redundant constrained sources. Plot *b* repeats the same experiment, but with the sources derived from the DBLP data. The sources are such that the experimental query returns upto 256 tuples. The experiment is conducted for recursion depth limits 1 and 3. We note once again, that LCW method remains fairly unaffected by the presence of redundant constrained sources, while the naive method degrades exponentially. Plot *c* considers DBLP data sources in a scenario where some constrained sources are left unsubsumed after the minimization. As expected, LCW performance degrades gracefully with increased number of constrained sources. Naive algorithm would not have shown such graceful degradation as no sources would be removed through subsumption.

### 7.3. Evaluating the effectiveness of source call ordering

In this section, we report on a set of experiments conducted on both artificial and real Internet sources to evaluate the effectiveness of our algorithm (referred to as HT). We compare the performance with two other greedy approaches for join ordering namely Bound-is-easier (BE) and Reduced Access (RA). For sources where HTBP information cannot be ascertained and hence is not available, HT will work as RA. The aim of our experiments is to demonstrate



that our algorithm outperforms algorithms that concentrate on reducing only tuple transfer cost (as in BE) or only source access cost (as in RA). BE corresponds roughly to the algorithm used in Information Manifold (Levy et al., 1996) while RA is similar to the algorithm proposed in Yerneni and Li (1999) for reducing source access costs. The experiments also indirectly show that the kind of coarse statistics that we assume in our algorithm are likely to be available easily in practice.

**7.3.1. Verifying the dominance of access cost.** An important assumption we made in developing the execution algorithm is about dominance of source access cost over tuple transfer cost for sources on the Internet. In this section we describe a simple empirical validation of this hypothesis. We consider access cost as the cost incurred in setting up a connection to a source. Transfer cost is incremental cost for each transaction. We considered two sources, an internet source <http://dvs1.dvllabs.com/adcritic/> and a local intranet source. We recorded the time taken to download various files with sizes ranging from 10KB to 25MB from both the sources. Each byte transfer is considered as a transaction. Average values were taken after 4 rounds of data transfer for each file.

Though access ( $a$ ) and transfer time ( $t$ ) for a source are not known, they can be calculated from a plot of downloading time vs. filesize as given in figure 8. Source access time ( $a$ ) is the  $y$ -intercept and tuple transfer time ( $t$ ) is the slope of the graph.

As can be seen from figure 8 the access cost for internet source (nearly 5 sec) is much higher than that for intranet source (92 msec). Tuple transfer time, on the other hand is nearly same for both types of sources. This validates our assumption that for data sources residing on the Internet, source access cost is considerably higher than tuple transfer cost.

**7.3.2. Source call ordering results for artificial sources.** Below we present results from the performance evaluation of HT, BE and RA algorithms for queries over sources mimicking Internet sources. We derived the sources from the relations present in the “Enterprise Schema” used in Elmasri and Navathe (1994). The sources are designed as Java servlets accessible from a server on the Intranet. The servlets accept a query and return relevant

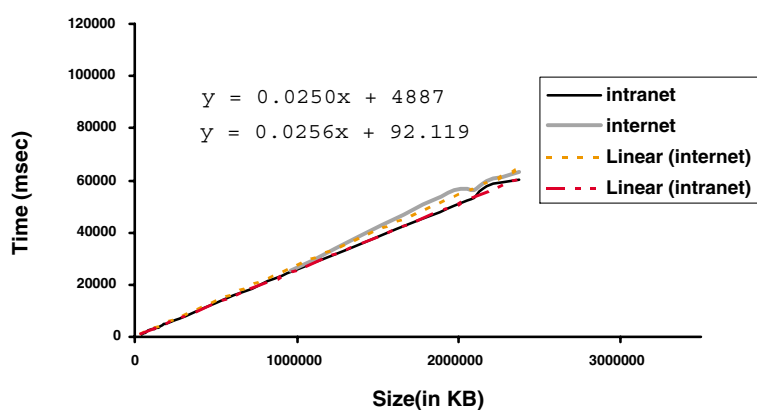


Figure 8. Comparison of access and transfer cost.

tuples extracted from data stored in flat files. To model the latency exhibited by Internet sources, we delay the response from the servlets. The delay is proportional to the number of tuples in the resultset. Some of the servlets were designed to also mimic the bursty behaviour (i.e., sending sets of tuples interspersed with delays (Urhan and Franklin, 1998)) shown by some Internet sources. A detailed description of sources in terms of their attributes and their forbidden binding patterns is given below:

*Employee*<sup>fffffff</sup> (*N* : *Name*, *S* : *SSN*, *B* : *Dateof Birth*, *A* : *Address*, *X* : *Sex*,  
*Y* : *Salary*, *U* : *Manager*, *D* : *DeptNo*)  
*Dept*<sup>fff</sup> (*Dn* : *DeptName*, *D* : *DeptNo*, *U* : *Manager*, *T* : *MrgStartDate*)  
*Deptloc*<sup>ff</sup> (*D* : *DeptNo*, *Dc* : *Location*)  
*Workson*<sup>fff</sup> (*S* : *SSN*, *P* : *ProjectID*, *H* : *Hours*)  
*Project*<sup>fff</sup> (*Pn* : *ProjectName*, *P* : *ProjectId*, *Dc* : *Location*, *D* : *DeptNo*).

The subscripts on source names describe forbidden binding pattern derived from “\$” and “%” annotations for the sources. Intuitively, FBP and HTBP will have fewer patterns than the set of feasible binding patterns for a source. Hence we use them to represent the infeasible/costly binding patterns thereby reducing the look up time for our experiments.

The Table 1 lists source binding patterns that generate large resultsets for most of the values given to the bound attributes. Each column in Table 1 is an HTBP table. Thus given Deptloc(D, “Houston”) and the HTBP table of Deptloc, one can see that pushing the value Location = “Houston” to the source Deptloc will result in high traffic. The source Dept. has an empty HTBP table. Hence any query with a binding pattern that is not forbidden for Dept can be issued to Dept. But we cannot assume a source with no HTBP to have low cost. Therefore we deliberately modeled Dept to simulate a source with high response time.

The graphs in figures 9–11 compare the average values of execution time, size of result sets and number of source calls for BE, HT and RA. The results were obtained by running queries Q1, Q2 and Q3 using these algorithms. The queries were generated with 3 different binding patterns per source and 5 different binding values per binding pattern. The times for running the source call ordering algorithms themselves were minute in comparison with the execution costs.

*Query1* : *Q(N, S, B, A, X, Y, U, D, Dn, U, T, Dc)*  
*Plan* : *Q(N, S, B, A, X, Y, U, D, Dn, U, T, Dc)* : –  
Employee(N,S,B,A,X,Y,U,D),Dept(Dn,D,U,T), Deptloc(D,Dc).

Table 1. HTBP for artificial sources.

Employee	Dept.	Deptloc	Workson	Project
f,f,f,f,f,f,f,f		b,b	b,f,f	f,f,b,f
f,f,f,f,f,f,f,b		b,f	f,b,f	f,f,b,b
f,f,f,f,f,f,b,b		f,b	f,f,b	

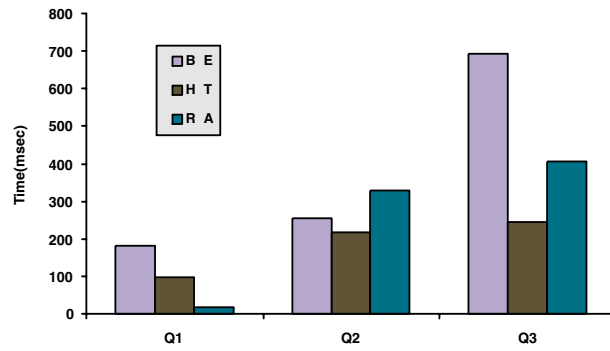


Figure 9. Comparison of algorithms w.r.t execution time for Q1,Q2,Q3.

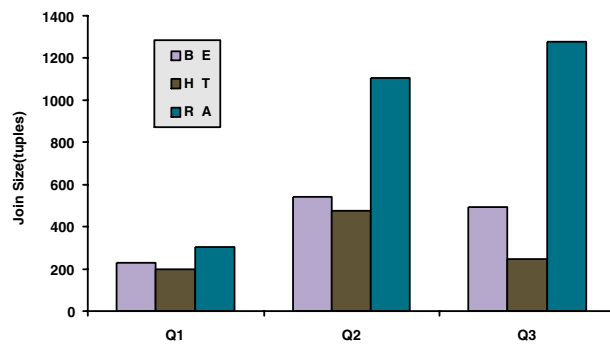


Figure 10. Comparison of algorithms w.r.t join size for Q1,Q2,Q3.

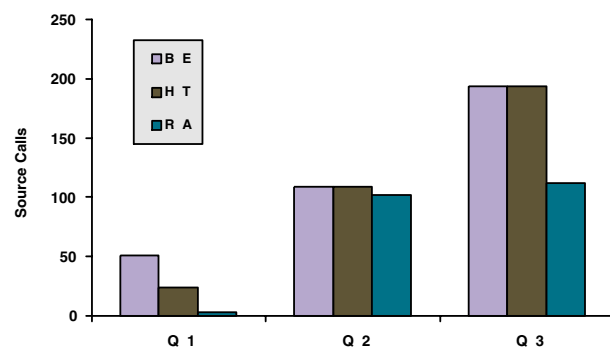


Figure 11. Comparison of algorithms w.r.t # of source calls for Q1,Q2,Q3.

*Query2*:  $Q(D, Dc, S, P, H, Pn)$   
*Plan*:  $Q(D, Dc, S, P, H, Pn)$ : –  
 Deptloc(D,Dc),Workson(S,P,H),Project(Pn,P,Dc,D).

*Query3*:  $Q(Dn, D, U, T, Dc, S, P, H, Pn)$   
*Plan*:  $Q(Dn, D, U, T, Dc, S, P, H, Pn)$ : –  
 Dept(Dn,D,U,T),Deptloc(D,Dc),Workson(S,P,H), Project(Pn,P,Dc,D).

From figure 11, we can see that RA always optimizes the number of source calls and has the least number of source calls for all 3 queries. BE on the other hand focuses on having smaller result sets and thus to reduce the transfer cost, as is evident from figure 10. But no approach is a clear winner when the total execution cost is considered (see figure 9). Our HT algorithm gives lowest execution cost for 2 out of the 3 cases considered. We can see that HT tries to reduce both number of source calls (figure 11) and/or resultset size (figure 10) while executing the queries. For Q2 and Q3, HT has source number of source calls equal to BE but smaller result sets compared to BE and RA and hence achieves lower execution cost. HT strikes a middle ground compared to BE and RA and tries to optimize both access cost and transfer cost. Thus HT generates low cost execution plans more often than BE and RA.

The case where RA is better than HT in figure 9 shows that HTBP is not always perfect. We consider a binding pattern as HTBP if it generates large resultsets for most binding values (without regard to attribute selectivities). It could well be the case that for a specific instantiation of the binding pattern (i.e. for particular value(s) of attribute(s)), a HTBP may not generate high traffic. RA which does not consider HTBP thus makes source calls using this binding pattern and emerges a winner.

**7.3.3. Source call ordering results for Internet sources.** The next set of tests were done on data sources derived from the DBLP Bibliography of AI and Database papers maintained by Michael Ley at <http://dblp.uni-trier.de>. We use a simple scenario with the execution plan using only 2 sources:

$Dp1^{fbff,ffff}$  (A:Author, Co:Co-Author, C:Conference, Y:Year)  
 $Dp2^{ffff}$  (A:Author, T:Title, C:Conference, Y:Year)

We derive two sources Dp1 and Dp2 shown above from the DBLP Bibliography by projecting the corresponding relations from DBLP. Specifically, we developed a wrapper over DBLP that accepts queries over relations of Dp1 and Dp2, forwards them to DBLP and extracts the relevant relation from the resultset given by DBLP. Queries accepted by the wrapper have to satisfy the FBP associated with the projected relation. The subscripts on source names give the forbidden binding patterns (FBP). HTBP for these sources are shown in Table 2. These HTBP statements can be determined by rudimentary probing techniques. Specifically, we execute queries on sources with various binding pattern/value combinations and log the resultset sizes and execution times. Given similar resultset sizes for two sources, we cannot deduce that both binding patterns are HTBP since the sources may have differing

Table 2. HTBP for real sources Dp1 and Dp2.

Dp1(A,Co,C,Y)	Dp2(A,T,C,Y)
F,B,B,B	F,F,F,B
	F,F,B,F
	F,B,F,F
	B,F,F,F

Table 3. Results of accessing Dp1, DP2 using BE and HT.

	BE	HT
Source calls	20.8/15.3	12.7/11.5
Join size	8.8/9.0	8.2/7.9
Time	30.5/21.8	18.8/16.4

processing power and database size. Hence we store the execution time for queries on a source and use these to determine the average response time and resultset size returned by the source. Any binding pattern that generates larger result sets than average resultset size or has considerably higher response time than the average case is considered HTBP. The binding pattern restriction for Dp1 is Dp1 (\$A, %Co, C, Y) and that for Dp2 is Dp2 (\$A, T, C, Y). The “%” annotation describes that the attribute has to be filtered locally. But both the sources must bind the attribute ‘Author’ to retrieve tuples. The experimental setup is thus:

*Query* :  $Q(A, Co, T, C, Y)$

*Plan* :  $Q(A, Co, T, C, Y) : - Dp2(A,T,C,Y), Dp1(A,Co,C,Y)$

Table 3 shows the performance of various algorithms for queries over Dp1 and Dp2. Performance is measured as the number of source calls made, resultset size (Joinsize) and the query response time (Time). The time incurred in executing the source call ordering algorithms were negligible compared to the execution costs. The results show that HT performs better than BE for these sources. The example also shows that the HTBP statistics can be collected for real Internet sources and that our algorithm does perform better than other existing source call algorithms.

## 8. Related work

As we mentioned, systems that consider integration in the context of information gathering use LAV approach to model sources. In the LAV approach, sources are seen as materialized views over the mediated schema. The user query, posed in terms of the mediator schema, has to be re-written solely in terms of source calls. Although this problem, on the surface, is similar to query rewriting in terms of materialized views (Chaudhuri et al., 1995), there are several important differences that complicate the query rewriting:

- We are only interested in rewritings that are entirely in terms of source relations.
- The sources may not contain all tuples satisfying their view definition. This leads to the so-called open-world assumption, and changes the objective of query planning from finding a “sound and complete” query plan to finding a “sound and maximally contained” (Duschka and Genesereth, 1997) query plan (a query plan  $P$  for the query  $Q$  is maximally contained if there is no other query plan  $P'$  for  $Q$  that can produce more answers for  $Q$  using the same set of sources).
- The materialized views represented by the sources may have a variety of access restrictions. In such a case, the maximally contained query plan may be a “recursive” query plan (or equivalently, an infinite union of conjunctive plans).

IM (Levy et al., 1996) and Occam (Kwok and Weld, 1996) are among the first systems to consider query planning in the LAV approach. Both these systems search through the space of conjunctive query plans, to find sound rewritings of the user query, and optimizing them for efficient execution.

There are two problems with the approaches used by IM and Occam, when some sources have access restrictions. To begin with, as shown by Duschka and Genesereth (1997); Duschka and Levy (1997), maximally contained plans will be recursive when we have sources with access restrictions. The approach of finding sound conjunctive query plans, used by IM and Occam essentially “unfolds” the recursion,<sup>14</sup> forcing them to handle infinite unions of conjunctive plans. IM gets around this by sacrificing guarantees of maximal containment. Specifically, as mentioned in Duschka and Levy (1997), while IM ensures that the query plans returned by the system are feasible (i.e., respect all access restrictions), it does not guarantee maximally contained plans. Occam (Kwok and Weld, 1996) was the first to formally recognize that the maximally contained plan may correspond to an infinite union of conjunctive query plans. It searches in the space of conjunctive query plans of increasing lengths, pruning candidate plans when they are found to be redundant.

The unfolding of recursion inherent in IM and Occam also leads to inefficient query plan execution. Specifically, the unfolded conjunctive query plans found by these algorithms tend to have a significant amount of overlapping structure, and executing them separately leads to significant redundant computation.

*Emerac* uses Duschka’s source inversion algorithm (Duschka, 1997) to generate a datalog query plan that is maximally contained with respect to the given query. The query minimization and optimization are done directly on the datalog query plan, without converting it to (a potentially infinite union of) conjunctive query plans. *Emerac* thus avoids the redundant processing involved in executing unfolded conjunctive query plans.

Friedman and Weld (1997) offer an efficient algorithm for minimizing a non-recursive query plan through the use of LCW statements. Their algorithm is based on pair-wise subsumption checks on conjunctive rules. Recursive rules correspond to infinite unions of conjunctive queries, and trying to prove subsumption through pair-wise conjunctive rule containment checks will not be decidable. The approach in Duschka (1997) also suffers from similar problems as it is based on the idea of conjunctive (un)foldings

of a query in terms of source relations (Qian, 1996). In the case of recursive queries or sources with binding restrictions, the number of such foldings is infinite. In contrast, our minimization algorithm is based on the notion of uniform containment for recursive datalog programs (Sagiv, 1988). This approach can check if sets of rules subsume a single rule. Thus it can minimize a much greater range of plans.

Execution optimization in *Emerac* involves ordering the calls to the sources so as to reduce the access and transfer costs. There are some similarities between this source call ordering and the join ordering in traditional databases. In contrast to traditional databases however, we cannot assume access to full statistics in the case of information gathering. For example, the execution ordering algorithm used in the Information Manifold (IM) system (Levy et al., 1996) generalizes the bound-is-easier style approach (first proposed as part of the Nail! system (Morris, 1988)) to work with Internet sources with multiple capability records (essentially, multiple feasible binding patterns per source), and to reduce tuple transfer costs by pushing selections to the sources. Our work can be seen as further extending the IM algorithm such that it uses coarse information about the result cardinality for each feasible binding pattern, as well as unselectable attribute limitations. Unlike the IM algorithm, we also explicitly consider optimizing both source access cost and tuple transfer cost. In contrast to IM which focuses on the tuple transfer costs, Yerneni and Li (1999) focus exclusively on minimizing access cost. The algorithm described in this paper may be seen as striking a middle ground between minimizing source access costs alone or minimizing tuple transfer costs alone. The experiments in Section 7.3 establish the importance of considering both types of costs. Finally, while use of heuristic algorithms for query optimization is one approach for dealing with lack of source statistics, another approach would be to *learn* the statistics. In Gruser and Zadorozhny (2000) describe an approach for online learning of response times for Web-accessible sources.

It should be noted that systems that address integration in the context of federated database systems do use more cost-based approaches. For example, the issue of join ordering in the context of heterogeneous distributed databases is considered in the DISCO (Raschid et al., 1998) and Garlic (Haas et al., 1997) projects. In contrast to our approach, both these projects assume availability of full statistics for the sources being integrated, and thus concentrate on cost-based optimization methods. For example, the Garlic optimizer assumes full knowledge of the statistics about the databases being integrated as well as their access and query support capabilities. The query processing capabilities are represented as a set of rules which are used by a Starburst-style optimizer (Lohman, 1989) to rewrite the mediator query. The statistics are used for cost-based optimization. In the DISCO (Raschid et al., 1998) approach, the optimizer assumes that the wrapper for each source provides a complete cost model for the source. The main difference between our approach and the Garlic approach is in terms of the granularity of knowledge available about the information sources being integrated. While the Garlic and DISCO approaches are well suited for federated database systems, where there is some level of central authority, they are less well suited for integration in the context of information gathering, where the sources are autonomous and decentralized, and are under no obligation to export source statistics. Our approach relies on more coarse-grained statistics and thus can get by without insisting on full knowledge of the source capabilities and statistics. In our current *Havasu* data integration project, we are pursuing a

complementary approach—that of *learning* the needed statistics, and then using them as the basis for cost-based query optimization (Nie and Kambhampati, 2001; Nie et al., 2002; Nie et al., 2001).

## 9. Conclusion

In this paper, we considered the query optimization problem for information gathering plans, and presented two novel techniques. The first technique makes use of LCW statements about information sources to prune unnecessary information sources from a plan. For this purpose, we have modified an existing method for minimizing datalog programs under uniform containment, so that it can minimize *recursive* information gathering plans with the help of source subsumption information. The second technique is a greedy algorithm for ordering source calls that respects source limitations, and takes both access costs and traffic costs into account, without requiring full source statistics. We have then discussed the status of a prototype implementation system based on these ideas called *Emerac*, and presented an evaluation of the effectiveness of the optimization strategies in the context of *Emerac*. We have related our work to other research efforts and argued that our approach is the first to consider end-to-end the issues of redundancy elimination and optimization in recursive information gathering plans.

We are currently exploring the utility of learning rudimentary source models by keeping track of time and solution quality statistics, and the utility of probabilistic characterizations of coverage and overlaps between sources (Nie et al., 2002; Nie et al., 2001). We are also working towards extending our current greedy plan generation methods, so as to search a larger space of feasible plans and to make the query optimization sensitive to both coverage and cost (Nie and Kambhampati, 2001).

## Notes

1. Notice that this is different from the capability provided by the existing search engines, which supply a list of pointers to the relevant sources, rather than return the requested data.
2. See (Lambrecht and Kambhampati, 1997) for a tutorial.
3. In other words, executing a maximally contained plan guarantees the return of every tuple satisfying the query that can be returned by executing any other query plan.
4. It would of course be interesting to try and “learn” the source statistics through judicious probing. See (Zhu and Larson, 1996) for a technique that does it in the context of multi-databases.
5. Things get a bit more complicated when there are variables in the body of the view that do not appear in the head. During inversion, every such variable is replaced with a new function term  $f_N(X_1, \dots, X_m)$ . The function symbols can then be eliminated by a flattening procedure, as there will be no recursion through them in the eventual plan, resulting in a datalog program in the end.
6. The next section contains an example illustrating this strategy.
7. Note that this step is safe because there is no recursion through global schema predicates. This step also removes any new predicates introduced through flattening of function symbols.
8. Even those sources that have large tables regulate their data output, by paginating the tuples and sending them in small quanta (e.g., first 10 tuples satisfying the query), which will avoid the network congestion if the users needed only a certain percentage of the tuples satisfying the query.
9. More elaborate languages for representing source access capabilities are proposed in (Garcia-Molina et al., 1999; Vassalos and Papakonstantinou, 1997).



10. Unless  $S_2$  accepts a list of possible values for X in a single call (Garcia-Molina et al., 1999).
11. Even those that have large tables regulate their data output, by paginating the tuples and sending them in small quanta (e.g., first 10 tuples satisfying the query), which will avoid the network congestion if the users needed only a certain percentage of the tuples satisfying the query.
12. The reason we keep an auxiliary variable  $V'$  to collect the bound variables within the second “for” loop and add them to  $V$  only outside of the inner loop is that we want any source calls being made in parallel within the same plan step to not depend on bindings that only become available in parallel branches of execution.
13. This example is inspired by the discussion in (Florescu et al., 1999).
14. IM’s bucket algorithm also uses a “generate-test” approach, generating candidate conjunctive query plans, and ensuring their soundness by testing if they are contained in the query. Duschka (1997) points out the disadvantages of this generate/test approach.

## References

- Raschid, L., Tomasic, A., and Valduriez, P. (1998). Scaling Access to Heterogeneous Data Sources with Disco. *IEEE TKDE*, 10(5).
- Abiteboul, S. and Duschka, O.M. (1998). Complexity of Answering Queries Using Materialized Views. In *Proceedings of the Twelfth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, PODS '98*.
- Adali, S., Candan, K.S., Papakonstantinou, Y., and Subrahmanian, V.S. (1996). Query Caching and Optimization in Distributed Mediator Systems. In *Proceedings of the 1996 ACM SIGMOD International Conference on Management of Data* (pp. 137–148).
- Adali, S., Candan, K.S., Papakonstantinou, Y., and Subrahmanian V.S. (1996). Query Caching and Optimization in Distributed Mediator Systems. In *Proceedings of the ACM Sigmod International Conference on Management of Data* (pp. 137–148).
- Chaudhuri, S. (1998). An Overview of Query Optimization in Relational Systems. In *Proceedings of the Sixteenth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, PODS 98* (pp. 34–43).
- Chaudhuri, S., Krishnamurthy, R., Potamianos, S., and Shim, K. (1995). Optimizing Queries with Materialized Views. In *Proceedings of the Eleventh International Conference on Data Engineering* (pp. 190–200). Los Alamitos, CA: IEEE Comput. Soc. Press.
- Chaudhuri, S. and Shim, K. (1993). Query Optimization in the Presence of Foreign Functions. In *Proc. 19th VLDB Conference*.
- Chawathe, S., Garcia-Molina, H., Hammer, J., Ireland, K., Papakonstantinou, Y., Ullman, J., and Widom, J. (1994). The TSIMMIS Project: Integration of Heterogeneous Information Sources. In *Proceedings of the 100th Anniversary Meeting* (pp. 7–18). Tokyo, Japan: Information Processing Society of Japan.
- Duschka, O.M. (1997). Query Optimization Using Local Completeness. In *Proceedings of the Fourteenth National Conference on Artificial Intelligence, AAAI-97* (pp. 249–255). Providence, RI.
- Duschka, O.M. and Genesereth, M.R. (1997). Answering Recursive Queries Using Views. In *Proceedings of the Sixteenth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, PODS '97* (pp. 109–116). Tucson, AZ.
- Duschka, O.M. and Levy, A.Y. (1997). Recursive Plans for Information Gathering. In *Proceedings of the Fifteenth International Joint Conference on Artificial Intelligence, IJCAI*. Nagoya, Japan.
- Duschka, O.M. (1997). Query Planning and Optimization in Information Integration. PhD thesis, Stanford University.
- Elmasri, R. and Navathe, S.B. (1994). *Fundamentals of Database Systems*, 2nd edn. The Benjamin/Cummings Publishing Company, Inc.
- Etzioni, O., Golden, K., and Weld, D. (1997). Sound and Efficient Closed-World Reasoning for Planning. *Artificial Intelligence*, 89(1/2), 113–148.
- Florescu, D., Koller, D., Levy, A.Y., and Pfeffer, A. (1997). Using Probabilistic Information in Data Integration. In *Proceedings of VLDB-97*.
- Florescu, D., Levy, A., Manolescu, I., and Suciu, D. (1999). Query Optimization in the Presence of Limited Access Patterns. In *Proc. SIGMOD Conference*.

- Florescu, D., Levy, A., and Mendelzon, A. (1998). Database Techniques for World-Wide Web: A Survey. *SIGMOD Record*.
- Friedman, M. and Weld, D.S. (1997). Efficiently Executing Information-Gathering Plans. In *Proceedings of the Fifteenth International Joint Conference on Artificial Intelligence, IJCAI*, Nagoya, Japan.
- Garcia-Molina, H., Labio, W., and Yerneni, R. (1999). Capability Sensitive Query Processing on Internet Sources. In *Proc. ICDE*.
- Garcia-Molina, H., Papanikolaou, Y., Quass, D., Rajaraman, A., Sagiv, Y., Ullman, J.D., Vassalos, V., and Widom, J. (1997). The TSIMMIS Approach to Mediation: Data Models and Languages. *Journal of Intelligent Information Systems*, 8(2), 117–132.
- Geddis, D.F., Genesereth, M.R., Keller, A.M., and Singh, N.P. (1995). Infomaster: A Virtual Information System. In *Intelligent Information Agents Workshop at CIKM '95*. Baltimore, MD.
- Haas, L.M., Kossmann, D., Wimmers, E.L., and Yang, J. (1997). Optimizing Queries Across Diverse Data Sources. In *Proc. VLDB*.
- Halevy, A. (2001). Answering Queries Using Views: A Survey. *VLDB Journal*.
- Hsu, C.-N. (1998). Initial Results on Wrapping Semistructured Web Pages with Finite-State Transducers and Contextual Rules. In *Proceedings of the AAAI Workshop on AI and Information Integration* (pp. 66–73).
- Gruser, L.R.J. and Zadorozhny, V. (2000). Learning Response Time for Websources Using Query Feedback and Application in Query Optimization. *The VLDB Journal*, 9.
- Kambhampati, S. and Gnanaprakasam, S. (1999). Optimizing Source-Call Ordering in Information Gathering Plans. In *Proc. IJCAI-99 Workshop on Intelligent Information Integration*.
- Kwok, C.T. and Weld, D.S. (1996). Planning to Gather Information. In *Proceedings of the AAAI Thirteenth National Conference on Artificial Intelligence*.
- Lambrecht, E. (1998). Optimizing Recursive Information Gathering Plans. Master's thesis, Arizona State University.
- Lambrecht, E. and Kambhampati, S. (1997). Planning for Information Gathering: A Tutorial Survey. Technical Report ASU CSE TR 97-017, Arizona State University. Available at: [rakaposhi.eas.asu.edu/ig-tr.ps](http://rakaposhi.eas.asu.edu/ig-tr.ps).
- Lambrecht, E. and Kambhampati, S. (1998). Optimizing Information Gathering Plans. In *Proc. AAAI-98 Workshop on Intelligent Information Integration*.
- Lambrecht, E., Kambhampati, S., and Gnanaprakasam, S. (1999). Optimizing Recursive Information Gathering Plans. In *Proc. IJCAI*.
- Levy, A.Y. (1996). Obtaining Complete Answers from Incomplete Databases. In *Proceedings of the 22nd International Conference on Very Large Databases* (pp. 402–412). Bombay, India.
- Levy, A.Y., Rajaraman, A., and Ordille, J.J. (1996). Querying Heterogeneous Information Sources Using Source Descriptions. In *Proceedings of the 22nd International Conference on Very Large Databases* (pp. 251–262). Bombay, India.
- Lohman, G., Haas, L.M., Freytag, J. and Pirahesh, H. (1989). Extensible Query Processing in Starburst. In *Proceedings of SIGMOD*.
- Morris, K.A. (1988). An Algorithm for Ordering Subgoals in Nail! In *Proceedings of PODS*.
- Nie, Z. and Kambhampati, S. (2001). Joint Optimization of Cost and Coverage of Query Plans in Data Integration. In *Proc. CIKM*.
- Nie, Z., Nambiar, U., Vaddi, S., and Kambhampati, S. (2002). Mining Coverage Statistics for Webservice Selection in a Mediator. In *Proc. CIKM*.
- Nie, Z., Kambhampati, S., Nambiar, U., and Vaddi, S. (2001). Mining Source Coverage Statistics for Data Integration. In *Proc. Web Information and Data Management (WIDM) Workshop*.
- Qian, X. (1996). Query Folding. In *Proceedings of the 12th International Conference on Data Engineering* (pp. 48–55), New Orleans, LA.
- Sagiv, Y. (1988). *Optimizing Datalog Programs*, ch. 17. M. Kaufmann Publishers.
- Tomasic, A., Raschid, L., and Valduriez, P. (1997). A Data Model and Query Processing Techniques for Scaling Access to Distributed Heterogeneous Databases in Disco. *IEEE Transactions on Computers, special issue on Distributed Computing Systems*.
- Ullman, J.D. (1989). *Principles of Database and Knowledgebase Systems*, vol. 2. Computer Science Press.

- Urhan, T. and Franklin, M. (1998). Cost-Based Query Scrambling for Initial Delays. In *Proceedings of SIGMOD*.
- Vassalos, V. and Papakonstantinou, Y. (1998). Using Knowledge of Tedundancy for Query Optimization in Mediators. In *Proceedings of the AAAI Workshop on AI and Information Integration* (pp. 29–35).
- Vassalos, V. and Papakonstantinou, Y. (1997). Describing and Using Query Capabilities of Heterogeneous Sources. In *Proc. VLDB*.
- Yerneni, R. and Li, C. (1999). Optimizing Large Join Queries in Mediation Systems. In *Proc. International Conference on Database Theory*.
- Zhu, Q. and Larson, P.-A. (1996). Developing Regression Cost Models for Multidatabase Systems. In *Proceedings of PDIS*.