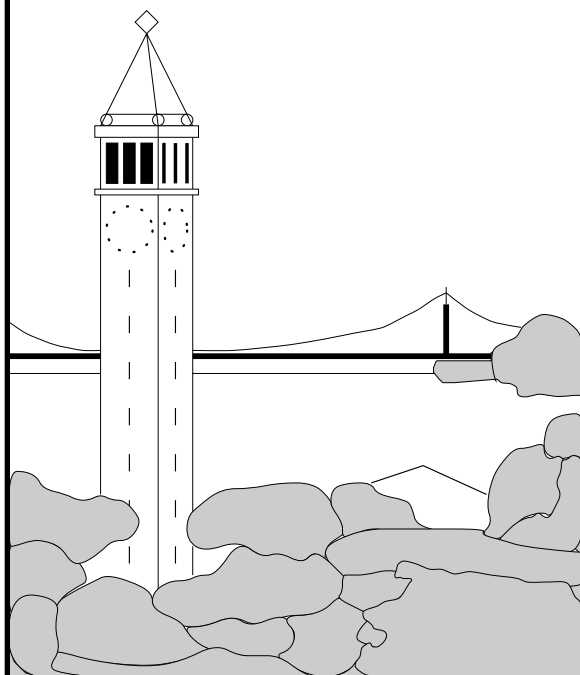


Using State Modules for Adaptive Query Processing

Vijayshankar Raman
IBM Almaden Research Center
rshankar@almaden.ibm.com

Amol Deshpande *Joseph M. Hellerstein*
University of California, Berkeley
{amol, jmh}@cs.berkeley.edu



Report No. UCB/CSD-03-1231

February 2003

Computer Science Division (EECS)
University of California
Berkeley, California 94720

Using State Modules for Adaptive Query Processing*

Vijayshankar Raman
IBM Almaden Research Center
rshankar@almaden.ibm.com

Amol Deshpande Joseph M. Hellerstein
University of California, Berkeley
{amol, jmh}@cs.berkeley.edu

Abstract

We present a query architecture in which join operators are decomposed into their constituent data structures (State Modules, or SteMs), and dataflow among these SteMs is managed adaptively by an Eddy routing operator. Breaking the encapsulation of joins serves two purposes. First, it allows the Eddy to observe multiple physical operations embedded in a join algorithm, allowing for better calibration and control of these operations. Second, the SteM on a relation serves as a shared materialization point, enabling multiple competing access methods to share results, which can be leveraged by multiple competing join algorithms. Our architecture extends prior work significantly, allowing continuously adaptive decisions for most major aspects of traditional query optimization: choice of access methods and join algorithms, ordering of operators, and choice of a query spanning tree.

SteMs introduce significant routing flexibility to the Eddy, enabling more opportunities for adaptation, but also introducing the possibility of incorrect query results. We present constraints on Eddy routing through SteMs that ensure correctness while preserving a great deal of flexibility. We also demonstrate the benefits of our architecture via experiments in the Telegraph dataflow system. We show that even a simple routing policy allows significant flexibility in adaptation, including novel effects like the automatic “hybridization” of multiple algorithms for a single join.

1 Introduction

It is often difficult to predict values of the parameters that govern database query execution. Cardinality estimates are highly imprecise [SLMK01, BC02], and competing demands on memory, system load, and network bandwidth are typically known only at runtime [P⁺93b, P⁺93a, ZL97]. In federated and web database systems, data distributions and rates often cannot be known in advance [ZR02, UF00, VN02]. Even for a single data source, statistical properties vary over time; this is of particular concern in continuous query systems [M⁺02, BBD⁺02]. Interactive query systems introduce another parameter that can vary during query execution: user preferences [H⁺99].

Such uncertainties have led to a focus on adaptive execution in many recent query systems, including Tukwila, Telegraph, Aurora, Query Scrambling, and STREAM [I⁺99, Tel, C⁺02, UFA98, BBD⁺02]. Perhaps the most adaptive of these approaches is the Eddy operator [AH00] of Telegraph, which executes queries by routing tuples between query modules such as selections and joins, dynamically reconsidering the ordering of such modules on a per-tuple basis.

*This work was supported by the NSF under Grants 0122599 and 0208588, a UC MICRO grant, a Microsoft Fellowship and gifts from Microsoft, IBM and Intel. Infrastructure for the research was provided by NSF grant EIA-9802069.

This paper presents an adaptation mechanism that substantially enhances the power of the Eddy, allowing continuously adaptive decisions for most of the major aspects of traditional query optimization: not only the ordering of operators, but also the choice of access methods, join algorithms, and the selection of a spanning tree in the query graph [IK84, KBZ86]. Our core idea is to refine the granularity of query modules, by breaking up join modules and elevating the data structures typically encapsulated within them into separate State Modules (SteMs).

The Join is a logical construct in the relational algebra; join algorithms typically involve multiple physical operations. The motivation behind splitting joins into SteMs is to decouple the physical operations that are typically encapsulated within join modules. This exposes these physical operations directly to the Eddy, for performance calibration, fine-grain routing adaptation, and work sharing.

Informally, a SteM is a half-join. It encapsulates a dictionary data structure over tuples from a table, and handles *build* (insert) and *probe* (lookup) requests on that dictionary. We show that all select-project-join queries can be executed by routing tuples carefully between access methods on data sources, SteMs, and selections. Join algorithms are *not explicitly programmed*, but are instead captured in the routing of tuples between SteMs and the access methods on the data sources.

The breaking of algebraic join encapsulation has two benefits. First, the Eddy can now monitor and control physical operations that are normally hidden within joins. By adapting the tuple routing to SteMs the Eddy adapts the order of these physical operations, and thereby the join algorithm itself. We will see an example in Section 4.2 where this allows the Eddy to distinguish between cached and uncached lookups in a networked index join, resulting in a simple routing policy with better performance than the corresponding join algorithm from the literature. In fact, by appropriate routing the Eddy can even simulate *hybrid* join algorithms that combine elements of different traditional algorithms. For example, we shall see an experiment in Section 4.5 where the Eddy “hybridizes” index and hash join algorithms, gradually converting one into the other during query execution.

Second, SteMs provide a shared data structure for materializing and probing the data accessed from a given table, regardless of the number of access methods or join algorithms involving that table. This sharing is especially useful for access method adaptation. The choice of access methods is difficult in federated systems [Tel, H⁺97, I⁺99], because a given table may be provided by multiple data sources, and a single source may support multiple access methods corresponding to different sets of bind-*fields*. An Eddy can run multiple access methods *concurrently*, and dynamically choose among them based on observed performance. The use of SteMs helps avoid redundant work during this competition; all access methods on a table build into the same SteM. Moreover, although an Eddy routing policy can effectively try out multiple competing join algorithms, all lookups on a table probe the same SteM, taking advantage of the shared materialization.

The flexibility enabled by SteMs comes with a challenge: arbitrary routing from multiple access methods through SteMs may not correspond to a valid query execution plan. Incorrect routing can lead to duplicate results, missing results, or infinite routing loops. Therefore we develop a set of constraints on the routing that guarantee correct query execution (Section 3), while preserving opportunities for the flexible kinds of adaptation described above.

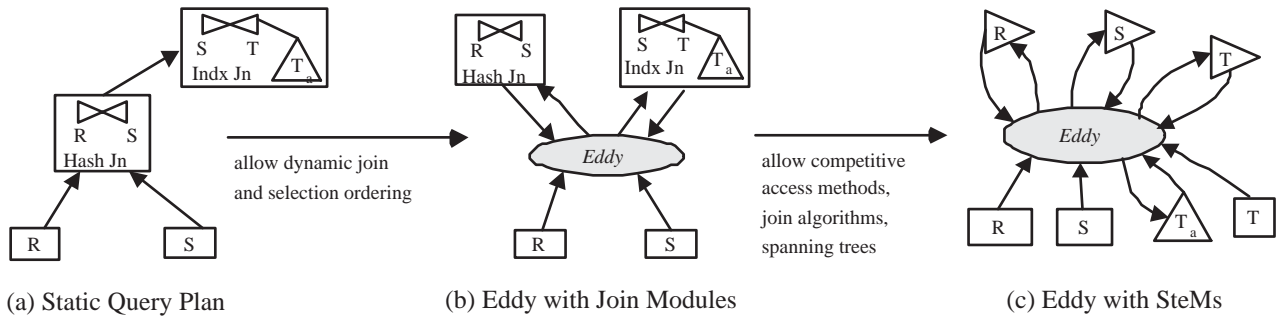


Figure 1: A three table join performed in three ways. The Eddy continually routes tuples between modules, which run as concurrent threads. Indexes are represented by triangles, and shown encapsulated within index join as suggested in [AH00]. SteMs are shown as sideways triangles (“half-bowties,” to signify that they are half-joins).

1.1 An Example

Consider a join of three tables R , S , and T , with equi-join predicates between R – S and S – T . Suppose there is a scan access method on each relation, and an index access method on T corresponding to the join attributes with S . Figure 1 shows three ways of running this query. Figure 1(a) is a traditional, statically chosen query plan involving a hash join and an index join. Figure 1(b) shows the approach of [AH00] where an Eddy is used to dynamically adapt the join order by controlling the tuple flow between the joins. Note that both these approaches make use of only the index access method on T , and a pre-chosen implementation for the RS and ST joins. Figure 1(c) shows the same query being executed with SteMs. All access methods over data sources are treated as query modules, and are used simultaneously. Tuples coming into the Eddy from these access methods are not routed to joins, but instead to SteMs and other access methods. This plan allows the use of all the access methods, and a variety of routing decisions that correspond to different join algorithms and join orders. We develop the details of this approach in the body of the paper.

1.2 Background

The setting for this work is Telegraph, an adaptive dataflow system for querying streams of networked data [Tel]. An early application of Telegraph was Federated Facts and Figures (FFF), a query system to combine data from diverse and distributed data sources. These include not only relational databases but also websites providing services and data backed by databases (the so-called “Deep Web” [RGM01, Lex]). For example, here is a query over three Web sources – a listing of contributors to election campaigns from the Federal Election Commission (FEC), a database of demographic information about neighborhoods from Yahoo (Demographics), and a database of crime ratings by area from APBNews.com (Crime).

```
SELECT  AVG(F.contribution), AVG(D.householdIncome), F.State, C.crimeRating
FROM    FEC as F, Demographics as D, Crime as C WHERE F.zip = D.zip and F.zip = C.zip
GROUP BY F.State, C.crimeRating
```

Among the various factors that we discussed earlier, our interest in adaptive query execution is motivated by two unpredictable properties in FFF:

Volatility of distributed data sources: Since Web sources are autonomously maintained, their speeds and availability are hard to estimate at optimization time, and could vary during query execution.

Volatility of user interests during online query processing: Since users often specify queries in an iterative, exploratory fashion, FFF uses an online performance metric [H⁺99, RH02] and gives out partial results during query execution. As the user sees these partial results, their interests in different parts of the result may change.

1.3 Outline of the paper

In the rest of the paper, we develop the SteM mechanism, and show how it helps in an environment like FFF. We begin with a description of the modules in our architecture (Section 2), and then describe how arbitrary select-project-join queries can be executed *correctly* using these modules (Section 3). Next, we present an experimental study that illustrates the various kinds of adaptations allowed by SteMs, and the performance benefits we get under an online query processing metric (Section 4). We discuss related work in Section 5, and conclude with a discussion of other implications of SteMs and directions for future work (Section 6).

2 Mechanics of Query Execution with SteMs

In this section, we first describe the modules in our architecture, including the State Modules (SteMs), and discuss how they are instantiated for an arbitrary query. We then illustrate a simple but important example of these modules in use: an n -ary version of the *symmetric hash join* operator.

2.1 Eddy, State, Access, and Selection Modules

Our architecture uses four kinds of modules: (1) selection modules that correspond to query predicates, (2) access modules that correspond to access methods over data sources, (3) State Modules (SteMs) that encapsulate data structures used in traditional join algorithms, and finally (4) an Eddy module that routes tuples between the other modules¹. Each module runs asynchronously in a separate thread, though this asynchrony can also be achieved in a single-threaded implementation [S⁺01].

We now describe the module functionality in detail. Simplified pseudo-code is given in Table 1. We start with some definitions.

Definition 1 (Base-table Component, Span) Consider a tuple t that belongs to the join of k base-tables T_1, T_2, \dots, T_k . The projections of $\{t\}$ on the columns from each of these base-tables form relations with a single row each. Each of these rows are called the base-table components, $t_{T_1}, t_{T_2}, \dots, t_{T_k}$, of t . We denote t by $\langle t_{T_1}, t_{T_2}, \dots, t_{T_k} \rangle$, and say that t spans the tables T_1, T_2, \dots, T_k .

Definition 2 (Singleton tuple) A singleton tuple is one that contains a single base-table component.

¹As in most discussions of query plans, we do not devote explicit modules to projection, but assume it is done on the fly by each module as aggressively as semantically possible. Group By, Aggregation, and complex SELECT-list expressions are implemented above the Eddy, before results are output to the user.

Module	Input tuple	Output tuple(s)	Action
SM	t	t or nothing	Bounce back t iff it matches predicate
AM	t	t matches for t EOT	Asynchronously bounce back t Asynchronously return all matches for t Return EOT after all matches have been returned.
SteM	$build_t$ EOT $probe_t$	— — $build_t$ or nothing — concatenated results $probe_t$ or nothing	Build $build_t$ into the SteM. Build EOT into the SteM. Asynchronously bounce back $build_t$ if needed for correctness (Section 3). Find matches for $probe_t$ among tuples in SteM. Concatenate these matches with $probe_t$ and return concatenated results. Asynchronously bounce back $probe_t$ if needed for correctness (Section 3).

Table 1: Functionality of the main query processing modules in our architecture.

2.1.1 Eddy Module

The Eddy’s role is to continuously route tuples among the rest of the modules, according to a *routing policy*. When a module other than the Eddy processes a tuple t , it can generate other tuples and send them back to the Eddy, for further routing. It can also optionally return (or *bounce back*) t to the Eddy if t requires additional processing. A tuple is removed from the Eddy’s dataflow and sent to the output if it spans all base tables and is verified to pass all predicates. The Eddy terminates the query when there are no tuples in the dataflow, and each module has finished processing all the tuples sent to it.

Each tuple also carries some state with it, called its *TupleState*, to track the work it has done in furthering query progress. The exact structure of TupleState depends on the routing policy. However, as a bare minimum, the TupleState must contain (a) the tables spanned by the tuple, and (b) the predicates that the tuple has passed (our implementation uses a bitmap, like the *done* bits of [AH00]). The former denotes the type of the tuple, and the latter is used by the Eddy to decide when the tuple is ready for output. In fact, this state alone suffices for all but one special class of *cyclic* queries; we will discuss the exception in Section 3.4.

2.1.2 Selection Modules (SMs)

Selection modules (SMs) are simple. When a selection module M receives an input tuple t , it returns t to the Eddy if t passes the selection predicate, and removes it from the dataflow otherwise. If t passes the predicate, M marks this fact in t ’s TupleState, so that the Eddy can track the progress made by t .

2.1.3 Access Modules (AMs)

An Access Module (AM) encapsulates a *single* access method over a data source – it can either be a scan, or an index on some set of columns. Each access method on a given relation is encapsulated in a separate AM.

A tuple t that is routed to an AM is called a *probe tuple*, and corresponds to a request for the AM to output tuples that “match” the probe tuple – the matches from an AM on table S are all $s \in S$ such that the concatenation of t and s satisfies all query predicates that are defined over the union of the columns spanned by s and t^2 . Note that the output schema of an AM is the same as that of the data source. In particular, the

²Some of these predicates will be enforced by the index lookup, the AM applies the others after the lookup.

AM does not concatenate the probe tuple to its output tuples. Such concatenation will be performed only by SteMs.

Scans are also treated as AMs, but only accept a special empty probe tuple we call a *seed tuple*, and in return, output all tuples in their data source. At query initialization, each scan AM is initialized by passing it a seed tuple, which informs it to begin returning the contents of the full scan to the Eddy.

In addition to returning matches, AMs asynchronously bounce back each probe tuple t to the Eddy. Intuitively the bounce back is required because the probe tuple is needed later, for eventual concatenation with each of its matches. This is discussed in more detail in Section 3.3.

Asynchronous Indexes and EOTs: As demonstrated in [GW00], the throughput of accesses to Web sources can be improved significantly by sending multiple asynchronous probes; similar arguments can be made about asynchronous random disk I/Os. In this spirit, we assume that all AM probes and responses are asynchronous. This asynchrony complicates issues somewhat, because the system needs to track when all matches have been returned for a given probe. We use the dataflow itself to pass this information. When an AM on a table T has returned all matches to a probe, it sends an *End-Of-Transmission (EOT) tuple* encoding the probing predicate (in the case of a scan AM, the predicate is simply “true”). In the common case of index lookups using equality predicates, the EOT tuple is a regular tuple with a special *EOT* value in all the non-bound fields (e.g., $\langle 15 \text{ John } EOT \ EOT \dots \rangle$ if the probe tuple binds the first two fields to 15 and John). For non-equality predicates, the EOT tuple contains pointers to the predicates, which are stored in a data structure created during query parsing. For scans, the EOT predicate contains the predicate “true”. The advantage of encoding EOTs as tuples rather than as control messages is that the EOTs can be stored in SteMs itself, alongside standard tuples, as we will see below.

2.1.4 State Modules (SteMs)

A SteM essentially corresponds to half of a traditional join operator. It stores homogeneous tuples (tuples spanning the same set of tables) formed during query processing, and supports insert (build), search (probe), and optionally delete (eviction) operations. In this paper, we only consider SteMs over base tables; *i.e.*, all tuples in a SteM are singleton tuples from the same table. As such, all joins on a given base table can and do use the same SteM for builds and probes involving that base table. For this purpose, we allow a SteM to perform searches on arbitrary predicates.

Two kinds of tuples can be routed to a SteM. When a *build tuple* $t \in T$ is routed to $SteM_T$, t is added to the set of tuples in $SteM_T$, and the indexes, if any, are updated accordingly. An EOT tuple from an AM on T is also routed as a build tuple to $SteM_T$. When a *probe tuple* p is routed to $SteM_T$, $SteM_T$ returns *concatenated matches* for it to the Eddy. These concatenated matches are all tuples in $\{p\} \bowtie SteM_T$ that satisfy all query predicates that can be evaluated on the columns in p and T .

Note that since the SteM is continually being built, it may not have all the tuples in $\pi_T(\{p\} \bowtie T)$. This is tracked by the presence of EOT tuples. If an EOT tuple in $SteM_T$ matches a probe p , then $SteM_T$ knows that it definitely contains all matches for a probe p . If not, the SteM might have to bounce back p so that it can be routed to other modules (to find the missing matches).³ The logic for when such bounce backs are

³This logic is simplified, and assumes that tuples from an AM arrive at the SteM in order. In fact, in our implementation the Eddy reorders tuples in the dataflow to match user interests [RH02]. So the EOT tuple for a probe could be built into a SteM before

needed is determined by the routing constraints, and will be developed in Section 3.

In our present implementation, we speed up join predicate lookups through indexes. A SteM on a table T (called $SteM_T$) has one main-memory index (hash table or binary tree) on each column of T that is involved in a join predicate. These are all secondary indexes having pointers to the same tuples in memory. We do not focus on disk-resident indexes in this paper because the datasets we have encountered in Web sources are typically small enough to fit in main memory. We defer discussion of multi-table SteMs and disk data management within SteMs to Section 6.

2.2 Query Planning

The use of Eddy and SteMs obviates the need for query optimization because there are no *a priori* decisions to be made. Unlike in [AH00], there is no need even for a “pre-optimizer” that chooses the join implementations, access methods, and query spanning tree. The query is instantiated as follows :

1. Check that the query is valid, *i.e.*, it can be executed given the bind-field constraints on the data sources (we use the algorithm from Nail [Mor88]).
2. Create an AM on each access method that can possibly be used in the query.
3. Create a SM on each predicate in the query.
4. Create a SteM on each base table in the query.
5. Create any seed tuples needed for scans (Section 2.1.3).

As described in the earlier section, only one SteM is created per data source. This SteM is shared not only among the join predicates involving that data source, but also among multiple instances of the source in the FROM clause, if any exist (*e.g.*, a self-join).

Though this paper focuses on execution of a single query, a SteM can also be used to share work and storage across concurrent queries. Related work in Telegraph uses SteMs in this way, in the context of continuous query processing [M⁺02, CF02].

2.3 Example: An N -way Symmetric Hash Join

We now give an example of how these modules can be used to implement an n -way version of the symmetric hash join (SHJ) [RS86, WA91]. The traditional, binary SHJ is a pipelining join that works by simultaneously building hash tables on both its inputs. Each input tuple is first built into a hash table on that input, and then immediately probed into a hash table on the other input. Due to its pipelining nature this operator is well-suited for interactive processing. Though originally designed as a memory-resident algorithm, it has subsequently been extended by [I⁺99] and [UF00] to spill to disk in memory-constrained environments.

There are two ways to extend the SHJ to multi-table queries. Consider an equi-join $R \bowtie_a S \bowtie_b T$.

Pipelining Binary Joins: Figure 2(i) shows how multiple binary SHJs can be pipelined to perform an n -way SHJ. To the best of our knowledge, this is the approach of choice in all current literature (*e.g.*, [UF01]).

n -ary SHJ Operator: Figure 2 (ii) shows how all the SHJs can be unified into a single operator that uses four hash indexes: one on R , one on T , and one on each join column of S (one of these is a secondary index). When a new R (T) tuple comes in, it is first built into the corresponding hash index H_{R_a} (H_{T_b}),

all the matches for that probe are built. To solve this, we tag the EOT tuple with the number of probe matches, which allows a SteM to verify if all matches have been built into it.

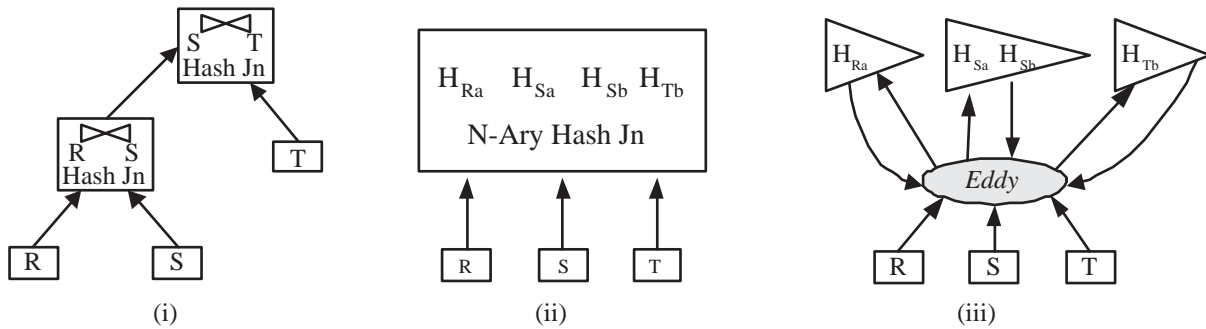


Figure 2: Three ways of doing a 3-table symmetric hash join (SHJ): (i) with pipelined binary SHJs, (ii) with a 3-ary SHJ operator, and (iii) with an Eddy and SteMs

and then probed into H_{S_a} (H_{S_b}). The resulting matches, if any, are then used to probe into H_{T_b} (H_{R_a}) and the result is output. When a new S tuple comes in, it is similarly built into H_{S_a} and H_{S_b} . At this point, we have a choice, corresponding to different join orders. We can either probe the S tuple into H_{R_a} and probe H_{T_b} with the resulting matches, or we can probe into H_{T_b} and then into H_{R_a} .

The initial Eddy paper [AH00] was based on the first approach – by connecting a set of pipelining binary join modules to an external Eddy module, the ordering of the join modules can be decided dynamically. In contrast, the SteMs mechanism is based on the second approach – it essentially places an Eddy *within* the n -ary SHJ operator, so that the ordering of the hashtable lookups can be decided dynamically. This is the core effect of SteMs – to give the Eddy access to the data structures typically stored inside join algorithms. However, the SteM approach is not implemented as part of the SHJ, and therefore becomes more generally applicable.

Figure 2 (iii) illustrates the translation from the unified n -ary SHJ operator to a routing through SteMs. We use a SteM on each source to encapsulate the hash indexes on that source, and an Eddy to route tuples between the SteMs. Each tuple is first built into a SteM on its source, and then immediately routed to the other SteMs. The Eddy can dynamically adapt the join ordering by changing the way it routes S tuples after it is built into $SteM_S$.

In addition to different routing opportunities, the n -ary hash join materializes different state than the traditional binary-SHJ scheme. Note that the n -way SHJ description above stores only singleton tuples in hash tables, whereas the traditional pipeline of binary SHJs materializes intermediate result tuples from joins below the root (e.g., tuples in $R \bowtie_a S$). SteMs can in principle support either scheme, or both, via a SteM to materialize each base-table or intermediate relation desired. This represents a tradeoff of performance for memory space – less memory is likely to be used if intermediate result tuples are not stored, but more probes may need to be made since the same intermediate results may need to be recomputed multiple times. In this paper (as in [M⁺02, CF02]), we choose not to store intermediate tuples in SteMs. In addition to the space/time tradeoffs, a secondary advantage of not materializing intermediate results is that tuple eviction is simplified. Each base-table component is stored in a single SteM, and so it can be easily evicted by the SteM if needed. Although not the focus of this paper, sliding-window queries and queries over unbounded data streams require tuple eviction, and [M⁺02, CF02] both use SteMs with eviction. We are currently investigating a hybrid approach that partially materializes intermediate results to the extent of available

memory (Section 6).

The n -ary SHJ can be used for any select-project-join query where all sources have scan access methods. In the next section, we generalize this simple operator to use other join algorithms as well as index access methods, and show how the Eddy can dynamically adapt the join algorithms, access method choices, and spanning tree choices.

3 Executing Arbitrary Select-Project-Join Queries with SteMs

Superficially, query execution with SteMs is simple. We only need to instantiate the AMs, SteMs, and SMs, as in Section 2.2, and let the Eddy route tuples through these operators. The problem is that arbitrary routing policies need not lead to correct results or terminating queries. Since we want the Eddy to adapt the routing dynamically, we now develop *constraints* on the routing policy that will ensure correctness.

The n -ary SHJ operator corresponds to one correct routing policy. We start by identifying the routing constraints that are implicit in this operator, and gradually generalize these constraints to a larger space of queries. Our presentation is intended to be intuitive and informal; proofs of correctness are in Appendix A.

3.1 Acyclic SPJ queries with a Single Scan AM on each Table

The n -ary SHJ is captured by two rules. The first is that the SteMs be implemented with hash indexes. The second is that the Eddy must obey the following routing constraints:

BuildFirst: A singleton tuple from a table T must first be routed to build into $SteM_T$.

SteM BounceBack: A SteM must always bounce back build tuples (so that they can probe the other SteMs for matches), and never bounce back probe tuples.

Atomicity: The building of a singleton tuple into a SteM must be atomically coupled to the probing of that tuple into the other SteMs.

BoundedRepetition: No tuple must be routed to the same module more than once.

The first three constraints capture the essence of the n -ary SHJ, and BoundedRepetition ensures query termination. Two relaxations of these constraints allow the Eddy to adapt over a much wider space of join algorithms.

Constraint Relaxation to allow other Join Algorithms

Our first relaxation removes the constraint that the SteMs must be implemented with hash indexes. For example, the SteM may use a linked list when it holds a small number of tuples, and switch to a hash-based implementation when the list size increases. This switch can be made independent of the other modules.

Our next relaxation is to remove the Atomicity constraint, and *decouple* the build and probe operations of each tuple. This allows the Eddy to interleave probes and builds of tuples in arbitrary ways, and thereby change join algorithms (in Section 3.5, we will relax this further by allowing the build to be completely avoided). Unfortunately, this build-probe decoupling can cause duplicate query results. For example, Figure 3 shows four steps in a SHJ. If $\langle r_1, s_1 \rangle$ satisfies the join predicate, $\langle r_1, s_1 \rangle$ is output at both step 3 and

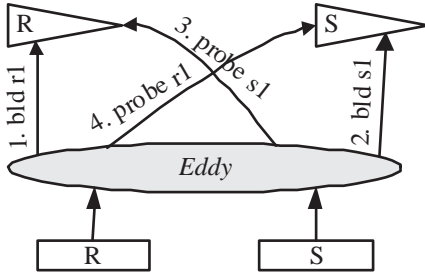


Figure 3: Duplicates arise because of decoupling build and probe of r_1

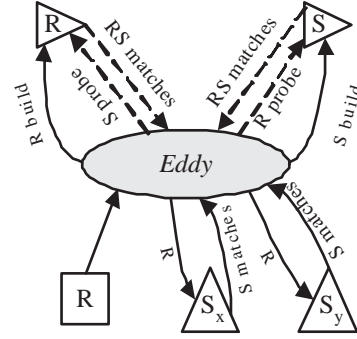


Figure 4: $R \bowtie S$ query with index AMs on S

step 4 because the builds and probes of r_1 and s_1 tuples are interleaved. To avoid such duplicates, we add a TimeStamp constraint [Ram01], to form the following set of constraints:

BuildFirst: A singleton tuple from a table T must first be routed to build into a SteM on T.

SteM BounceBack: SteMs must always bounce back build tuples, and never bounce back probe tuples.

BoundedRepetition: No tuple must be routed to the same module more than once.

TimeStamp:

- Each singleton tuple t is assigned a global *Timestamp* $TS(t)$ (wall-clock time) when it builds into a SteM. Before building, $TS(t)$ is defined to be ∞ . For other tuples $TS((t_1, \dots, t_n))$ is defined to be $\max(TS(t_1), \dots, TS(t_n))$ i.e., the timestamp of its last base-table component.
- When a tuple r probes into a SteM and finds a match s , the result $\langle r, s \rangle$ is returned to the Eddy iff $TS(r) > TS(s)$.

The TimeStamp constraint says that only the last arriving base-table component of a result tuple will generate that tuple, by probing into other SteMs to join with previously-arrived components.

Simulating and Hybridizing Non-Pipelined Join Algorithms

These relaxed constraints allow the Eddy to simulate several join algorithms besides the SHJ. Consider a two table join of R and S. The following policy can simulate many non-pipelining join algorithms:

1. Route all R tuples to build into $SteM_R$
2. Route all S tuples to build into $SteM_S$
3. Route all S tuples to probe into $SteM_R$
4. Route all R tuples to probe into $SteM_S$

The SteM implementation decides exactly which join algorithm will be simulated. E.g., the following “asynchronous” hash index implementation simulates a *Grace Hash Join* [FKT86]. While build tuples are routed to $SteM_R$ and $SteM_S$, the SteMs create hash partitions on disk. But instead of bouncing back these build tuples immediately, they do so asynchronously, *clustered by the hash partition*. Therefore in Step 3, when the bounced-back S tuples probe $SteM_R$, $SteM_R$ gets very good I/O locality. Because of the TimeStamp constraint, Step 4 does not produce any results. It can be completely avoided by maintaining in each SteM the minimum timestamp of all tuples the SteM contains – the Eddy need only route to a SteM

probe tuples with timestamp greater than this minimum timestamp.

It is unusual to describe Grace Hash Join in terms of a routing policy. But the advantage is that the Eddy can now dynamically *hybridize between SHJ and Grace Hash Join*, by changing its routing as follows. Rather than do all of Steps 1 and 2 before Steps 3 and 4, the Eddy can dynamically decide to interleave them. Specifically, when a tuple r is bounced back after building into $SteM_R$, the Eddy may choose to immediately probe r into $SteM_S$. This choice is based on the level of interactivity desired by the user. For instance, the Eddy can start with frequent probes to give interactive responses early on, and later degenerate to occasional probes in order to reduce completion time (when probes are infrequent more probes for the same partition are clustered together, so I/O cost is lesser). The frequent probe phase simulates an SHJ, and the occasional probe phase is similar to Grace.

An exactly analogous implementation of SteMs with tournament trees that spill sorted runs to disk will simulate a Sort-Merge join [Knu73]. The Hybrid-Hash Join [D⁺84] is simulated if the SteMs maintain a full in-memory hash table on some of the partitions and bounce back build tuples for these partitions ahead of others. The Eddy can then route S tuples from these in-memory partitions to probe into $SteM_R$ even before all S tuples have been built.

Note that one part of the join logic – choosing whether the indexes are hash indexes or tournament trees – is captured in the SteMs implementation. It is up to the SteM implementation to internally adapt this if needed. But the remaining part, i.e., the interleaving of builds and probes, is captured in the routing policy, and can be dynamically adapted by the Eddy.

3.2 Competitive AMs

We now expand our class of queries to include those over tables with more than one AM. Such alternate AMs are very common for Web sources in Telegraph FFF. Different websites often provide the same data, and a single website may support multiple AMs corresponding to different sets of fields that can be chosen as the lookup key. We address tables with multiple scan AMs in this section, and discuss index AMs in the next section.

Traditional database systems typically pick one AM per data source at optimization time. We want to be able to run multiple AMs on a single source in competition with one another, and let the Eddy dynamically choose one AM, or switch between AMs. For example, if a particular AM stalls because the underlying source is delayed, the Eddy should be able to use the alternate AMs. In our architecture, this is quite straightforward to do since all the access methods are exposed to the Eddy. The main problem turns out to be duplicates; the same tuple can be generated by different AMs. However because of the BuildFirst constraint, such duplicates can be easily removed when they build into the SteM on the source itself. We only need a simple enhancement to the SteM BounceBack constraint:

SteM BounceBack: A $SteM_S$ must bounce back a build tuple s unless it is a duplicate of another s' that is already in $SteM_S$.

Aside on duplicate semantics: Many rules can be applied to deal with duplicates, including set and bag semantics (e.g., see [Alb91]). When identical tuples can be identified through their value in a key column (that may be projected out in the result), the projection can be postponed till result output so that the SteM

can preserve the exact number of duplicates in the source [P⁺92]. This issue is further complicated when the AMs involve different, possibly inconsistent, Web sources. We currently adopt a set semantics, where a SteM removes any build tuple that is identical to another tuple already present in the SteM.

3.3 Index AMs

When a data source has an index AM, we encounter another problem. Figure 4 shows the execution of a simple two-table join query in this class. Recall that our indexes are allowed to return matches asynchronously. A tuple r from R is first built in $SteM_R$, and then probed into $SteM_S$ to see if matches for r have been already cached there. But unless *all* matches are already cached, r must be bounced-back by $SteM_S$, so that it can probe into one of the AMs on S . The difference from the previous section is that there is no scan AM on S , so r must probe into an index AM to seed the generation of its matches.

Subsequently, the index AMs on S will return matches for r , say s_1 and s_2 . These matches will be first built into $SteM_S$ and then probed into $SteM_R$. It is only during this probe that s_1 and s_2 will join with r (and possibly with other R tuples as well). Thus $SteM_R$'s role is as a *rendezvous buffer* [GW00] to hold pending probe tuples until matches arrive.

Since s_1 and s_2 are built into $SteM_S$, subsequent R tuples with the same bind column values as r will find index matches in $SteM_S$ itself. So $SteM_S$ will not bounce back these R tuples ($SteM_S$ verifies that it has *all* relevant matches by checking its EOT tuples). Thus $SteM_S$'s role is that of a cache on index lookups into S . In fact, when there are multiple AMs on a source, they all cooperate in building the same cache, and the work of probing alternate AMs is not wasted. We will see experimentally in Section 4.3 that this reduces the cost of competition.

When a data source has both scan and index AMs, the tuple routing determines whether an index join is performed or a hash join is performed. We will see an experiment in Section 4.5 where the Eddy dynamically adapts its routing to switch between the two during query execution.

To summarize, the enhanced SteM BounceBack constraint is as follows:

SteM BounceBack:

- A $SteM_S$ must bounce back a build tuple s unless it is a duplicate of another s' that was previously in $SteM_S$.
- A $SteM_S$ must bounce back a probe tuple r unless S has a scan AM, or $SteM_S$ already contains all matches for r .

As mentioned in Section 2.1.4, $SteM_S$ uses the presence of EOT tuples from probes into AMs on S to verify whether $SteM_S$ already contains all matches for a given probe tuple.

3.4 Cyclic Queries

Cyclicity in the query join graph complicates matters still further. Traditionally, the plan chosen by the query optimizer contains join modules only over a spanning tree of the query join graph. This spanning tree is determined before query execution, even for prior adaptive query processing schemes like the initial Eddy paper [AH00]. Static spanning tree choices hurt in two ways:

- The spanning tree choice is typically made based on selectivities, which are hard to estimate for queries over Web sources. So the resulting execution strategy can be arbitrarily sub-optimal.
- A static spanning tree choice can also constrain the generation of partial query results. Consider a three way join of R, S, T where there are join predicates between each pair of tables. If we choose $R \bowtie S \bowtie T$ as the spanning tree and source S stalls during query execution, the entire query blocks. If the spanning tree could be changed dynamically, RT tuples could be generated. These partial results with missing values for S columns could be very valuable in interactive querying environments [RH02].

The problem with not fixing a spanning tree a priori is that duplicates can arise even after timestamping. Consider the following sequence of events in the above 3-way join query: (1) a tuple t probes into $SteM_S$ to find a match $\langle s, t \rangle$, (2) $\langle s, t \rangle$ probes into $SteM_R$ to find a match $\langle r, s, t \rangle$, (3) $SteM_S$ bounces back t as per the SteM BounceBack constraint, (4) t probes into $SteM_R$ to produce $\langle r, t \rangle$ which probes into $SteM_S$ to produce $\langle r, s, t \rangle$ again.⁴

To avoid such duplicates, we must ensure that previously bounced-back tuples (like t) cannot probe other SteMs.

ProbeCompletion Constraint: A tuple t that has been bounced back after probing into a $SteM_S$ must not probe into any other SteM afterwards. The routing policy must however maintain t in the dataflow, routing it to other modules, until it has been probed into an AM on S .

Definition 3 (Prior Probers, Probe Completion Tables) *Tuples like t that have been bounced back after probing into SteMs are called prior probers. The corresponding table S is called the probe completion table of t , and the AMs on S are called the probe completion AMs of t . The identity of the probe completion table is marked in the TupleState of t .*

3.5 Relaxing the BuildFirst Constraint

The constraints developed so far guarantee that all select-project-join queries will be executed correctly. But one of these constraints, the BuildFirst constraint, is particularly restrictive and could result in highly inefficient execution in situations where one of the input tables is much larger than the others. Suppose that the R table was much larger than both S and T tables in the example of Figure 2(iii). In that case, it might be better to build SteMs on the S and T tuples and probe the R tuples directly into these two SteMs, without building into $SteM_R$. This is equivalent to building temporary index on only one side of the join.

We can enable such optimizations by allowing the Eddy to *not build* a $SteM$ on a table R as long as there is only one access method on R and that access method is *scan*. If there multiple access methods on R or if there is an index AM on R , the $SteM$ is required to avoid duplicate results.

Now if an R tuple is bounced back from a $SteM_S$, it means that all S matches for this R tuple could not be found at that time. So this R tuple needs to be routed back to $SteM_S$ to find the remaining matches. So we relax the BoundedRepetition constraint to allow the Eddy to route a given tuple repeatedly to the same module. To ensure that these repeated probes do not produce duplicates, we assign every R tuple a *LastMatchTimeStamp*. This is initially set to 0. Every time the R tuple is routed to $SteM_S$, the

⁴Note that this only happens if there is no scan AM on S , because otherwise $SteM_S$ does not bounce back the t tuples sent to it (Section 3.3).

Constraints to be enforced by Routing Policy Implementor	
BoundedRepetition	- No tuple can be routed to a given module more than a finite number of times.
BuildFirst	- A singleton tuple from a table T must first be built into $SteM_T$ iff T has multiple AMs or, T has an index AM
ProbeCompletion	- A prior prober t must not be routed to any SteM other than that on its probe completion table. - The Eddy can remove a prior prober t from the dataflow only after t has been probed into one of t 's probe completion AMs.

Constraints enforced within SteM and AM implementation	
SteM BounceBack	- A $SteM_S$ must bounce back a build tuple s unless it is a duplicate of another tuple s' that is already in $SteM_S$. - A $SteM_S$ must bounce back a probe tuple r unless $SteM_S$ already contains all matches for r , or S has a scan AM, and all base-tuple components of r have been cached in other SteMs
TimeStamp	- When a tuple r probes into a SteM and finds a match s , the result $\langle r, s \rangle$ is returned to the Eddy iff $TS(r) > TS(s) > LastMatchTS(r)$.

Table 2: Routing constraints that ensure correct query execution

$LastMatchTimeStamp$ is updated to the maximum of the timestamps of all tuples in $SteM_S$.

The constraints we have developed so far are summarized in Table 2. Notice that the SteM BounceBack and Timestamp rules are implemented internally to the AMs and SteMs, and the routing policy implementor need not be aware of them at all. The following correctness theorems are proved in Appendix A.

Theorem 1 (Duplicate Avoidance) *If the Eddy follows a routing policy that satisfies the constraints of Table 2, there will not arise duplicate versions of any tuple in the dataflow, other than singleton tuples that have not yet been built into SteMs.*

Theorem 2 (Correctness) *If the Eddy follows a routing policy that satisfies the constraints of Table 2, it will not output any tuple that is not in the query result, and will output all query result tuples in a finite number of routing steps.*

4 Experimental Results

We now illustrate the kinds of adaptation that SteMs enable, through an experimental study. Our focus is on the online metric of maximizing the rate at which result tuples are generated, though some of the experiments also demonstrate the effectiveness of our system for the traditional metric of completion time. All our experiments are based on an implementation of SteMs in Telegraph [Tel], and were run on a lightly loaded machine with dual 666MHz Pentium-III processors and 768MB RAM, running Redhat Linux 6.0. The salient points of our experimental study are as follows :

- Even a simple join algorithm like the index join encapsulates multiple physical operations, and this causes a *head-of-line blocking* problem. This problem can be avoided by breaking the join module into SteMs.

- SteMs allow the Eddy to efficiently learn between competitive access methods, while doing almost no redundant work.
- SteMs allow the Eddy to dynamically choose the join spanning tree for cyclic queries.
- SteMs allow the Eddy to dynamically switch between an index join algorithm and a symmetric hash join algorithm during query execution.
- With SteMs, the Eddy can adaptively choose the way it reorders tuples in interactive environments.

We use synthetic data sources for our experiments so that the source properties can be easily controlled. The data sources that we use are as shown in Table 3.

4.1 Eddy routing policy

Our implementation uses a routing policy designed to maximize the value of the partial results output to the user [RH02]. The details of this policy are not needed to understand the advantages of SteMs in our experiments. We briefly summarize it here for completeness.

When a tuple t with a TupleState T is routed to a module M , the benefit $B(t, M)$ is the value of the partial result that will be output by M . This benefit depends on the expected number of matches that M will return and the user’s preferences for the matches.⁵ M also takes an expected time $C(t, M)$ to process t . To maximize the value to the user over time, the Eddy continually routes so as to maximize $B(t, M)/C(t, M)$. Clearly it is not feasible to do this optimization across all tuples. As discussed in [RH02], though, this ratio depends largely on M and the tuplestate T of t . So we only optimize at this granularity.

To this policy we add the constraints of Section 3, specialized as follows:

BuildFirst: Singleton tuples are always first built into their corresponding SteMs, regardless of whether they come from sources with multiple AMs. This simplifies our implementation, and is inexpensive because Web sources typically have data sizes much smaller than memory sizes.

SteM BounceBack: In addition to the bounce back circumstances of Table 2, we set SteMs on tables with index AMs to also bounce back any probe tuple that satisfies a predicate prioritized by the user. Notice that in the case where a $SteM_S$ has both an index AM and a scan AM, this bounce back is redundant. But, if the prioritized probes are bounced back, they can subsequently probe into AM_S . This speeds up the entry of matches for these tuples into the dataflow and thereby the output of prioritized results to the user.

4.2 Index join improvement through SteMs

We start with an experiment that shows the effect of decoupling physical operations within a join. We consider an Index join.

Consider the first query Q1 of Figure 5 that joins tables R and S . The join is an equi-join between $S.x$, a key column of S , and $R.a$. Table R has a total of 1000 tuples, with 250 distinct values of $R.a$. In a traditional query processor, this query will be executed using an index join module as shown in the Figure 6. In contrast,

⁵Unless the tuple returned by the module contains the key columns of the result, it cannot be output as a partial result at all. However, such tuples are still given a value because they can subsequently generate partial results by joining with other tuples. For details please see [RH02].

Source	Schema	Description
R	{key: integer, a: integer}	R is a table with 1000 tuples, and has a scan access method. key is its primary key, and a is a field with 250 distinct values, randomly assigned.
S	{x: integer, y:integer}	S has two keys, x and y , and has asynchronous index access methods on both of them. All S tuples have identical values of x and y .
T	{x: integer, y:integer}	T has x as its primary key, and has an asynchronous index access method on x . All T tuples have identical values of x and y .
U	{x: integer, y:integer}	U has x as its primary key, and has an asynchronous index access method on x . The set of tuples in U is the same as the set of tuples in T .
W	{key:integer}	W has an asynchronous index access method on its primary key key , and a scan access method.

Table 3: Data sources used in our experiments. Index lookups are implemented as sleeps of identical duration.

Q1	SELECT * FROM R, S WHERE $R.a = S.x$
Q2	SELECT * FROM R, S WHERE $R.a = S.x$ and $R.a = S.y$
Q3	SELECT * FROM R, T, U WHERE $R.a = T.x$ and $R.a = U.x$ and $T.y = U.x$
Q4	SELECT * FROM R, W WHERE $R.key = W.key$
Q5	SELECT AVG($R.a$) FROM R, W WHERE $R.key = W.key$ GROUP BY $W.key \% 20$

Figure 5: Queries used in our experiments (% is the modulus sign)

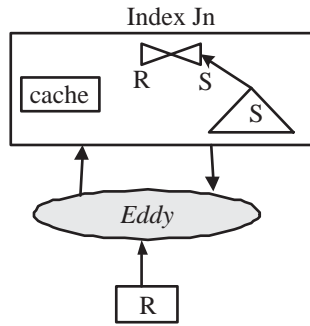


Figure 6: Executing query Q1 of Figure 5 with join modules

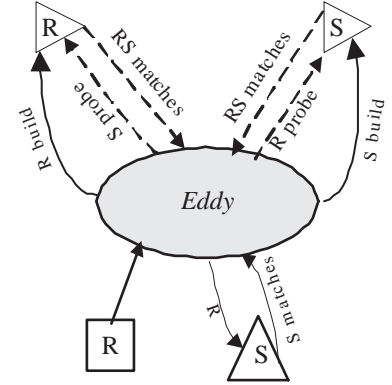


Figure 7: Executing query Q1 of Figure 5 with SteMs

our system will use a SteM on R and S , a scan AM on R , and an index AM on S (Figure 7). $SteM_R$ holds the pending R probe tuples while AM_S processes the probes, and $SteM_S$ caches probe results.

Figure 8(i) plots the number of RS results output over time in the two schemes. The curve for the plan using the index join is parabolic, as expected. The cost of probing into the index join decreases continually over time as the cache size, and hence the probability of cache hits increases. In contrast, the plan using the SteMs takes about the same time overall, but is almost linear in shape. It rises comparatively faster in the initial stages of the processing and as such, does better on our online processing metric.

To understand this behavior, we plot the number of probes into the remote source, S , for the two approaches (Figure 8(ii)). Notice that this curve is almost identical with and without SteMs. Thus the lookup caches on S build up at the same rate in both cases. The difference is that with SteMs, the probes into the caches happen much more quickly.

In the first approach (without SteMs), every tuple coming out of the scan on R does not immediately probe into the index join on S . Since all queues between the Eddy and the modules are finite in size, these probes can only happen at the speed of the index join, which in turn is bottlenecked by the speed at which the S index can handle R probes. This is unfortunate, because many of the R tuples may not need to probe into the S index at all – they may find matches in the S cache itself. With SteMs, this “head-of-line blocking” does not happen, because probes into the cache and the index have separate queues.

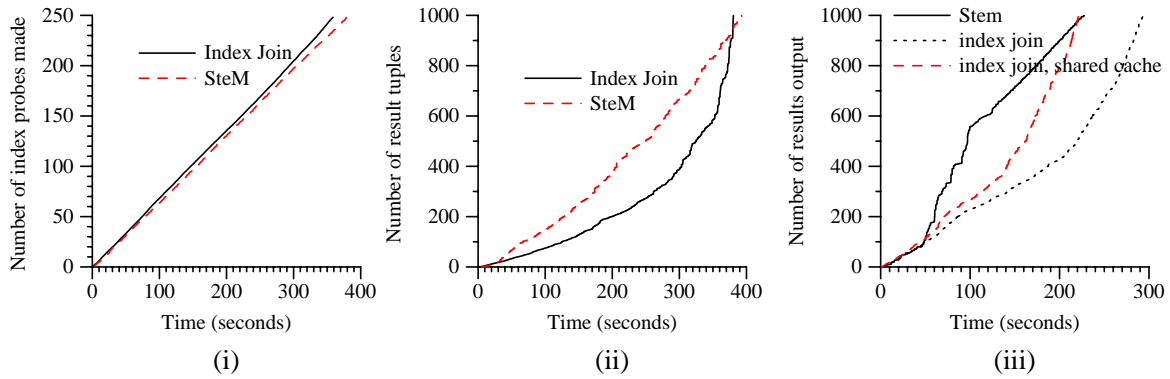


Figure 8: Number of (i) tuples output over time, (ii) probes into the S index, by the SteMs and Index Join approaches for query Q1; (iii) Performance of query Q2 with SteMs, with index joins, and with modified index joins that share caches between joins on the same source

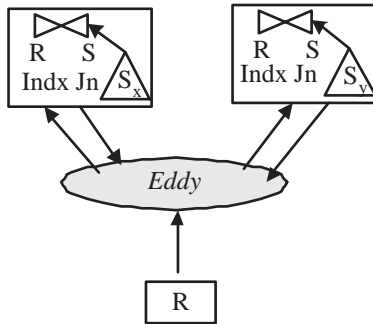


Figure 9: Executing query Q2 with traditional join modules

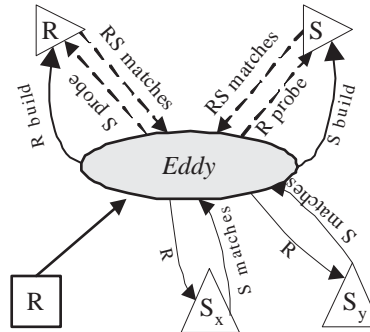


Figure 10: Executing query Q2 with SteMs

This experiment illustrates our point that even simple join operators encapsulate multiple physical operations. In this example, the index join comprises two operations, cache lookup and index lookup, that have different performance characteristics. These performance characteristics could also vary with time; *e.g.*, cache lookups may become expensive if the cache runs out of memory and starts paging to disk. Therefore it is important to avoid encapsulating such operations within the join modules.

4.3 Learning Between Competitive Access Methods – Avoiding Cache Fragmentation

Our second experiment looks at sources with alternate access methods. Our aim with this experiment is to demonstrate how SteMs can effectively reuse the work that was done while learning between the competing access methods. The query that we use for this experiment, Q2, is obtained by adding an equality predicate between $R.a$ and $S.y$ to Q1. This does not alter the query result in any way, since $S.x$ and $S.y$ values are identical for all S tuples. But as a result of adding this predicate, Q2 can now use either of the two index access methods on S , on $S.x$ or on $S.y$.

As before, we consider two ways of executing this query, one using two index join modules (Figure 9) and one using SteMs (Figure 10). The index join approach is akin to the approach taken in Oracle RDB [AZ96].

In both the approaches, we let the Eddy use first 200 R tuples to learn which of the two index methods

is better. This is done by routing these tuples randomly to the two index joins with equal probability. After this learning phase, the Eddy chooses one of the index joins and executes the rest of the query using it. Since we are mainly concerned in this experiment with the state-management overhead of learning, and not how the Eddy learns which of the two index joins is better, we set the probe costs for both the indexes to be equal and constant throughout the execution of the query.

Figure 8(iii) shows the number of RS tuples generated over time. Consider the two curves corresponding to the SteMs approach and the index join approach. Both curves have two phases. The curve for the plan using index joins rises rapidly until about 90 seconds, after which it rises more slowly. This is the point where the Eddy stops routing to the index join on $S.x$. The same effect is seen on the SteMs curve also, except that the SteM has output many more tuples at the time of the phase shift.

The SteMs curve is better than the index join curve for two reasons. First, as in the previous experiment, the SteMs approach benefits from having separate queues into the cache and the remote indexes. But note that unlike in the last experiment the SteMs curve even completes faster than the index join curve. This happens because the SteMs approach unifies the two caches in the two alternative index joins (on $S.x$ and $S.y$) that are fragmented in the alternative approach. Essentially, after the Eddy stops sending probes to the $S.x$ index join, all the tuples cached in that index join are unusable. In contrast when SteMs are used all probes into S indexes are stored in a single, unified cache, resulting in a higher proportion of cache hits.

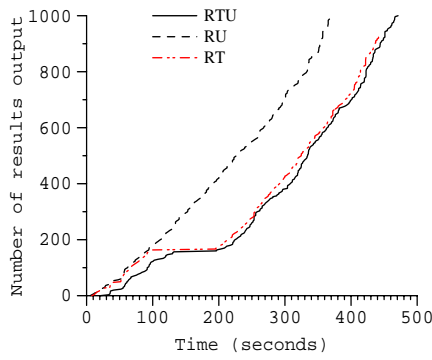
To isolate the effect of this cache fragmentation we artificially set the two index joins on S to share their caches, and rerun the query. The third curve in Figure 8(iii) (“index joins, shared cache”) shows that this approach has the same overall completion time as the SteMs approach, though the effect of “head-of-line” blocking can still be seen.

4.4 Spanning tree selection

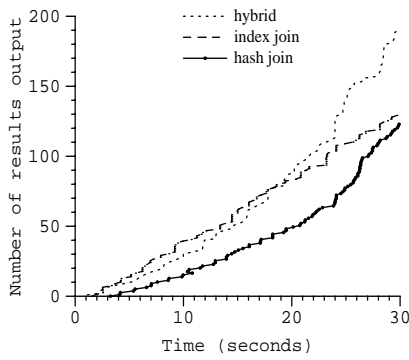
We next consider a cyclic query Q3. It involves 3 tables R, T, U with equality predicates between all three pairs of them. The traditional way to execute this query is to select a spanning tree of its join graph and create join modules along its edges, and enforce the predicate for the remaining edge with a selection module. When SteMs are used, there is no need to choose a spanning tree up front.

To see the advantage of this flexibility, we consider a situation where source T experiences a 100 second delay after the query has been run for 100 seconds. We compare two approaches: one where the Eddy creates join modules on the RT and TU join predicates, and one where the Eddy creates SteMs on R, S , and T . Figure 11 shows the number of partial and full result tuples generated over time with these two approaches. The generation of full result tuples is affected during the delay in both cases. But without SteMs, even partial result generation is severely affected since the RU join was de-selected during query optimization. Whereas with SteMs, the Eddy is able to generate RU join results during the delay.

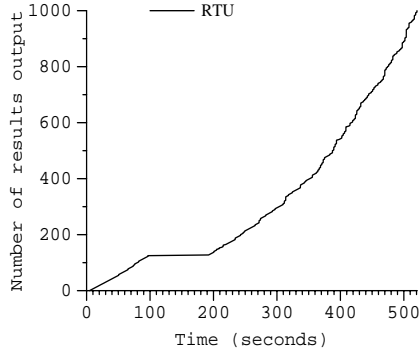
In this experiment, the advantage provided by SteMs is directly reflected in the generation of partial result tuples (RU tuples). The other experiments of this section mainly demonstrate advantages for full result tuple generation. But these advantages apply equally well to partial result tuple generation as well. For instance, the index join improvement we saw in query Q1 (Section 4.2) arises in the current experiment



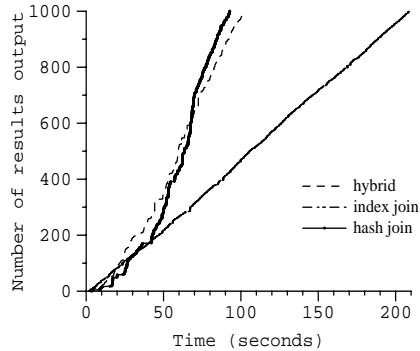
(i)



(i)



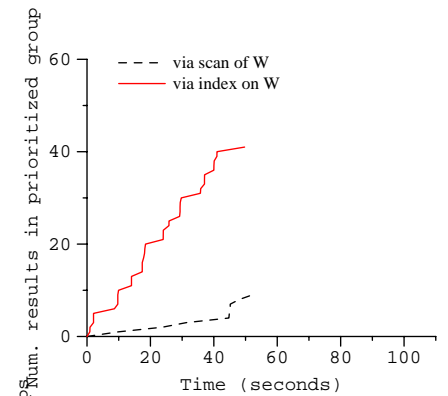
(ii)



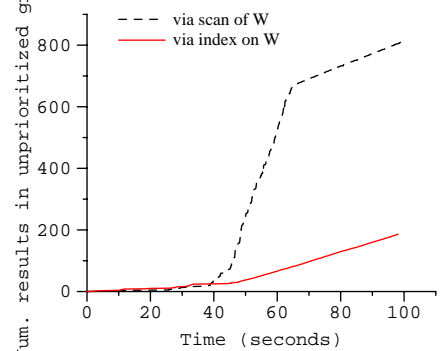
(ii)

Figure 11: Number of full and partial result tuples output over time for query Q3: (i) with SteMs, and (ii) with traditional join operators

Figure 12: Number of tuples output for Q4 using index join, hash join, and the hybrid approach during: (i) first 30 seconds, and (ii) first 200 seconds



(i)



(ii)

Figure 13: Number of (i) prioritized and (ii) unprioritized results generated by the Eddy for Q5, in the two ways

as well, for the partial RT result tuples – notice the sublinearity of the RTU curve of Figure 11(ii) compared to that of Figure 11(i).

4.5 Index/Hash join hybridization based on costs

Our next experiments studies the ability of our system to choose and hybridize among alternative join algorithms based on their costs.

We use query Q4, which joins R with a table W that has both an index and scan access method. To ensure that our results are not affected by cache effects, we use an equijoin between the key columns of R and W . This means that there are two natural ways of joining R and W : using the scans on R and W in a symmetric hash join, and using the scan on R and the index on W in an index join. A third way is for the Eddy to use both access methods on W , with SteMs on R and W , and choose a hybrid join algorithm.

Figure 12 (i) plots the number of result tuples generated over time in all these three approaches, during the first few seconds of the query execution. We see that the index join initially outperforms the hash join. This happens because the W index outputs the exact matching W tuple for each R probe tuple, whereas the W scan outputs all W tuples in an arbitrary order – only some of the R probes will find matches in the tuples scanned from W . The symmetric hash join however catches up with the index join quickly, as the R and W hash tables are filled. Figure 12 (ii) plots the same graph over the entire query execution period.

In this overall analysis, the hash join beats the index join handily because the scan on W is a faster access method than the index on W ⁶.

As we can see, the approach using SteMs performs well in both these plots. In the early stages, it performs much like the index join. Most of the R tuples are routed to the W index because the fanout of probes into $SteM_W$ is very low. But as the W tuples scan into $SteM_W$, most of the R tuples find matches in $SteM_W$ itself. The overall completion time of the hybrid approach is slightly more than that of the hash join, because the Eddy keeps sending a small fraction of the R tuples to probe into the W index throughout the processing (this is an artifact of the Eddy routing policy that we use which continues to explore alternative approaches for executing the query).

4.6 Index/Hash join hybridization based on user preferences

In our final experiment, we investigate join algorithm hybridization when the user has given preferences for different kinds of tuples. We will demonstrate the advantages of SteMs in interactive processing environments where the user gives preferences for various rows in the result, while the query is running (as in the Control project [H⁺99, RRH00]).

We modify query Q4 to form Q5 by adding a GROUP BY clause on $W.key \% 20$ (% is the sign for modulus in Telegraph). We model a scenario where the user has prioritized tuples in the group with $W.key \% 20 = 0$ alone.

Figure 13 (i) plots the number of results generated over time in the prioritized group using the hybrid approach from the previous section. Notice that very few of these prioritized results are generated by W tuples coming out of scans, because the prioritized tuples occur infrequently (5% of the tuples from the base tables are prioritized). Most of these results are generated by prioritized R tuples probing into the index on W . On the other hand, the scan on W contributes most of the *unprioritized* output tuples as can be seen in Figure 13 (ii). The number of unprioritized results generated in either way is low until all the prioritized results have been generated, because the Eddy prioritizes the routing of prioritized tuples.

This differential behavior arises for the following reason. The Eddy can only occasionally send tuples to probe into the W index because the throughput of these probes is low. Whenever it can send a probe to the W index, the Eddy always prefers to send prioritized tuples because the resulting matches will have higher benefit to the user⁷. In contrast, probes into the W SteM have high throughput and so the Eddy can send probe tuples to this SteM frequently. Prioritized tuples are much rarer than unprioritized ones, so most of the tuples probing into the SteMs are unprioritized tuples.

Thus the Eddy has combined the hash join and index join algorithms according to user interests. Using the index to get W tuples with $W.key$ divisible by 20 is like using the Index Stride reordering algorithm of [HHW97] to prioritize them. Whereas, using the scan on W to get the other W tuples is like using the Juggle reordering algorithm [RRH00]. In effect, the Eddy decides which reordering algorithm is appropriate based on the user preferences. Neither the Juggle nor the Index Stride algorithm alone is effective for this

⁶Both the hash join curve and the hybrid curve are quadratic until about 59 seconds because the R and W tuples are both being scanned in. At this point the scan from R stops, so the curves becomes linear, with a reduced slope.

⁷**Note:** Both R tuples and W tuples with $R.key$ (or $W.key$) divisible by 20 are prioritized even though the group by is only on $R.key$, because the Eddy automatically uses transitivity among the query's equality predicates to infer tuple priorities.

query. When Index Stride is used, the overall performance is poor, and when Juggle is used the W tuples with $W.key$ divisible by 20 are not prioritized well. A hybridization of Juggle and Index Stride is the best approach, and arises naturally when SteMs are used.

5 Related Work

There has long been interest in adapting query optimization decisions on the fly. Due to space constraints, we only discuss the most relevant work here – for detailed surveys, see [H⁺00] or [Ram01].

There has been much work on adapting join and selection ordering. In early work, the Ingres query processor [SWK76] did not use a query optimizer at all, but instead used a heuristic approach where each tuple could be routed through a collection of index (or nested-loops) joins in a different order. Graefe and Cole [GC94, I⁺92] study ways of optimizing queries in a parametrized fashion, so that the actual execution plan can be chosen just before execution. Other recent work allows the operator ordering to be changed even after a query has commenced execution. Kabra and DeWitt [K⁺98] reoptimize a running query at every block in the query plan. Tukwila [I⁺99] uses a similar approach, and has a rule-based language for specifying when reoptimization should occur. Query Scrambling [UFA98] reoptimizes queries running over a WAN whenever there is a delay in accessing a source.

There have also been attempts to develop join algorithms that can internally adapt to some changing properties. The Ripple Join [HH99] adapts to changing statistical properties of the data, to optimize for user feedback in online aggregation. The XJoin [UF00] is a variant of the symmetric hash join that dynamically changes its execution strategy to work with previously scanned tuples whenever there is a delay in one of its inputs. There has also been some work on making hash join and sort operators adaptive to memory fluctuations [P⁺93b, P⁺93a, ZL97]. The DEC (now Oracle) RDB system introduced a strategy of running multiple alternative access methods simultaneously for a short while and then stopping all but the best access method [AZ96].

We depart from this prior work in two important aspects. First, we adapt query execution at a very fine, per-tuple granularity. Second, while prior work focuses primarily on adapting join orders, our architecture allows much greater flexibility in adaptation, including the choice of access paths to data sources, join algorithms, join spanning trees, and join orders.

SteMs were developed as a part of the Telegraph project at Berkeley, and build on the Eddy tuple routing operator of [AH00]. The first (and complete) presentation of SteMs occurs in [Ram01]; this paper is meant to be a concise description. Since then, SteMs have been used in other Telegraph work on continuous query processing, with a focus on sharing SteMs across queries [M⁺02, CF02].

6 Conclusions and Future Work

Join operators constitute an important part of traditional query processors. These operators typically encapsulate complex algorithms that maintain much state about the tables involved in the join. The routing of a tuple to a join often results in a chain of steps within the join operator, that constitute multiple physical operations. In this paper, we have developed a way of executing queries by routing tuples not through join

operators but instead through State Modules that encapsulate data structures for holding intermediate query processing state. With this mechanism, most of the decisions involved in query optimization, including the ordering of joins and selections, the choice of access methods on the tables, the choice of reordering mechanism, the choice of join algorithms, and the choice of join spanning tree are determined by the routing of tuples, and are thus made dynamically by the Eddy. We have designed a set of restrictions on the Eddy's routing policy that ensure correct query execution. Our experiments demonstrate that the SteMs mechanism allows powerful adaptation by the Eddy in various situations.

We plan to extend this work in several directions. An important restriction of this paper is that it does not consider SteMs that span multiple tables. Though this reduces memory overheads, it can be inefficient in more traditional query execution scenarios as it leads to repeated probes that can be avoided by storing intermediate results. We are currently investigating extensions of our architecture that allow intermediate results, while retaining the adaptivity that SteMs provide.

Since SteMs encapsulate the data structure, and communicate directly with the Eddy, they enable the Eddy to observe and control memory resource utilization across *all* modules in the query. The Eddy can make memory allocation decisions in a globally optimal manner, possibly based on overall memory availability as well as relative frequency of probes into each SteM. This can be extended to let the Eddy control spilling of tuples to the disk as well. It will be interesting to see if such adaptive control of spilling can help the Eddy simulate join algorithms such as the XJoin [UF00] algorithm that dynamically adapt disk spilling. In presence of multi-table SteMs, this opens up a new set of optimization opportunities, where the Eddy can dynamically decide whether to materialize intermediate results or not based on memory availability.

Another important research direction is to formally study the space of join processing strategies opened up by the decoupling of state management and routing logic. We believe this will lead to better adaptive routing policies for learning many kinds of hybrid join strategies, which may be appropriate in particular circumstances but are not common enough to justify programming new join operators.

References

- [AH00] R. Avnur and J. M. Hellerstein. Eddies: Continuously adaptive query processing. In *SIGMOD*, 2000.
- [Alb91] J. Albert. Algebraic properties of bag data types. In *VLDB*, 1991.
- [AZ96] G. Antoshnekov and M. Ziauddin. Query processing and optimization in Oracle Rdb. *VLDB Journal*, 5(4), 1996.
- [BBD⁺02] B. Babcock, S. Babu, M. Datar, R. Motwani, and J. Widom. Models and issues in data stream systems. In *PODS*, 2002.
- [BC02] N. Bruno and S. Chaudhuri. Exploiting statistics on query expressions for optimization. In *SIGMOD*, 2002.
- [C⁺02] D. Carney et al. Monitoring streams: A new class of data management applications. In *VLDB*, 2002.

- [CF02] S. Chandrasekaran and M. J. Franklin. Streaming queries over streaming data. In *VLDB*, 2002. to appear.
- [D⁺84] D. J. DeWitt et al. Implementation techniques for main memory database systems. In *SIGMOD*, 1984.
- [FKT86] S. Fushimi, M. Kitsuregawa, and H. Tanaka. An overview of the system software of a parallel relational database machine GRACE. In *VLDB*, 1986.
- [GC94] G. Graefe and R. Cole. Optimization of dynamic query evaluation plans. In *SIGMOD*, 1994.
- [GW00] R. Goldman and J. Widom. WSQ/DSQ: a practical approach for combined querying of databases and the web. In *SIGMOD*, 2000.
- [H⁺97] L. M. Haas et al. Optimizing queries across diverse data sources. In *VLDB*, 1997.
- [H⁺99] J. M. Hellerstein et al. Interactive data analysis: The Control project. *IEEE Computer*, 32(8), 1999.
- [H⁺00] J. M. Hellerstein et al. Adaptive query processing: technology in evolution. *IEEE Data Engg. Bull.*, 23(2), 2000.
- [HH99] P. J. Haas and J. M. Hellerstein. Ripple joins for Online Aggregation. In *SIGMOD*, 1999.
- [HHW97] J. M. Hellerstein, P. J. Haas, and Helen J. Wang. Online aggregation. In *SIGMOD*, 1997.
- [I⁺92] Y. E. Ioannidis et al. Parametric query optimization. In *VLDB*, 1992.
- [I⁺99] Z. G. Ives et al. An adaptive query execution system for data integration. In *SIGMOD*, 1999.
- [IK84] T. Ibaraki and T. Kameda. Optimal nesting for computing N-relational joins. *TODS*, 9(3), 1984.
- [K⁺98] N. Kabra et al. Efficient mid-query reoptimization of sub-optimal query execution plans. In *SIGMOD*, 1998.
- [KBZ86] R. Krishnamurthy, H. Boral, and C. Zaniolo. Optimization of nonrecursive queries. In *VLDB*, 1986.
- [Knu73] D. E. Knuth. *Sorting and Searching*, volume 3 of *The Art of Computer Programming*. Addison-Wesley, 1973.
- [Lex] BrightPlanet LexiBot. www.brightplanet.com.
- [M⁺02] S. Madden et al. Continuously adaptive continuous queries over streams. In *SIGMOD*, 2002.
- [Mor88] K. A. Morris. An algorithm for ordering subgoals in NAIL! In *PODS*, 1988.
- [P⁺92] H. Pirahesh et al. Extensible/rule-based query rewrite optimization in Starburst. In *SIGMOD*, 1992.
- [P⁺93a] H. Pang et al. Memory-adaptive external sorting. In *VLDB*, 1993.
- [P⁺93b] H. Pang et al. Partially preemptive hash joins. In *SIGMOD*, 1993.
- [Ram01] V. Raman. *Interactive Query Processing*. PhD thesis, U.C.Berkeley, 2001.
- [RGM01] S. Raghavan and H. Garcia-Molina. Crawling the hidden web. In *VLDB*, 2001.
- [RH02] V. Raman and J.M. Hellerstein. Partial results for online query processing. In *SIGMOD*, 2002.
- [RRH00] V. Raman, B. Raman, and J. M. Hellerstein. Online dynamic reordering. *VLDB Journal*, 9(3), 2000.
- [RS86] L. Rashid and S. Su. A parallel processing strategy for evaluating recursive queries. In *VLDB*,

1986.

- [S⁺01] M. A. Shah et al. Java support for data-intensive systems. *SIGMOD Record*, 4(30), 2001.
- [SLMK01] M. Stillger, G. M. Lohman, V. Markl, and M. Kandil. LEO – DB2’s LEarning Optimizer. In *VLDB*, 2001.
- [SWK76] M.R. Stonebraker, E. Wong, and P. Kreps. The design and implementation of INGRES. *TODS*, 1(3):189–222, September 1976.
- [Tel] The Telegraph project. <http://db.cs.berkeley.edu/telegraph>.
- [UF00] T. Urhan and M. J. Franklin. XJoin: A Reactively-Scheduled Pipelined Join Operator. *IEEE Data Engineering Bulletin*, 23(2), 2000.
- [UF01] T. Urhan and M. J. Franklin. Dynamic pipeline scheduling for improving interactive query performance. In *VLDB*, 2001.
- [UFA98] T. Urhan, M. J. Franklin, and L. Amsaleg. Cost-based query scrambling for initial delays. In *SIGMOD*, 1998.
- [VN02] S. Viglas and J. Naughton. Rate-based query optimization for streaming information sources. In *SIGMOD*, 2002.
- [WA91] A. N. Wilschut and P. M. G. Apers. Dataflow query execution in a parallel main-memory environment. In *PDIS*, 1991.
- [ZL97] W. Zhang and P. Larson. Dynamic memory adjustment for external mergesort. In *VLDB*, 1997.
- [ZR02] V. Zadorozhny and L. Raschid. Query optimization to meet performance targets for wide area applications. In *ICDCS*, 2002.

A Proof of Correctness of Routing Constraints

We first show that the Eddy will not output duplicate query results, or even duplicate partial query results.

We need some definitions :

Definition 4 (SubTuple and SuperTuple) We define a tuple s to be a subtuple of a tuple t if its base-table components form a subset of the base-table components of t , i.e., $s = \langle s_1, s_2, \dots, s_k \rangle$ is a subtuple of $t = \langle t_1, t_2, \dots, t_m \rangle$ if there exist distinct base-tables c_1, c_2, \dots, c_k , $1 \leq c_1 < c_2 < \dots < c_k \leq m$ such that $s_i = t_{c_i} \forall 1 \leq i \leq k$. Conversely, we call t a supertuple of s .

Definition 5 (T-Arity of a tuple) The t -arity of a tuple is the number of base-table components it has.

Definition 6 (Generator) A tuple in the query dataflow is called a generator if it is a subtuple of a query result tuple.

Theorem 3 If an Eddy follows a routing policy that satisfies the constraints of Table 2, there will not arise duplicate versions of any tuple in the dataflow, other than singleton tuples not yet been built into SteMs.

Proof: We proceed by contradiction. Consider a query over n sources s_1, s_2, \dots, s_n , with a query result QR. Let t be the tuple of minimum t -arity for which duplicate versions occur in the dataflow during query execution. Call two of these versions t_a and t_b . Suppose, without loss of generality, that $t = \langle t_{s_1}, t_{s_2}, \dots, t_{s_k} \rangle$.

If t is a singleton tuple. t_a and t_b cannot have been built into $SteM_{s_1}$, because then only one of the two would have been bounced back after building. So we only consider the situation where $k > 1$.

As discussed before, any tuple is created during query execution by a singleton tuple probing into various SteMs and accumulating components from other base-tables. Arrange the base-table components of t_a and t_b in the order in which they were accumulated; $t_a = \langle t_{s_{a_1}}, t_{s_{a_2}}, \dots, t_{s_{a_k}} \rangle$ and $t_b = \langle t_{s_{b_1}}, t_{s_{b_2}}, \dots, t_{s_{b_k}} \rangle$, where a_1, a_2, \dots, a_k and b_1, b_2, \dots, b_k are both permutations of $1, 2, \dots, k$.

By the TimeStamp constraint, $\langle t_{s_{a_1}} \rangle$ and $\langle t_{s_{b_1}} \rangle$ must both have the highest timestamp among the base-table components of t , and so $a_1 = b_1$ (i.e., t_a and t_b were generated by $\langle t_{s_{a_1}} \rangle$ probing into other SteMs). Let $t_{s_{a_1}} \dots t_{s_{a_l}}$, $l > 0$, be the longest common prefix of $\langle t_{s_{a_1}}, t_{s_{a_2}}, \dots, t_{s_{a_k}} \rangle$ and $\langle t_{s_{b_1}}, t_{s_{b_2}}, \dots, t_{s_{b_k}} \rangle$. Two cases arise:

Case 1 $l = k$: We know $k = l > 1$. Hence t_a and t_b are both generated by $\langle t_{s_{a_1}} \dots t_{s_{a_{l-1}}} \rangle$ probing into $SteM_{s_{a_l}}$. $SteM_{s_{a_l}}$ cannot have duplicate versions of $t_{s_{a_l}}$. Even if $\langle t_{s_{a_1}} \dots t_{s_{a_{l-1}}} \rangle$ probes into $SteM_{s_{a_l}}$ multiple times, the TimeStamp constraint ensures that it can match with $t_{s_{a_l}}$ only once. So there must have been duplicate versions of $\langle t_{s_{a_1}} \dots t_{s_{a_{l-1}}} \rangle$ itself. This contradicts our hypothesis that t is a duplicate of minimum t -arity.

Case 2 $l < k$: The generation of t_a involves $\langle t_{s_{a_1}} \dots t_{s_{a_l}} \rangle$ probing into $SteM_{s_{a_{l+1}}}$ and the generation of t_b involves $\langle t_{s_{a_1}} \dots t_{s_{a_l}} \rangle$ probing into $SteM_{s_{b_{l+1}}}$. The ProbeCompletion constraint ensures that $\langle t_{s_{a_1}} \dots t_{s_{a_l}} \rangle$ can only probe into a single SteM. We know that $a_{l+1} \neq b_{l+1}$, because we have chosen the longest common prefix. Therefore there must have been duplicate versions of $\langle t_{s_{a_1}} \dots t_{s_{a_l}} \rangle$. This again contradicts our hypothesis that t is a duplicate of minimum t -arity. \square

Theorem 4 *If the Eddy follows a routing policy that satisfies the constraints of Table 2, only a finite number of tuples will arise in the query dataflow over the course of query execution.*

Proof: Observe that since we are dealing with finite relations, and because of Theorem 3, only a finite number of tuples can arise in the query dataflow, other than singleton tuples that have not yet been built into SteMs. If the latter set were infinite, the set of tuples coming out of some AM must be infinite. Clearly this must be an index AM, say AM_R (all tables are finite, so scans can only generate a finite number of tuples). For an infinite number of tuples to come out of AM_R , an infinite number of tuples must probe into AM_R . Again, because of Theorem 3, these probe tuples must all be singleton tuples. Therefore there must exist some query relation $S \neq R$, such that there are an infinite number of S tuples probing into R . But the scan AMs on S can only generate a finite number of tuples. Also, tuples coming out of the index AMs on S can only probe into AM_R after they have built into $SteM_S$. But $SteM_S$ will only bounce back distinct S tuples to probe into R , so there can only be a finite number of probe tuples generated this way. \square

Theorem 5 *If the Eddy follows a routing policy that satisfies the constraints of Table 2, it will not output any tuple that is not in the query result, and will output all query result tuples in a finite number of routing steps.*

Proof: Consider a query over n sources s_1, s_2, \dots, s_n , with a query result QR.

Notice first that the Eddy will not output any tuple not in QR as a query result tuple, because it will first route the tuple to all selection modules. (Section 2.1.1).

Consider the set of all tuples that arise in the dataflow during query execution, and order them in increasing order of the *latest* time when they were routed by the Eddy (the same tuple might be routed by the Eddy

multiple times if it gets bounced back from modules). From Theorem 4, this sequence has no duplicates except for singleton tuples. For each singleton tuple that occurs in this sequence, remove all occurrences except the first (the rest will be absorbed by the corresponding SteM). Let the resulting sequence of tuples be $G = \{g_1, g_2, \dots, g_m\}$.

Consider a query result tuple $p = \langle p_{s_1} p_{s_2} \dots p_{s_n} \rangle \in QR$. Clearly, at least one of the sources, say s_1 , must be scannable, for the query to be executable. Therefore there must exist a generator for p in G – when the scan AM on s_1 emits $\langle p_{s_1} \rangle$.

Among all the generators of p in G , choose the ones with highest Timestamp, and among these the ones with maximum t -arity. Let g be one of these generators. We will now show that $g = p$.

We proceed by contradiction. Suppose that g is not equal to p . Clearly, $g \notin QR$. Notice that g must eventually be routed by the Eddy, because there are only a finite number of tuples ever in the dataflow (Theorem 4), and query execution will not terminate unless all these tuples are removed. Consider the ways in which g can be routed by the Eddy.

Case 1: g is routed to a SM. Since g is a subtuple of p , it will satisfy the selection predicate and be bounced back to the Eddy, and we would have selected this later occurrence of g when constructing G . Hence we have a contradiction.

Case 2: g is routed to build into a SteM. If g is bounced back, the same argument as Case 1 holds. If g is absorbed, then a duplicate g' of g has built into the SteM earlier, and we would have selected g' over g when constructing the sequence G . Hence we have a contradiction.

Case 3: g is routed to probe into a SteM, say $SteM_{s_w}$. Four sub-cases arise:

Case 3a: $p_{s_w} \in SteM_{s_w}$ and $TS(p_{s_w}) < TS(g)$. Then the concatenation of g and p_{s_w} will arise in the dataflow: if $LastMatchTS(g) < TS(p_{s_w})$, $SteM_{s_w}$ will output the concatenation on this probe by g ; otherwise the SteM has already output the concatenation on a previous probe by g . This concatenation is a generator with higher t -arity than and the same timestamp as g , contradicting our scheme for choosing generators.

Case 3b: $p_{s_w} \in SteM_{s_w}$ and $TS(p_{s_w}) > TS(g)$. Then $\langle p_{s_w} \rangle$ would have entered the dataflow after building into $SteM_{s_w}$. It is a generator for p with higher timestamp than g , contradicting our scheme for choosing generators.

Case 3c: $p_{s_w} \notin SteM_{s_w}$ and there exists a scan AM on $SteM_{s_w}$. Then $\langle p_{s_w} \rangle$ will enter the dataflow and get built into $SteM_{s_w}$ at some later time. This $\langle p_{s_w} \rangle$ will be a generator with higher timestamp than g , contradicting our scheme for choosing generators.

Case 3d: $p_{s_w} \notin SteM_{s_w}$ and there does not exist a scan AM on $SteM_{s_w}$. g will be bounced back, and we would have selected this later occurrence of g when constructing G . Hence we have a contradiction.

Case 4: g is routed to an AM, say AM_{s_w} . Two subcases arise:

Case 4a: g is not a prior prober, or g is a prior prober with a different probe completion AM than AM_{s_w} . g will be bounced back, and we would have selected this later occurrence of g for G . Hence we have a contradiction.

Case 4b: g is a prior prober and AM_{s_w} is its probe completion AM. g must have probed into $SteM_{s_w}$ at some time, to have been transformed into a prior prober. We know that p_{s_w} was not in $SteM_{s_w}$ at that time, because then the returned match will have a higher t -arity than g (as in Case 3a). Whenever p_{s_w} enters

the dataflow (it must enter at some time, at the latest as an asynchronous match from the $SteM_{sw}$) and builds into $SteM_{sw}$, it will be bounced back with a higher timestamp than g . This again contradicts our scheme for choosing generators.

Thus, $g = p$, and will be generated and output by the Eddy. Since p was chosen as some arbitrary result tuple, all result tuples will be generated by the Eddy. \square