

Predicate Migration: Optimizing Queries with Expensive Predicates

Joseph M. Hellerstein*

Michael Stonebraker

Computer Science Division, EECS Department
University of California, Berkeley, CA 94720
joey@cs.wisc.edu, mike@postgres.berkeley.edu

Abstract. The traditional focus of relational query optimization schemes has been on the choice of join methods and join orders. Restrictions have typically been handled in query optimizers by “predicate pushdown” rules, which apply restrictions in some random order before as many joins as possible. These rules work under the assumption that restriction is essentially a zero-time operation. However, today’s extensible and object-oriented database systems allow users to define time-consuming functions, which may be used in a query’s restriction and join predicates. Furthermore, SQL has long supported subquery predicates, which may be arbitrarily time-consuming to check. Thus restrictions should not be considered zero-time operations, and the model of query optimization must be enhanced.

In this paper we develop a theory for moving expensive predicates in a query plan so that the total cost of the plan — including the costs of both joins and restrictions — is minimal. We present an algorithm to implement the theory, as well as results of our implementation in POSTGRES. Our experience with the newly enhanced POSTGRES query optimizer demonstrates that correctly optimizing queries with expensive predicates often produces plans that are orders of magnitude faster than plans generated by a traditional query optimizer. The additional complexity of considering expensive predicates during optimization is found to be manageably small.

1 Introduction

Traditional relational database (RDBMS) literature on query optimization stresses the significance of choosing an efficient order of joins in a query plan. The placement of the other standard relational operators (selection and projection) in the plan has typically been handled by “pushdown” rules (see *e.g.*, [Ul189]), which state that restrictions and projections should be pushed down the query plan tree as far as possible. These rules place no importance on the ordering of projections and restrictions once they have been pushed below joins.

The rationale behind these pushdown rules is that the relational restriction and projection operators take essentially no time to carry out, and reduce subsequent join costs. In today’s

systems, however, restriction can no longer be considered to be a zero-time operation. Extensible database systems such as POSTGRES [SR86] and Starburst [HCL⁺90], as well as various Object-Oriented DBMSs (*e.g.*, [MS87], [WLH90], [D⁺90], [ONT92], etc.) allow users to implement predicate functions in a general-purpose programming language such as C or C++. These functions can be arbitrarily complex, potentially requiring access to large amounts of data, and extremely complex processing. Thus it is unwise to choose a random order of application for restrictions on such predicates, and it may not even be optimal to push them down a query plan tree. Therefore the traditional model of query optimization does not produce optimal plans for today’s queries, and as we shall see, the plans that traditional optimizers generate can be many orders of magnitude slower than a truly optimal plan.

To illustrate the significance of ordering restriction predicates, consider the following example:

Example 1.

```
/* Find all maps from week 17 showing more than
1% snow cover. Channel 4 contains images
from the frequency range that interests us. */
retrieve (maps.name)
where maps.week = 17 and maps.channel = 4
and coverage(maps.picture) > 1
```

In this example, the function `coverage` is a complex image analysis function that may take many thousands of instructions to compute. It should be quite clear that the query will run faster if the restrictions `maps.week = 17` and `maps.channel = 4` are applied before the restriction `coverage(maps.picture) > 1`, since doing so minimizes the number of calls to `coverage`.

While restriction ordering such as this is important, correctly ordering restrictions within a table-access is not sufficient to solve the general problem of where to place predicates in a query execution plan. Consider the following example:

Example 2.

```
/* Find all channel 4 maps from weeks starting
in June that show more than 1% snow cover.
Information about each week is kept in the
weeks table, requiring a join. */
retrieve (maps.name)
where maps.week = weeks.number
and weeks.month = "June"
and maps.channel = 4
and coverage(maps.picture) > 1
```

*Current address: Computer Sciences Department, University of Wisconsin, Madison, WI 53706. This material is based upon work supported under a National Science Foundation Graduate Fellowship.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

SIGMOD /5/93/Washington, DC, USA

© 1993 ACM 0-89791-592-5/93/0005/0267...\$1.50

Traditionally, a DBMS would execute this query by applying all the single-table restrictions in the *where* clause before performing the join of maps and weeks, since early restriction can lower the complexity of join processing. However in this example the cost of evaluating the expensive restriction predicate may outweigh the benefit gained by doing restriction before join. In other words, this may be a case where “predicate pushdown” is precisely the wrong technique. What is needed here is “predicate pullup”, namely postponing the restriction coverage (`maps.picture`) > 1 until after computing the join of maps and weeks.

In general it is not clear how joins and restrictions should be interleaved in an optimal execution plan, nor is it clear whether the migration of restrictions should have an effect on the join orders and methods used in the plan. This paper describes and proves the correctness of the *Predicate Migration Algorithm*, which produces a minimal-cost query plan for queries with expensive predicates. Predicate Migration modestly increases query optimization time: the additional cost factor is polynomial in the number of operators in a query plan. This compares favorably to the exponential join enumeration schemes used by most query optimizers, and is easily circumvented when optimizing queries without expensive predicates — if no expensive predicates are found while parsing the query, the techniques of this paper need not be invoked. For queries with expensive predicates, the gains in execution speed should offset the extra optimization time. We have implemented Predicate Migration in POSTGRES, and have found that with modest overhead in optimization time, the execution time of many practical queries can be reduced by orders of magnitude. This will be illustrated below.

1.1 Application to Existing Systems: SQL and Subqueries

It is important to note that expensive predicate functions do not exist only in next-generation research prototypes. Current relational languages, such as the industry standard, SQL [ISO91], have long supported expensive predicate functions in the guise of *subquery predicates*. A subquery predicate is one of the form “expression operator query”. Evaluating such a predicate requires executing an arbitrary query and scanning its result for matches — an operation that is arbitrarily expensive, depending on the complexity and size of the subquery. While some subquery predicates can be converted into joins (thereby becoming subject to traditional join-based optimization strategies) even sophisticated SQL rewrite systems, such as that of Starburst [PHH92], cannot convert all subqueries to joins. When one is forced to compute a subquery in order to evaluate a predicate, then the predicate should be treated as an expensive function. Thus the work presented in this paper is applicable to the majority of today’s production RDBMSs, which support SQL.

1.2 Related Work

Stonebraker first raised the issue of expensive predicate optimization in the context of the POSTGRES multi-level store [Sto91]. The questions posed by Stonebraker are directly addressed in this paper, although we vary slightly in the definition of cost metrics for expensive functions.

One of the main applications of the system described in [Sto91] is Project Sequoia 2000 [SD92], a University of California project that will manage terabytes of Geographic Infor-

mation System (GIS) data, to support global change researchers. It is expected that these researchers will be writing queries with expensive functions to analyze this data. A benchmark of such queries is presented in [SFGM93].

Ibaraki and Kameda [IK84], Krishnamurthy, Boral and Zaniolo [KBZ86], and Swami and Iyer [SI92] have developed and refined a query optimization scheme that is built on the notion of *rank* that we will use below. However, their scheme uses *rank* to reorder joins rather than restrictions. Their techniques do not consider the possibility of expensive restriction predicates, and only reorder nodes of a single path in a left-deep query plan tree, while the technique presented below optimizes all paths in an arbitrary tree. Furthermore, their schemes are a proposal for a completely new method for query optimization, while ours is an extension that can be applied to the plans of any query optimizer. It is possible to fuse the technique we develop in this paper with those of [IK84, KBZ86, SI92], but we do not focus on that issue here since their schemes are not widely in use.

The notion of expensive restrictions was considered in the context of the *LDL* logic programming system [CGK89]. Their solution was to model a restriction on relation *R* as a join between *R* and a virtual relation of infinite cardinality containing the entire logical predicate of the restriction. By modeling restrictions as joins, they were able to use a join-based query optimizer to order all predicates appropriately. Unfortunately, most traditional DBMS query optimizers have complexity that is exponential in the number of joins. Thus modelling restrictions as joins can make query optimization prohibitively expensive for a large set of queries, including queries on a single relation. The scheme presented here does not cause traditional optimizers to exhibit this exponential growth in optimization time.

Caching the return values of function calls will prove to be vital to the techniques presented in this paper. Jhingran [Jhi88] has explored a number of the issues involved in caching procedures for query optimization. Our model is slightly different, since our caching scheme is value-based, simply storing the results of a function on a set of argument values. Jhingran’s focus is on caching complex object attributes, and is therefore instance-based.

1.3 Structure of the Paper

The following section develops a model for measuring the cost and selectivity of a predicate, and describes the advantages of caching for expensive functions. Section 3 presents the Predicate Migration Algorithm, a scheme for optimally locating predicates in a given join plan. Section 4 details the results of our implementation experience in POSTGRES. Section 5 summarizes and provides directions for future research.

2 Background: Expenses and Caching

Query optimization schemes typically attempt to find a query plan of minimal estimated cost. To develop our optimizations, we must enhance the traditional model for analyzing query plan cost. This will involve some modifications of the usual metrics for the expense of relational operators, and will also require the introduction of *function caching* techniques. This preliminary discussion of our model will prove critical to the analysis below.

A relational query in a language such as SQL or Postquel [RS87] may have a *where* clause, which contains an

arbitrary Boolean expression over constants and the range variables of the query. We break such clauses into a maximal set of conjuncts, or “Boolean factors” [SAC⁺79], and refer to each Boolean factor as a distinct “predicate” to be satisfied by each result tuple of the query. When we use the term “predicate” below, we refer to a Boolean factor of the query’s *where* clause. A *join predicate* is one that refers to multiple tables, while a *restriction predicate* refers only to a single table.

Traditional query optimizers compute *selectivities* for both joins and restrictions. That is, for any predicate p (join or restriction) they estimate the value

$$selectivity(p) = \frac{card(output(p))}{card(input(p))}$$

and make the assumption that selectivities of different predicates are independent. Typically these estimations are based on default values and system statistics [SAC⁺79], although recent work suggests that accurate and inexpensive sampling techniques can be used [LNSS93, HOT88].

2.1 Cost of User-Defined Functions in POSTGRES

In an extensible system such as POSTGRES, arbitrary user-defined functions may be introduced into both restriction and join predicates. These functions may be written in a general programming language such as C, or in the database query language, e.g. SQL or Postquel. In this section we discuss programming language functions; we handle query language functions below.

Given that user-defined functions may be written in a general purpose language such as C, there is little hope for the database to correctly estimate the cost and selectivity of predicates containing these functions, at least not initially.¹ In this section we extend the POSTGRES function definition syntax to capture a function’s expense. Selectivity modeling for user-defined operators in POSTGRES has been described in [Mos90].

To introduce a function to POSTGRES, a user first writes the function in C and compiles it, and then issues Postquel’s `define function` command. To capture expense information, the `define function` command accepts a number of special flags, which are summarized in Table 1.

The cost of a predicate in POSTGRES is computed by adding up the costs for each expensive function in the expression. Given a POSTGRES predicate $p(a_1, \dots, a_n)$, the expense per tuple is recursively defined as:

$$e_p = \begin{cases} \sum_{i=1}^n e_{a_i} + percall_cpu(p) \\ \quad + perbyte_cpu(p) * (byte_pct(p)/100) \\ \quad * \sum_{i=1}^n bytes(a_i) + access_cost & \text{if } p \text{ is a function} \\ 0 & \text{if } p \text{ is a constant or tuple variable} \end{cases}$$

where e_{a_i} is the recursively computed expense of argument a_i , $bytes$ is the expected (return) size of the argument in bytes, and $access_cost$ is the cost of retrieving any data necessary to compute the function. This data may be stored anywhere in the various levels of the POSTGRES multi-level store, but unlike [Sto91] we do not require the user to define constants specific

¹After repeated applications of a function, one could collect performance statistics and use curve-fitting techniques to make estimates about the function’s behavior. Such techniques are beyond the scope of this paper.

to the different levels of the multi-level store. Instead, this can be computed by POSTGRES itself via system statistics, thus providing more accurate information about the distribution and caching of data across the storage levels.

2.2 Cost of SQL Subqueries and Other Query Language Functions

SQL allows a variety of subquery predicates of the form “expression operator query”. Such predicates require computation of an arbitrary SQL query for evaluation. Simple *uncorrelated* subqueries have no references to query blocks at higher nesting levels, while *correlated* subqueries refer to tuple variables in higher nesting levels.

In principle, the cost to check an uncorrelated subquery restriction is the cost e_m of materializing the subquery once, and the cost e_s of scanning the subquery once per tuple. However, we will need these cost estimates only to help us reorder operators in a query plan. Since the cost of initially materializing an uncorrelated subquery must be paid regardless of the subquery’s location in the plan, we ignore the overhead of the materialization cost, and consider an uncorrelated subquery’s cost per tuple to be e_s .

Correlated subqueries must be materialized for each value that is checked against the subquery predicate, and hence the per-tuple expense for correlated subqueries is e_m . We ignore e_s here since scanning can be done during each materialization, and does not represent a separate cost. Postquel functions in POSTGRES have costs that are equivalent to those of correlated subqueries in SQL: an arbitrary access plan is executed once per tuple of the relation being restricted by the Postquel function.

The cost estimates presented here for query language functions form a simple model and raise some issues in setting costs for subqueries. The cost of a subquery predicate may be lowered by transforming it to another subquery predicate [LDH⁺87], and by “early stop” techniques, which stop materializing or scanning a subquery as soon as the predicate can be resolved [Day87]. Incorporating such schemes is beyond the scope of this paper, but including them into the framework of the later sections merely requires more careful estimates of the subquery costs.

2.3 Join Expenses

In our subsequent analysis, we will be treating joins and restrictions uniformly in order to optimally balance their costs and benefits. In order to do this, we will need to measure the expense of a join per tuple of the join’s input, i.e. per tuple of the cartesian product of the relations being joined. This can be done for any join method whose costs are linear in the cardinalities of the input relations, including the most common algorithms: nested-loop join, hash join, and merge join. Note that sort-merge join is not linear in the cardinalities of the input relations. However, most systems, including POSTGRES, do not use sort-merge join, since in situations where merge join requires sorting of an input, either hash join or nested-loop join is almost always preferable to sort-merge.

A query may contain many join predicates over the same set of relations. In an execution plan for a query, some of these predicates are used in processing a join, and we call these *primary join predicates*. If a join has expensive primary join predicates, then the cost per tuple of a join should reflect the expensive function costs. That is, we add the expensive functions’ costs,

<i>flag name</i>	<i>description</i>
<i>percall_cpu</i>	execution time per invocation, regardless of the size of the arguments
<i>perbyte_cpu</i>	execution time per byte of arguments
<i>byte_pct</i>	percentage of argument bytes that the function will need to access

Table 1: Function Expense Parameters in POSTGRES

as described in Section 2.1, to the join costs per tuple.

Join predicates that are not applicable while processing the join are merely used to restrict its output, and we refer to these as *secondary join predicates*. Secondary join predicates are essentially no different from restriction predicates, and we treat them as such. These predicates may be reordered and even pulled up above higher join nodes, just like restriction predicates. Note, however, that a secondary join predicate must remain above its corresponding primary join. Otherwise the secondary join predicate would be impossible to evaluate.

2.4 Function Caching

The existence of expensive predicates not only motivates research into richer optimization schemes, it also suggests the need for DBMSs to cache the results of expensive predicate functions. In this paper, we assume that the system caches the return values of all functions for at least the duration of a query.² This lowers the cost of a function, since with some probability the function can be evaluated simply by checking the cache. Given the distribution of the data in a function's cache, and the distribution of the inputs to a function, one can derive a ratio of cache misses to cache lookups for the function. This ratio serves as the probability of a cache miss for a given tuple, and should be factored into the per-tuple cost for a function.

In addition to lowering function cost, caching will also allow us to pull expensive restrictions above joins without modifying the total cost of the restriction nodes in the plan. In general, a join may produce as many tuples as the product of the cardinalities of the inner and outer relations. However, it will produce no new *values* for attributes of the tuples; it will only recombine these attributes. If we move a restriction in a query plan from below a join to above it, we may dramatically increase the number of times we evaluate that restriction. However by caching expensive functions we will not increase the number of expensive function calls, only the number of cache lookups, which are quick to evaluate. This results from the fact that after pulling up the restriction, the same set of function calls on distinct arguments will be made. In most cases the primary join predicates will in fact *decrease* the number of distinct values passed into the function. Thus we see that with function caching, pulling restrictions above joins does not increase the number of function calls, and often will decrease that number.

Current SQL systems do not support arbitrary caching of the results of evaluating subquery predicates. To benefit from the techniques described in this paper, an SQL system must be enhanced to do this caching, at least for the duration of a query. It is interesting to note that in the original paper on optimizing

²As discussed in [Hel92], this cannot be done for some functions, e.g. functions that calculate the time of day. Such functions are rather unusual, though, since they result in ill-defined queries: the answer to such queries is dependent on the order in which tuples are scanned, something that is non-deterministic in relational-based systems.

Table	Tuple Size	#Tuples
maps	1 040 424	932
weeks	24	19
emp	32	10 000
dept	44	20

Table 2: Benchmark Database

SQL queries in System R [SAC⁺79], there is a description of a limited form of caching for correlated subqueries. System R saved the materialization of a correlated subquery after each evaluation, and if the subsequent tuple had the same values for the columns referenced in the subquery, then the predicate could be evaluated by scanning the saved materialization of the subquery. Thus System R would cache a single materialization of a subquery, but did not cache the result of the subquery predicate. That is, for a subquery of the form "expression operator query", System R cached the result of "query", but not "expression operator query".

2.5 Environment for Performance Measurements

It is not uncommon for queries to take hours or even days to complete. The techniques of this paper can improve performance by several orders of magnitude — in many cases converting an over-night query to an interactive one. We will be demonstrating this fact during the course of the discussion by measuring the performance effect of our optimizations on various queries. In this section we present the environment used for these measurements.

We focus on a complex query workload (involving subqueries, expensive user-defined functions, etc), rather than a transaction workload, where queries are relatively simple. There is no accepted standard complex query workload, although several have been proposed ([SFGM93, TOB89, O'N89], etc.) To measure the performance effect of Predicate Migration, we have constructed our own benchmark database, based on a combined GIS and business application. Each tuple in maps contains a reference to a POSTGRES large object [Ols92], which is a map picture taken by a satellite. These map pictures were taken weekly, and the maps table contains a foreign key to the weeks table, which stores information about the week in which each picture was taken. The familiar emp and dept tables store information about employees and their departments. Some physical characteristics of the database are shown in Table 2.

Our performance measurements were done in a development version of POSTGRES, similar to the publicly available version 4.1 (which itself contains the Predicate Migration optimizations). POSTGRES was run on a DECStation 5000/200 workstation, equipped with 24Mb of main memory and two 300Mb DEC RZ55 disks, running the Ultrix 4.2a operating system. We measured the elapsed time (total time taken by system), and CPU time (the time for which CPU is busy) of optimizing

and executing each example query, both with and without Predicate Migration. These numbers are presented in the examples which appear throughout the rest of the paper.

3 Min-Cost Plans for Queries With Expensive Predicates

At first glance, the task of correctly optimizing queries with expensive predicates appears exceedingly complex. Traditional query optimizers already search a plan space that is exponential in the number of relations being joined; multiplying this plan space by the number of permutations of the restriction predicates could make traditional plan enumeration techniques prohibitively expensive. In this section we prove the reassuring results that:

1. Given a particular query plan, its restriction predicates can be optimally interleaved based on a simple sorting algorithm.
2. As a result of the previous point, we need merely enhance the traditional join plan enumeration with techniques to interleave the predicates of each plan appropriately. This interleaving takes time that is polynomial in the number of operators in a plan.

3.1 Optimal Predicate Ordering in Table Accesses

We begin our discussion by focusing on the simple case of queries over a single table. Such queries may have an arbitrary number of restriction predicates, each of which may be a complicated Boolean function over the table's range variables, possibly containing expensive subqueries or user-defined functions. Our task is to order these predicates in such a way as to minimize the expense of applying them to the tuples of the relation being scanned.

If the access path for the query is an index scan, then all the predicates that match the index and can be applied during the scan are applied first. This is because such predicates are essentially of zero cost: they are not actually evaluated, rather the indices are used to retrieve only those tuples which qualify. Note that it is possible to index tables on function values as well as on table attributes [MS86, LS88]. If a scan is done on such a "function" index, then predicates over the function may be applied during the scan, and are considered to have zero cost, regardless of the function's expense.

We will represent each of the subsequent non-index predicates as p_1, \dots, p_n , where the subscript of the predicate represents its place in the order in which the predicates are applied to each tuple of the base table. We represent the expense of a predicate p_i as e_{p_i} , and its selectivity as s_{p_i} . Assuming the independence of distinct predicates, the cost of applying all the non-index predicates to the output of a scan containing t tuples is

$$e_1 = e_{p_1} t + s_{p_1} e_{p_2} t + \dots + s_{p_1} s_{p_2} \dots s_{p_{n-1}} e_{p_n} t.$$

The following lemma demonstrates that this cost can be minimized by a simple sort on the predicates.

Lemma 1 *The cost of applying expensive restriction predicates to a set of tuples is minimized by applying the predicates in ascending order of the metric*

$$\text{rank} = \frac{\text{selectivity} - 1}{\text{cost-per-tuple}}$$

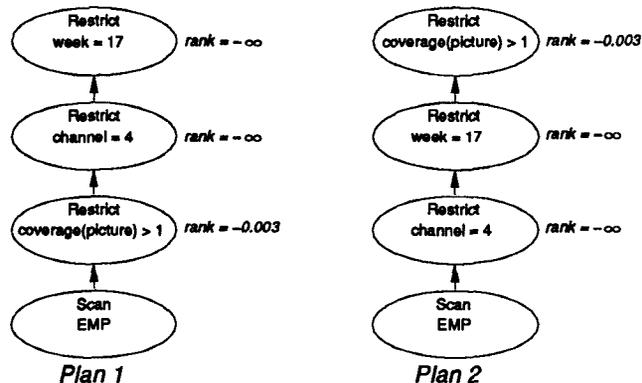


Figure 1: Two Execution Plans for Example 1

Proof. This results directly from work done by W. E. Smith [Smi56] on job scheduling. It has been reviewed in a database context in [Han77], [IK84], [KBZ86], and [Hel92]. Intuitively, the above ordering gives priority to the execution of restrictions with low selectivity and low cost. This reduces the number of tuples that will have to be processed by more expensive predicates. ■

Thus we see that for single table queries, predicates can be optimally ordered by simply sorting them by their rank. Swapping the position of predicates with equal rank has no effect on the cost of the sequence.

To see the effects of reordering restrictions, we return to Example 1 from the introduction. We ran the query in POSTGRES without the rank-sort optimization, generating Plan 1 of Figure 1, and with the rank-sort optimization, generating Plan 2 of Figure 1. As we expect from Lemma 1, the first plan has higher cost than the second plan, since the second is correctly ordered by rank. The optimization and execution times were measured for both runs, as illustrated in Table 3. We see that correctly ordering the restrictions can improve query execution time by orders of magnitude.

3.2 Predicate Migration: Moving Restrictions Among Joins

In the previous section, we established an optimal ordering for restrictions. In this section, we explore the issue of ordering restrictions among joins. Since we will eventually be applying our optimization to each plan produced by a typical join-enumerating query optimizer, our model here is that we are given a fixed join plan, and want to minimize the plan's cost under the constraint that we may not change the order of the joins. This section develops a polynomial-time algorithm to optimally place restrictions and secondary join predicates in a join plan.

3.2.1 Definitions

The thrust of this section is to handle join predicates in our ordering scheme in the same way that we handle restriction predicates: by having them participate in an ordering based on rank. However, since joins are binary operators, we must generalize our model for single-table queries to handle both restrictions and joins. We will refer to our generalized model as a *global* model, since it will encompass the costs of all inputs to a query, not just the cost of a single input to a single node.

Execution Plan	Optimization Time		Execution time	
	CPU	Elapsed	CPU	Elapsed
Plan 1	0.12 sec	0.24 sec	20 min 34.36 sec	20 min 37.69 sec
Plan 2 (ordered by rank)	0.12 sec	0.24 sec	0 min 2.66 sec	0 min 3.26 sec

Table 3: Performance of Example 1

Definition 1 A plan tree is a tree whose leaves are scan nodes, and whose internal nodes are either joins or restrictions. Tuples are produced by scan nodes and flow upwards along the edges of the plan tree.

Some optimization schemes constrain plan trees to be within a particular class, such as the *left-deep* trees, which have scans as the right child of every join. Our methods will not require this limitation. We do not, however, consider non-tree queries, *i.e.* queries with common subexpressions or recursion.

Definition 2 A stream in a plan tree is a path from a leaf node to the root.

Figure 2 below shows a tree with one of its plan streams outlined. Within the framework of a single stream, a join node is simply another predicate; although it has a different number of inputs than a restriction, it can be treated in an identical fashion. We do this by considering each predicate in the tree — restriction or join — as an operator on the *entire* input set to the query. That is, we consider the input to the query to be the cartesian product of all relations referenced in the query, and we model each node as an operator on that cartesian product. By modeling each predicate in this global fashion, we can naturally compare restrictions and joins in different streams. However, to do this correctly, we must modify our notion of the per-tuple cost of a predicate:

Definition 3 Given a query over relations a_1, \dots, a_n , the global cost of a predicate p over relations a_1, \dots, a_k is defined as:

$$\text{global-cost}(p) = \frac{\text{cost-per-tuple}(p)}{\text{card}(a_{k+1}) \cdots \text{card}(a_n)}$$

where *cost-per-tuple* is the cost attribute of the predicate, as described in Section 2.

That is, to define the cost of a predicate over the entire input to the query, we must divide out the cardinalities of those tables that do not affect the predicate. As an illustration, consider the case where p is a single-table restriction over relation a_1 . If we push p down to directly follow the table-access of a_1 , the cost of applying p to that table is $\text{cost-per-tuple}(p)\text{card}(a_1)$. But in our new global model, we consider the input to each node to be the cartesian product of a_1, \dots, a_n . However, note that the cost of applying p in both the global and single-table models is the same, *i.e.*,

$$\begin{aligned} &\text{global-cost}(p)\text{card}(a_1 \times \dots \times a_n) \\ &= \text{cost-per-tuple}(p)\text{card}(a_1). \end{aligned}$$

Recall that because of function caching, even if we pull p up to the top of the tree, its cost should not reflect the cardinalities of relations a_2, \dots, a_n . Thus our global model does not change the cost analysis of a plan. It merely provides a framework in which we can treat all predicates uniformly.

The selectivity of a predicate is independent of the predicate's location in the plan tree. This follows from the fact that

$\text{card}(a_1 \times a_2) = \text{card}(a_1)\text{card}(a_2)$. Thus the global rank of a predicate is easily derived:

Definition 4 The global rank of a predicate p is defined as

$$\text{rank} = \frac{\text{selectivity}(p) - 1}{\text{global-cost}(p)}$$

Note that the global cost of a predicate in a single-table query is the same as its user-defined *cost-per-tuple*, and hence the global rank of a node in a single-table query is the same as its *rank* as defined previously. Thus we see that the global model is a generalization of the one presented for single-table queries. In the subsequent discussion, when we refer to the *rank* of a predicate, we mean its global rank.

In later analysis it will prove useful to assume that all nodes have distinct *ranks*. To make this assumption, we must prove that swapping nodes of equal *rank* has no effect on the cost of a plan.

Lemma 2 Swapping the positions of two equi-rank nodes has no effect on the cost of a plan tree.

Proof. Note that swapping two nodes in a plan tree only affects the costs of those two nodes. Consider two nodes p and q of equal rank, operating on input of cardinality t . If we order p before q , their joint cost is $e_1 = te_p + ts_p e_q$. Swapping them results in the cost $e_2 = te_q + ts_q e_p$. Since their ranks are equal, it is a matter of simple algebra to demonstrate that $e_1 = e_2$, and hence the cost of a plan tree is independent of the order of equi-rank nodes. ■

Knowing this, we could achieve a unique ordering on *rank* by assigning unique ID numbers to each node in the tree and ordering nodes on the pair (*rank*, ID). Rather than introduce the ID numbers, however, we will make the simplifying assumption that *ranks* are unique.

In moving restrictions around a plan tree, it is possible to push a restriction down to a location in which the restriction cannot be evaluated. This notion is captured in the following definition:

Definition 5 A plan stream is semantically incorrect if some predicate in the stream refers to attributes that do not appear in the predicate's input.

Streams can be rendered semantically incorrect by pushing a secondary join predicate below its corresponding primary join, or by pulling a restriction from one input stream above a join, and then pushing it down below the join into the other input stream. We will need to be careful later on to rule out these possibilities.

In our subsequent analysis, we will need to identify plan trees that are equivalent except for the location of their restrictions and secondary join predicates. We formalize this as follows:

Definition 6 Two plan trees T and T' are join-order equivalent if they contain the same set of nodes, and there is a one-to-one mapping g from the streams of T to the streams of T' such that for any stream s of T , s and $g(s)$ contain the same join nodes in the same order.

3.2.2 The Predicate Migration Algorithm: Optimizing a Plan Tree By Optimizing its Streams

Our approach in optimizing a plan tree will be to treat each of its streams individually, and sort the nodes in the streams based on their *rank*. Unfortunately, sorting a stream in a general plan tree is not as simple as sorting the restrictions in a table access, since the order of nodes in a stream is constrained in two ways. First, we are not allowed to reorder join nodes, since join-order enumeration is handled separately from Predicate Migration. Second, we must ensure that each stream remains semantically correct. In some situations, these constraints may preclude the option of simply ordering a stream by ascending *rank*, since a predicate p_1 may be constrained to precede a predicate p_2 , even though $rank(p_1) > rank(p_2)$. In such situations, we will need to find the optimal ordering of predicates in the stream subject to the precedence constraints.

Monma and Sidney [MS79] have shown that finding the optimal ordering under a large class of precedence constraints can be done fairly simply. Their analysis is based on two key results:

1. A stream can be broken down into *modules*, where a module is defined as a set of nodes that have the same constraint relationship with all nodes outside the module. An optimal ordering for a module forms a subset of an optimal ordering for the entire stream.
2. For two predicates p_1, p_2 such that p_1 is constrained to precede p_2 and $rank(p_1) > rank(p_2)$, an optimal ordering will have p_1 directly preceding p_2 , with no other unconstrained predicates in between.

Monma and Sidney use these principles to develop the *Series-Parallel Algorithm Using Parallel Chains*, an $O(n \log n)$ algorithm for optimizing streams under a large class of constraints. The algorithm repeatedly isolates modules in a stream, optimizing each module individually, and using the resulting orders for modules to find a total order for the stream. Since the constraints which can appear in a query plan stream are subsumed by those considered by Monma and Sidney, we use their algorithm as a subroutine in our optimization algorithm.

“Predicate pushdown” is traditionally considered a good heuristic, and most systems construct plan trees with restriction and secondary join predicates pushed down as far as possible. Thus our algorithm was designed to work on plan trees with predicates already pushed down. For completeness, we include the pushdown step in the algorithm, although it would be unnecessary in most RDBMS implementations.

Predicate Migration Algorithm: *To optimize a plan tree, we push all predicates down as far as possible, and then repeatedly apply the Series-Parallel Algorithm Using Parallel Chains [MS79] to each stream in the tree, until no more progress can be made.*

Upon termination, the Predicate Migration Algorithm produces a tree in which each stream is *well-ordered* (i.e. optimally ordered subject to the precedence constraints). We proceed to prove that the Predicate Migration Algorithm is guaranteed to terminate in polynomial time, and we also prove that the resulting tree of well-ordered streams represents the optimal choice of predicate locations for the given plan tree.

Theorem 1 *Given any plan tree as input, the Predicate Migration Algorithm is guaranteed to terminate in polynomial time,*

producing a join-order equivalent tree in which each stream is semantically correct and well-ordered.

Proof. The proof, which appears in [Hel92], has been deleted due to space constraints. It develops a conservative upper bound of $O(n^4 \log n)$ for the algorithm’s running time, where n is the number of nodes in the plan tree. ■

Theorem 1 demonstrates that the Predicate Migration Algorithm terminates, and [MS79] assures us that each stream in the resulting tree is well-ordered. This is not sufficient, however, to establish the optimality of the algorithm’s output — we must also prove that the resulting tree of well-ordered streams is a minimal-cost tree. This is guaranteed by the following:

Theorem 2 *For every plan tree T there is a unique join-order equivalent plan tree T' with only semantically correct, well-ordered streams. Moreover, T' is a minimal cost tree that is join-order equivalent to T and semantically correct.*

Proof. Deleted due to space constraints. It appears in full in [Hel92]. ■

Theorems 1 and 2 demonstrate that the Predicate Migration Algorithm produces our desired minimal-cost interleaving of predicates. As a simple illustration of the efficacy of Predicate Migration, we go back to Example 2 from the introduction. Figure 2 illustrates plans generated for this query by POSTGRES running both with and without Predicate Migration. The performance measurements for the two plans appear in Table 4.

4 Implementation and Further Measurement

The Predicate Migration Algorithm, as well as pruning optimizations described in [Hel92], were implemented in the POSTGRES next-generation DBMS, which has an optimizer based on that of System R. The addition of Predicate Migration to POSTGRES was fairly straightforward, requiring slightly more than one person-month of programming. The implementation consists of two files containing a total of about 2000 lines, or 600 statements, of C language code. It should thus be clear that enhancing an optimizer to support Predicate Migration is a fairly manageable task.

Given the ease of implementation, and the potential benefits for both standard SQL and extensible query languages, it is our belief that Predicate Migration is a worthwhile addition to any DBMS. To further motivate this, we present two more examples, which model SQL queries that would be natural to run in most commercial DBMSs. We simulate an SQL correlated subquery with a Postquel query language function, since POSTGRES does not support SQL. As noted above, SQL’s correlated subqueries and Postquel’s query language functions require the same processing to evaluate, namely the execution of a subplan per value. The only major distinction between our Postquel queries and an SQL system is that Postquel may return a different number of duplicate tuples than SQL, since Postquel assigns no semantics to the duplicates in a query’s output. In our benchmark database the example queries return no tuples, and hence this issue does not affect the performance of our examples.

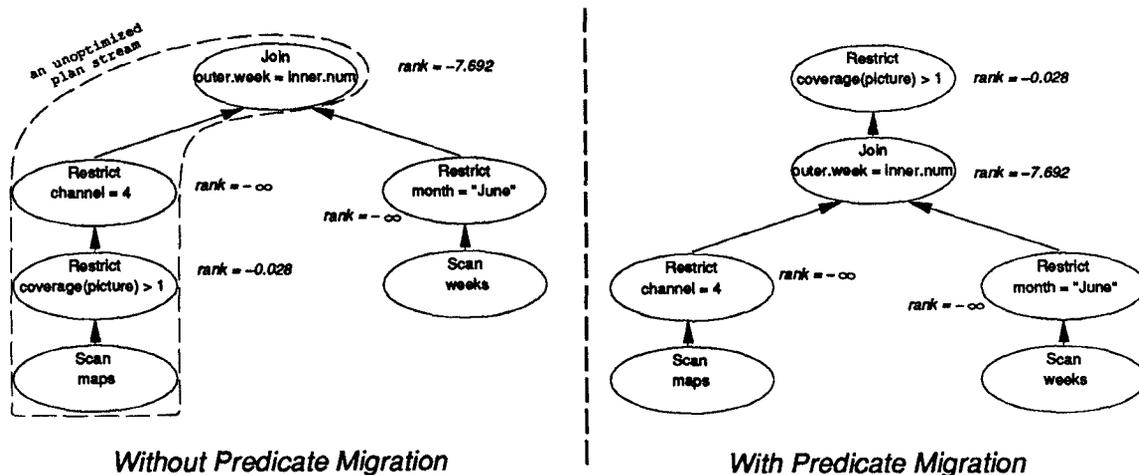


Figure 2: Plans For Example 2, With and Without Predicate Migration

Execution Plan	Optimization Time		Execution time	
	CPU	Elapsed	CPU	Elapsed
Without Predicate Migration	0.29 sec	0.30 sec	20 min 29.79 sec	21 min 12.98 sec
With Predicate Migration	0.36 sec	0.57 sec	0 min 3.46 sec	0 min 6.75 sec

Table 4: Performance of Plans for Example 2

Example 3. This query finds all technical departments with either low budgets or an employee over the age of 65. In SQL, the query is:

```
SELECT name FROM dept d1
WHERE d1.category = 'tech'
AND (d1.budget < 1000
OR EXISTS (SELECT 1 FROM emp
WHERE emp.dno = d1.dno
AND emp.age > 65));
```

Since the existential subquery is nested within an OR, the subquery cannot be converted to a join [PHH92]. To simulate this query in Postquel, we define a function `seniors`, which takes one argument (`$1`) of type integer, and executes the Postquel query:

```
retrieve (x = "t")
where emp.dno = $1 and emp.age > 65
```

Given this function, the SQL query is simulated by the following Postquel query:

```
retrieve (dept.name)
where dept.category = "tech"
and (dept.budget < 1000
or seniors(dept.dno))
```

Predicate Migration ensures that the expensive OR clause containing `seniors` is applied after the restriction `dept.category = "tech"`.³ As shown in Table 5, Predicate Migration speeds up execution time by orders of magnitude, while affecting optimization time only marginally.

³As an additional optimization, POSTGRES orders the operands of OR by *rank*, and quits evaluating the OR expression as soon as any operand evaluates to true. This issue was left out of the discussion previously in order to simplify matters. It is a straightforward extension to the techniques presented here.

Example 4. Our final example uses a subquery and a join to find the managers of the departments found in the previous example. The SQL version of the query is:

```
SELECT dept.name, mgr.name
FROM dept d1, emp mgr
WHERE d1.category = 'tech'
AND d1.dno = mgr.dno
AND (d1.budget < 1000
OR EXISTS (SELECT 1 FROM emp e1
WHERE e1.dno = d1.dno
AND e1.age > 65));
```

Since this uses the same subquery as the previous example, the equivalent Postquel query can reuse the function `seniors`:

```
retrieve(dept.name, mgr.name) from mgr in emp
where dept.category = "tech"
and dept.dno = mgr.dno
and (dept.budget < 1000
or seniors(dept.dno))
```

Predicate Migration in this query pulls the expensive OR clause above the join of `dept` and `emp`, resulting in the dramatic execution speedup shown in Table 6. Once again, the increase in optimization time is comfortably low.

These examples demonstrate that even for short queries in standard SQL, the techniques presented in this paper can improve execution time by orders of magnitude.

5 Conclusions and Future Work

In this paper we highlight the fact that database query optimization has up until now ignored the costs associated with restriction. We present a framework for measuring these costs, and we argue the necessity of caching expensive functions in a DBMS. We develop the Predicate Migration Algorithm, which

Execution Plan	Optimization Time		Execution time	
	CPU	Elapsed	CPU	Elapsed
Unoptimized Plan	0.34 sec	0.75 sec	2 min 25.61 sec	2 min 26.32 sec
Optimized Plan	0.34 sec	0.88 sec	0 min 0.06 sec	0 min 0.39 sec

Table 5: Performance of Plans for Example 3

Execution Plan	Optimization Time		Execution time	
	CPU	Elapsed	CPU	Elapsed
Unoptimized Plan	0.13 sec	0.42 sec	2 min 24.51 sec	2 min 25.69 sec
Optimized Plan	0.16 sec	0.52 sec	0 min 0.06 sec	0 min 0.39 sec

Table 6: Performance of Plans for Example 4

is proven to transform query plans in a way that optimally interleaves restriction and join predicates. This was implemented in POSTGRES, and measurements show that Predicate Migration is a low-overhead optimization that can produce query plans that run orders of magnitude faster than those produced by systems without Predicate Migration. This work can be applied not only to advanced research DBMSs such as POSTGRES, but also to any DBMS that supports SQL. There are not many additions to current DBMSs that can produce dramatic performance gains with modest implementation cost. Predicate Migration is one such addition.

The optimization schemes in this paper are useful for run-time re-optimization. That is, if a query is optimized and the resulting plan is stored for a period of time, the statistics that shaped the choice of the optimal plan may have changed. Predicate Migration can be re-applied to the stored plan at runtime with little difficulty. This may not produce an optimal plan, since the join orders and methods may no longer be optimal. But it will optimize the stored plan itself, without incurring the exponential costs of completely re-optimizing the query. This could be particularly beneficial for queries with subqueries, since the costs of the subqueries are likely to change over time.

This paper represents only an initial effort at optimizing queries with expensive predicates, and there is substantial work remaining to be done in this area. The first and most important question is whether the assumptions of this paper can be relaxed without making query optimization time unreasonably slow. The two basic assumptions in the paper are (1) that function caching is implemented, and (2) that join costs are linear in the size of the inputs. Without either of these assumptions, there are no obvious directions to pursue a polynomial-time algorithm for Predicate Migration. If one does not have function caching, then our cost model no longer applies, since a restriction function will be called once for every *tuple* that flows through its predicate, rather than once per *value* of the attributes on which it is defined. If one does not assume linear join costs, then the algorithm of [MS79] no longer applies. It would be interesting to discover whether the problem of Predicate Migration can be solved in polynomial time in general, or whether the assumptions made here are in fact crucial to a polynomial-time solution.

The implementation of function caching in POSTGRES has not been completed. Once that is accomplished, we will be able to perform more complex experiments than the ones presented here, which were carefully tailored to produce no duplicate function calls after pullup. A more comprehensive performance

study could develop a test suite of queries with expensive functions, and compare the performance of the Predicate Migration Algorithm against more naive predicate pullup heuristics.

It would be interesting to attempt to extend this work to handle queries with common subexpressions and recursion. Pulling up a restriction from a common subexpression may require duplication of the restriction, while in cyclic (*i.e.* recursive) query plans it is not even clear what “predicate pullup” means, since “up” is not well defined.

Finally, our cost analyses for user-defined functions could be dramatically improved by techniques to more correctly assess the expected running-time of a function on a given set of inputs. Particularly, the POSTGRES define function command includes an implicit assumption that users’ functions will have complexity that is linear in the size of their data objects. This simplifying assumption was made to ease implementation, but it is certainly possible to add curve-fitting algorithms to better model a function’s running time and complexity.

6 Acknowledgments

Wei Hong was an invaluable resource, providing extensive and regular feedback on this work. Jeff Naughton’s encouragement, patience and support helped to bring this project to completion. Thanks to Mike Olson, Mark Sullivan, and Kurt Brown for their comments on earlier drafts of this paper. This work could not have been completed without the assistance, suggestions, and friendly support of the entire POSTGRES research group.

References

- [CGK89] Danette Chimenti, Ruben Gamboa, and Ravi Krishnamurthy. Towards an Open Architecture for LDL. In *Proc. 15th International Conference on Very Large Data Bases*, Amsterdam, August 1989.
- [D⁺90] O. Deux et al. The Story of *O₂*. *IEEE Transactions on Knowledge and Data Engineering*, 2(1), March 1990.
- [Day87] Umeshwar Dayal. Of Nests and Trees: A Unified Approach to Processing Queries that Contain Nested Subqueries, Aggregates, and Quantifiers. In *Proc. VLDB 87 [Pro87]*, pages 197–208.
- [Han77] Michael Z. Hanani. An Optimal Evaluation of Boolean Expressions in an Online Query System. *Communications of the ACM*, 20(5):344–347, may 1977.
- [HCL⁺90] L.M. Haas, W. Chang, G.M. Lohman, J. McPherson, P.F. Wilms, G. Lapis, B. Lindsay, H. Pirahesh, M. Carey, and E. Shekita. Starburst Mid-Flight: As the Dust Clears. *IEEE*

- Transactions on Knowledge and Data Engineering*, pages 143–160, March 1990.
- [Hel92] Joseph M. Hellerstein. Predicate Migration: Optimizing Queries With Expensive Predicates. Technical Report Sequoia 2000 92/13, University of California, Berkeley, December 1992.
- [HOT88] Wen-Chi Hou, Gultekin Ozsoyoglu, and Baldeao K. Taneja. Statistical Estimators for Relational Algebra Expressions. In *Proc. 7th ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems*, pages 276–287, Austin, March 1988.
- [IK84] Toshihide Ibaraki and Tiko Kameda. Optimal Nesting for Computing N-relational Joins. *ACM Transactions on Database Systems*, 9(3):482–502, October 1984.
- [ISO91] ISO_ANSI. Database Language SQL ISO/IEC 9075:1992, 1991.
- [Jhi88] Anant Jhingran. A Performance Study of Query Optimization Algorithms on a Database System Supporting Procedures. In *Proc. VLDB 88* [Pro88].
- [KBZ86] Ravi Krishnamurthy, Haran Boral, and Carlo Zaniolo. Optimization of Nonrecursive Queries. In *Proc. 12th International Conference on Very Large Data Bases*, pages 128–137, Kyoto, August 1986.
- [LDH⁺87] Guy M. Lohman, Dean Daniels, Laura M. Haas, Ruth Kistler, and Patricia G. Selinger. Optimization of Nested Queries in a Distributed Relational Database. In *Proc. VLDB 87* [Pro87].
- [LNSS93] Richard J. Lipton, Jeffrey F. Naughton, Donovan A. Schneider, and S. Seshadri. Efficient Sampling Strategies for Relational Database Operations. To appear in *Theoretical Computer Science*, 1993.
- [LS88] C. Lynch and M. Stonebraker. Extended User-Defined Indexing with Application to Textual Databases. In *Proc. VLDB 88* [Pro88].
- [Mos90] Claire Mosher (ed.). The POSTGRES Reference Manual, Volume 2. Technical Report M90/53, Electronics Research Laboratory, University of California, Berkeley, July 1990.
- [MS79] C. L. Monma and J.B. Sidney. Sequencing with Series-Parallel Precedence Constraints. *Mathematics of Operations Research*, 4:215–224, 1979.
- [MS86] D. Maier and J. Stein. Indexing in an Object-Oriented DBMS. In Klaus R. Dittrich and Umeshwar Dayal, editors, *Proc. Workshop on Object-Oriented Database Systems*, Asilomar, September 1986.
- [MS87] D. Maier and J. Stein. Development and Implementation of an Object-Oriented DBMS. In Bruce Shriver and Peter Wegner, editors, *Research Directions in Object-Oriented Programming*. MIT Press, 1987.
- [Ols92] Michael A. Olson. Extending the POSTGRES Database System to Manage Tertiary Storage. Master's thesis, University of California, Berkeley, May 1992.
- [O'N89] P. O'Neil. Revisiting DBMS Benchmarks. *Datamation*, pages 47–54, September 15, 1989.
- [ONT92] ONTOS, Inc. *ONTOS Object SQL Guide*, February 1992. For the ONTOS DB database, Release 2.2.
- [PHH92] Hamid Pirahesh, Joseph M. Hellerstein, and Waqar Hasan. Extensible/Rule-Based Query Rewrite Optimization in Starburst. In *Proc. ACM-SIGMOD International Conference on Management of Data*, pages 39–48, San Diego, June 1992.
- [Pro87] *Proc. 13th International Conference on Very Large Data Bases*, Brighton, September 1987.
- [Pro88] *Proc. 14th International Conference on Very Large Data Bases*, Los Angeles, August-September 1988.
- [RS87] L.A. Rowe and M.R. Stonebraker. The POSTGRES Data Model. In *Proc. VLDB 87* [Pro87], pages 83–96.
- [SAC⁺79] Patricia G. Selinger, M. Astrahan, D. Chamberlin, Raymond Lorie, and T. Price. Access Path Selection in a Relational Database Management System. In *Proc. ACM-SIGMOD International Conference on Management of Data*, Boston, June 1979.
- [SD92] Michael Stonebraker and Jeff Dozier. Sequoia 2000: Large Capacity Object Servers to Support Global Change Research. Technical Report Sequoia 2000 91/1, University of California, Berkeley, March 1992.
- [SFGM93] Michael Stonebraker, James Frew, Kenn Gardels, and Jeff Meredith. The Sequoia 2000 Storage Benchmark. In *Proc. ACM-SIGMOD International Conference on Management of Data*, Washington, D.C., May 1993.
- [SI92] Arun Swami and Balakrishna R. Iyer. A Polynomial Time Algorithm for Optimizing Join Queries. Research Report RJ 8812, IBM Almaden Research Center, June 1992.
- [Smi56] W. E. Smith. Various Optimizers For Single-Stage Production. *Naval Res. Logist. Quart.*, 3:59–66, 1956.
- [SR86] M.R. Stonebraker and L.A. Rowe. The Design of POSTGRES. In *Proc. ACM-SIGMOD International Conference on Management of Data*, Washington, D.C., May 1986.
- [Sto91] Michael Stonebraker. Managing Persistent Objects in a Multi-Level Store. In *Proc. ACM-SIGMOD International Conference on Management of Data*, pages 2–11, Denver, June 1991.
- [TOB89] C. Turbyfill, C. Orji, and Dina Bitton. AS3AP - A Comparative Relational Database Benchmark. In *Proc. IEEE Comcon Spring '89*, February 1989.
- [Ull89] Jeffrey D. Ullman. *Principles of Database and Knowledge-Base Systems*, volume 2. Computer Science Press, 1989.
- [WLH90] K. Wilkinson, P. Lyngbaek, and W. Hasan. The Iris Architecture and Implementation. *IEEE Transactions on Knowledge and Data Engineering*, 2(1), March 1990.