

Processing Sliding Window Multi-Joins in Continuous Queries over Data Streams*

Lukasz Golab

M. Tamer Özsu

School of Computer Science
University of Waterloo
Waterloo, Canada
{lgolab, tozsu}@uwaterloo.ca

Abstract

We study sliding window multi-join processing in continuous queries over data streams. Several algorithms are reported for performing continuous, incremental joins, under the assumption that all the sliding windows fit in main memory. The algorithms include multi-way incremental nested loop joins (NLJs) and multi-way incremental hash joins. We also propose join ordering heuristics to minimize the processing cost per unit time. We test a possible implementation of these algorithms and show that, as expected, hash joins are faster than NLJs for performing equi-joins, and that the overall processing cost is influenced by the strategies used to remove expired tuples from the sliding windows.

1 Introduction

A data stream is a real-time, continuous, ordered (explicitly by timestamp or implicitly by arrival time) sequence of items. Applications where information naturally occurs as a stream of data values include sensor data processing [7, 32], Internet traffic analysis [13, 24], financial tickers [9, 33], and analysis of various transaction logs such as Web server logs and telephone call records [10]. Due to their continuous and dynamic nature, querying data streams involves running a query continually over a period of time and generating new

answers as new items arrive. First proposed in [25], these types of queries are known in the literature as *continuous*, *standing*, or *persistent* queries [9, 21].

Several issues arise in on-line stream processing. Firstly, unbounded streams may not be wholly stored in bounded memory. Secondly, because data streams are temporally ordered, new items are often more accurate or more relevant than older items. Finally, streaming query plans may not use blocking operators that must consume the entire input before any results are produced.

A common solution to these issues is to restrict the range of continuous queries to a sliding window that contains the last T items or those items that have arrived in the last t time units. The former is called a *count-based*, or a *sequence-based* sliding window, while the latter is called a *time-based* or a *timestamp-based* sliding window [5]. Constraining all queries by sliding window predicates allows continuous queries over unbounded streams to be executed in finite memory and in an incremental manner by generating new results as new items arrive. In particular, while joining two or more infinite streams is seldom feasible (see Arasu et al. for a discussion of memory requirements of streaming queries [2]), evaluating windowed joins over stream excerpts is practical and useful in many applications. For example, an Internet traffic engineer may pose a query that joins traffic traces from various links in a network, with equality comparison of the source and destination IP addresses of each packet header as the join predicate. This windowed join may be used to trace packets through the network and identify sessions whose packets follow different paths to reach the same destination. The latter could be used to study load balancing in a network with many redundant links.

Processing continuous queries over sliding windows introduces two issues that affect the design of windowed algorithms: re-execution strategies and tuple invalidation procedures. An *eager re-evaluation* strategy generates new results after each new tuple arrives,

*This research is partially supported by the Natural Sciences and Engineering Research Council (NSERC) of Canada.

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 VLDB copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Very Large Data Base Endowment. To copy otherwise, or to republish, requires a fee and/or special permission from the Endowment.

but may be infeasible in situations where streams have high arrival rates. A more practical solution—*lazy re-evaluation*—is to re-execute the query periodically. The downside of lazy re-evaluation is an increased delay between the arrival of a new tuple and the generation of new results based on this tuple—long delays may be unacceptable in streaming applications that must react quickly to unusual patterns in the data. Similarly, *eager expiration* proceeds by scanning the sliding windows and removing old tuples upon arrival of each new tuple (this is trivial in count-based windows as each new tuple simply replaces the oldest tuple in its window). In contrast, *lazy expiration* involves removing old tuples periodically and requires more memory to store tuples waiting to be expired.

1.1 Problem Statement

Given n data streams and n corresponding sliding windows, our goal in this work is to continually evaluate the exact join of all n windows. We assume that each stream consists of relational tuples with the following schema: a timestamp attribute ts and an attribute $attr$ containing values from an ordered domain. We also assume that all windows fit in main memory and we require that all query plans use extreme right-deep join trees that do not materialize any intermediate results. Furthermore, we do not permit time-lagged windows, i.e. all windows are assumed to start at the current time and expire at the current time minus the corresponding window size.

We define the semantics of sliding window joins as monotonic queries over append-only relations; therefore once a result tuple is produced, it remains in the answer set indefinitely. Hence, new results may be streamed directly to the user. In particular, for each newly arrived tuple k , the join is to probe all tuples present in the sliding windows precisely at the time of k 's arrival (i.e. those tuples which have not expired at time equal to k 's timestamp), and return all (composite) tuples that satisfy every join predicate. Moreover, we impose a time bound τ on the time interval from the arrival of a new tuple k until the time when all join results involving k and all tuples older than k have been streamed to the user. This means that a query must be re-executed at least every τ time units. However, it does not mean that all join results containing k will be available at most τ units after k 's arrival; k will remain in the window until it expires, during which time it may join with other (newer) tuples.

1.2 Summary of Contributions

In this paper, we present a solution to the multi-join processing problem over sliding windows residing in main memory. Our algorithms execute multiple joins together in a series of nested for-loops and process newly arrived tuples from each window separately (possibly using different join orders). We first propose

λ_i	Arrival rate of stream i in tuples per unit time
S_j	Sliding window corresponding to stream j
T_j	Time size of the j^{th} time-based window
C_j	Number of tuples in S_j
v_j	Number of distinct values in S_j
b_j	Number of hash buckets in the hash index of S_j , if such an index exists
τ	Continuous query re-execution interval
$a \circ b$	Concatenation of tuples a and b
θ	Arithmetic comparison predicate, e.g. =

Table 1: Explanations of symbols used in this paper.

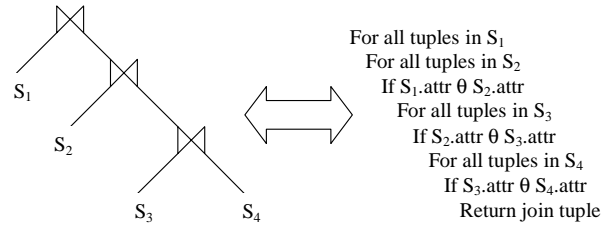


Figure 1: Join order $S_1 \bowtie (S_2 \bowtie (S_3 \bowtie S_4))$ expressed as a join tree (left) and as a series of for-loops (right).

a multi-way NLJ for eager query re-evaluation that always iterates over newly arrived tuples in the outer for-loop. We demonstrate that this strategy may not work well for lazy re-evaluation and propose another NLJ-based operator that is guaranteed to perform as well or better than the initial algorithm. We also extend our NLJ-based algorithms to work with hash indices on some or all sliding windows. Moreover, we propose join ordering heuristics for our algorithms that attempt to minimize the number of intermediate tuples that are passed down to the inner for-loops. These heuristics are based on a main-memory per-unit-time cost model. Finally, we test an implementation of our algorithms and investigate the effects of re-evaluation and expiration strategies on the overall processing cost.

Table 1 lists the symbols used in this paper and their meanings. In Figure 1, we give an example to explain our convention for describing join ordering. In this example, the join order $S_1 \bowtie (S_2 \bowtie (S_3 \bowtie S_4))$ is expressed as a join tree on the left and as a series of nested for-loops on the right. We refer to S_1 as being “on top of the plan” or “ordered first”, S_2 as “ordered second” and so on. For brevity, we may omit parentheses and write S_1, S_2, S_3, S_4 to represent the join order shown in Figure 1. Note that the pseudocode on the right of Figure 1 shows the general strategy used in our multi-way joins, whereby the join is effectively evaluated “from the top down”; the join predicate is evaluated inside each for-loop in order to minimize the number of tuples that are passed down to the inner for-loops.

1.3 Roadmap

The remainder of this paper is organized as follows. In Section 2, we review related work. Section 3 defines our multi-way incremental join algorithms, whose processing cost is analyzed in Section 4. Section 5 outlines join ordering heuristics. We discuss experimental results in Section 6. Section 7 concludes the paper and outlines directions for future research.

2 Related Work

As surveyed in [4, 15], there has been a great deal of recent interest in developing novel data management techniques and adapting traditional database technology to the data stream model; Cougar [32], Aurora [7], TelegraphCQ [8, 22], and STREAM [4, 23] are some examples. The first two focus on processing sensor data. Cougar concentrates on distributed query processing inside the sensor network, while Aurora allows users to create query plans by visually arranging query operators using a boxes (corresponding to query operators) and links (corresponding to data flow) paradigm. TelegraphCQ is designed for adaptive query processing, where query plans are re-ordered throughout the lifetime of a continuous query in response to changes in the execution environment (e.g. fluctuating stream arrival rates). STREAM addresses all aspects of stream data management, including memory management, operator scheduling, and approximate query answering via summary information. A continuous query language (CQL) has also been proposed within the STREAM project [3].

Windowing is one solution proposed in the literature for bounding the memory requirements of continuous queries and unblocking streaming operators. Another alternative is to maintain compact stream summaries and provide approximate answers based on the summaries. Many summary structures have been proposed in the literature, e.g. [1, 12, 14, 17]. In some cases, even a sliding window is too large to fit in memory and both windowing and summary structures are necessary. For instance, Datar et al. [11] give an approximate algorithm for bit counting in a sliding window. The third method is to exploit any constraints that may hold in a data stream. For example, Babu and Widom [6] show that foreign key constraints and ordered or clustered arrival (i.e. stream items arrive in some known order, or duplicate items arrive in one contiguous batch) may be exploited to lower the memory usage in continuous query processing. Moreover, assertions, referred to as *punctuations* in [26], could be inserted into a stream to specify a restriction on subsequently arriving items. For instance, a punctuation may arrive stating that all future items shall have the A attribute larger than ten. This punctuation could be used to partially unblock a group-by query on A —all those groups whose A attribute is larger than ten

are guaranteed not to change.

Relevant work on join processing over unbounded streams includes non-blocking binary join algorithms such as the XJoin [27], which is a variant of the Symmetric Hash Join [30], and the Ripple Join [18]. Viglas et al. [29] have developed a multi-way version of the XJoin called the MJoin. Moreover, Viglas and Naughton [28] propose a rate-based query optimization model for continuous queries over data streams; if the stream input rates change over time, then so will the output rate of a join.

The Symmetric Hash Join has been extended to a binary sliding window join by Kang et al. [20], who also introduce a per-unit-time cost model for insertion, invalidation, and join processing costs. Hammad et al. deal with multi-way windowed joins in [19], but their focus is on scheduling which tuples to process when two or more continuous queries are interested in sliding windows of various sizes over the join result.

While a multi-way hash join has been proposed in the literature, to the best of our knowledge this work is the first to consider multi-way joins designed explicitly for sliding windows. We also know of no previous work on the join ordering problem in the context of sliding windows over data streams, although this problem is identified in the context of optimizing for the highest output rate of queries over infinite streams [28, 29]. Generally, main-memory join ordering techniques focus on pushing expensive predicates to the top of the plan (see, e.g. [31]).

3 Sliding Window Join Algorithms

3.1 Motivation

We begin the discussion of sliding window multi-way joins with examples that motivate our strategies, starting with the simplest case of eager re-evaluation and expiration. We initially concentrate on time-based windows and defer extensions of our algorithms to count-based windows to Section 3.6.

A binary incremental NLJ has been proposed by Kang et al. [20] and proceeds as follows. Let S_1 and S_2 be two sliding windows to be joined. For each newly arrived S_1 -tuple, we scan S_2 and return all matching tuples. We then insert the new tuple into S_1 and invalidate expired tuples. We follow the same procedure for each newly arrived S_2 -tuple. Extending the binary NLJ to more than two windows is straightforward: for each newly arrived tuple k , we execute the join sequence in the order prescribed by the query plan, but we only include k in the join process (not the entire window that contains k). For example, suppose that we wish to join three windows, S_1 , S_2 , and S_3 , using the plan $S_1 \bowtie (S_2 \bowtie S_3)$. Upon arrival of a new S_1 -tuple, we invalidate expired tuples in S_2 and S_3 and then probe all tuples in $S_2 \bowtie S_3$. If a new S_2 -tuple arrives, we first expire old tuples from S_1 and S_3 , then

S_1		S_2		S_3	
ts	$attr$	ts	$attr$	ts	$attr$
90	1	150	1	195	1
100	1	180	1	205	1

Table 2: Partial contents of sliding windows S_1 , S_2 , and S_3 .

Incorrect Results					
ts	$attr$	ts	$attr$	ts	$attr$
90	1	150	1	195	1
90	1	150	1	205	1
90	1	180	1	195	1
90	1	180	1	205	1
100	1	150	1	195	1
100	1	150	1	205	1
100	1	180	1	195	1
100	1	180	1	205	1

Correct Results					
ts	$attr$	ts	$attr$	ts	$attr$
100	1	150	1	195	1
100	1	180	1	195	1

Table 3: Partial results returned by the naive algorithm (top) and correct partial join results (bottom).

for each tuple currently in S_1 , we compute the join of the new tuple with S_3 and probe the result set. Similarly, upon arrival of a new S_3 -tuple, we expire S_1 and S_2 -tuples, then for each S_1 -tuple, we evaluate the join of S_2 with the new tuple and probe the result set. We call this algorithm the NAIVE MULTI-WAY NLJ.

To perform lazy re-evaluation with the naive algorithm, we re-execute the join every τ time units, first joining new S_1 -tuples with other sliding windows, then new S_2 -tuples, and so on. However, we must ensure that expired tuples are not included in the result, as illustrated by the following example. Suppose that we wish to process a batch of newly arrived tuples from S_3 given the join order $S_1 \bowtie (S_2 \bowtie S_3)$. For each tuple in S_1 , we re-compute the join $S_2 \bowtie S_3$, but using only the newly arrived tuples from S_3 . Assume that all windows are 100 seconds long and that their partial contents are as shown in Table 2. All join predicates are equalities on the attribute $attr$. If the current time is 210, algorithm NAIVE MULTI-WAY NLJ will return incorrect results when processing new S_3 -tuples. Table 3 shows the incremental result set returned for all new S_3 -tuples on the top and the actual result set that should be returned on the bottom. As can be seen, some tuples from S_1 should not join with any new S_3 -tuples as those S_1 -tuples should have expired by the time the new S_3 -tuples have arrived. Thus, we need to verify that a composite $S_2 \circ S_3$ tuple (call the S_2 -tuple j and the S_3 tuple k) joins with an S_1 tuple, call it i , only if $i.ts \geq k.ts - T_1$ and $j.ts \geq k.ts - T_2$. These two conditions ensure that i and j have not expired at time equal to k 's timestamp.

Origin of new tuple	Join order
S_1	$S_1 \bowtie (S_2 \bowtie S_3)$
S_2	$S_2 \bowtie (S_1 \bowtie S_3)$
S_3	$S_3 \bowtie (S_1 \bowtie S_2)$

Table 4: Join orderings used in our multi-way NLJ given that the global join ordering is $S_1 \bowtie (S_2 \bowtie S_3)$.

3.2 Improved Eager Multi-Way NLJ

In the above examples, when a new S_3 -tuple arrives we re-compute the join of S_2 with the new S_3 -tuple for each tuple in S_1 . Generally, this results in unnecessary work whenever new tuples arrive from a window that is not ordered first in the join tree. We propose a more efficient technique for handling new S_3 -tuples, in which we initially select only those tuples in S_1 which join with the new S_3 -tuple (suppose there are c such tuples), and make c scans of S_2 . In contrast, the previous approach requires a number of scans of S_2 equal to the number of tuples in S_1 . Unless all tuples in S_1 satisfy the join predicate (as in a Cartesian product), c is less than the size of S_1 . When a new S_1 -tuple or a new S_2 -tuple arrives, we proceed similarly by selecting those tuples from the window on top of the join order that join with the new tuple, and for each match, scanning the remaining window. In effect, the window at the top of the join order always consists of only one tuple and the join order changes in response to the origin of the incoming tuple (Table 4). This is possible because we have assumed a common join attribute across all streams. We define the global join order of a query plan as the static order that our algorithm would follow had it not ordered new tuples first. In the above example, the global order is S_1, S_2, S_3 .

The pseudocode of the algorithm is given below. Without loss of generality, we let the global join order be $S_1, S_2 \dots S_n$.

Algorithm EAGER MULTI-WAY NLJ

If a new tuple k arrives on stream i
 Insert new tuple in window S_i
 COMPUTEJOIN($k, (S_1, \dots, S_{i-1}, S_{i+1}, \dots, S_n)$)

Algorithm COMPUTEJOIN

Input: new tuple k from window S_i and a join order $(S_1, \dots, S_{i-1}, S_{i+1}, \dots, S_n)$.

$\forall u \in S_1$ and $k.ts - T_1 \leq u.ts \leq k.ts$

If $k.attr \theta u.attr$

... \ loop through S_2 up to S_{i-2}

$\forall v \in S_{i-1}$ and $k.ts - T_{i-1} \leq v.ts \leq k.ts$

If $k.attr \theta v.attr$

$\forall w \in S_{i+1}$ and $k.ts - T_{i+1} \leq w.ts \leq k.ts$

If $k.attr \theta w.attr$

... \ loop through S_{i+2} up to S_{n-1}

$\forall x \in S_n$ and $k.ts - T_n \leq x.ts \leq k.ts$

If $k.attr \theta x.attr$

Return $k \circ u \circ v \circ \dots \circ x$

Early projection may be performed on-the-fly by removing unwanted attributes in the Return statement, though this is not necessary in main-memory query processing. Note that algorithm MULTI-WAY NLJ may be used with eager or lazy expiration and is guaranteed not to join any expired tuples because of the timestamp comparisons done in each for-loop.

3.3 Lazy Multi-Way NLJs

A straightforward adaptation of algorithm EAGER MULTI-WAY NLJ to the lazy re-evaluation scenario is to process in the outer-most for-loop all the new tuples which have arrived since the last re-execution. Let NOW be the time at which re-execution begins. The pseudocode is given below.

Algorithm LAZY MULTI-WAY NLJ

Insert each new tuple into its window as it arrives
 Every time the query is to be re-executed

```

  For  $i = 1 \dots n$ 
     $\forall k \in S_i$  and  $NOW - \tau \leq k.ts \leq NOW$ 
      COMPUTEJOIN( $k, (S_1, \dots, S_{i-1}, S_{i+1}, \dots, S_n)$ )
  
```

Observe that the above algorithm can be made more general if newly arrived tuples are not restricted to the outer-most for-loop, leading to a general lazy NLJ algorithm that accepts arbitrary join orders. Let $O_i = O_{i,1}, \dots, O_{i,n}$ be the join order for performing lazy re-evaluation involving newly arrived tuples in window i . For instance, if $O_1 = S_1 \bowtie (S_3 \bowtie S_2)$ then $O_{1,1} = S_1$, $O_{1,2} = S_3$, and $O_{1,3} = S_2$. We also allow syntax of the form $\lambda_{i,j}$ and $T_{i,j}$ to represent the rate and time size of the window that is j^{th} in the join order when joining with new tuples from window i . Furthermore, let p be the position within the ordering of the window with the newly arrived tuples. In the preceding example, $p = 1$ because $O_{1,1} = S_1$. The general algorithm is as follows.

Algorithm GENERAL LAZY MULTI-WAY NLJ

Insert new tuples into windows as they arrive
 Every time the query is to be re-executed

```

  For  $i = 1 \dots n$ 
    GENERALCOMPUTEJOIN( $i, O_i$ )
  
```

Algorithm GENERALCOMPUTEJOIN

Input: window subscript i and a join order O_i

```

 $\forall u \in O_{i,1}$ 
 $\forall v \in O_{i,2}$ 
  If  $u.attr \theta v.attr$ 
    ... \ loop through  $O_{i,3}$  up to  $O_{i,p-1}$ 
     $\forall k \in O_{i,p}$  and  $NOW - \tau \leq k.ts \leq NOW$  and
     $k.ts - T_{i,1} \leq u.ts \leq k.ts$  and
     $k.ts - T_{i,2} \leq v.ts \leq k.ts$  and ...
    If  $u.attr \theta k.attr$ 
      ... \ loop through  $O_{i,p+1}$  up to  $O_{i,n-1}$ 
       $\forall x \in O_{i,n}$  and  $k.ts - T_{i,n} \leq x.ts \leq k.ts$ 
      If  $u.attr \theta x.attr$ 
        Return  $u \circ v \circ \dots \circ k \circ \dots \circ x$ 
  
```

Note that timestamp comparisons with new tuples from window S_k can only be done in S_k 's for-loop and below. This is why S_k 's for-loop must check that the composite tuple produced so far has not expired relative to a new tuple k . Algorithms LAZY MULTI-WAY NLJ and GENERAL MULTI-WAY NLJ may only be used with lazy expiration, which must be performed at most as frequently as the re-evaluation frequency.

3.4 Multi-Way Hash Joins

We proceed as in the NLJ, except that at each for-loop, we only scan one hash bucket instead of the entire sliding window. An eager version of the hash join is given below. We use the notation $B_{i,k}$ to represent the hash bucket in the i^{th} window to which attribute $attr$ of tuple k maps, i.e. $B(i, k) = h_i(k.attr)$ where h_i is the hash function used for the i^{th} window.

Algorithm MULTI-WAY HASH JOIN

If a new tuple k arrives on stream i

Insert new tuple in window S_i

COMPUTEHASHJOIN($k, (S_1, \dots, S_{i-1}, S_{i+1}, \dots, S_n)$)

Algorithm COMPUTEHASHJOIN

Input: new tuple k from window S_i and a join order $(S_1, \dots, S_{i-1}, S_{i+1}, \dots, S_n)$.

$\forall u \in B_{1,k}$ and $k.ts - \lambda_1 T_1 \leq u.ts \leq k.ts$

If $k.attr \theta u.attr$

... \ loop through $B_{2,k}$ up to $B_{i-2,k}$

$\forall v \in B_{i-1,k}$ and $k.ts - \lambda_{i-1} T_{i-1} \leq v.ts \leq k.ts$

If $k.attr \theta v.attr$

$\forall w \in B_{i+1,k}$ and $k.ts - \lambda_{i+1} T_{i+1} \leq w.ts \leq k.ts$

If $k.attr \theta w.attr$

... \ loop through $B_{i+2,k}$ up to $B_{n-1,k}$

$\forall x \in B_{n,k}$ and $k.ts - \lambda_n T_n \leq x.ts \leq k.ts$

If $k.attr \theta x.attr$

Return $k \circ u \circ v \circ w \circ \dots \circ x$

The lazy version of the multi-way hash join can be similarly obtained from algorithm LAZY MULTI-WAY NLJ or GENERAL LAZY MULTI-WAY NLJ. Because the sliding windows are assumed to fit in main memory, the multi-way NLJ may be considered a special case of the multi-way hash join with $b_i = 1$ for all i .

3.5 Hybrid NLJ-Hash Join

If a query requires an equi-join of a set of windows and a general theta-join of another set, all on the same attribute, then we may combine the join predicates into one hybrid multi-way join operator. We proceed in a nested-loops fashion as before, at each step scanning hash buckets if hash tables are available and scanning whole windows otherwise.

3.6 Extensions to Count-Based Windows

Our algorithms can be easily modified for use with count-based windows. Firstly, note that eager expiration in count-based windows is easy: if we implement

such windows (and hash buckets) as circular arrays, then we can perform insertion and invalidation in one step by overwriting the oldest tuple. Thus, algorithms MULTI-WAY NLJ and MULTI-WAY HASH JOIN carry over and in fact subroutines COMPUTEJOIN and COMPUTEHASHJOIN can be made simpler by omitting the timestamp comparisons. The more interesting case is that of lazy re-evaluation as eager expiration may not be performed in some applications. A possible solution is to maintain a circular counter and assign positions to each element in the sliding window (call them cnt). The join could then be performed as outlined in our algorithms with one exception: when probing for tuples to join with a new tuple k , instead of verifying that each tuple’s timestamp has not expired at time $k.ts$, we ensure that each tuple’s counter attribute cnt has not expired at time $k.ts$. To do this, for each sliding window we find the counter value of the element with the largest timestamp not exceeding $k.ts$, subtract the window length from this counter (call the counter value obtained so far tmp), and ensure that we join only those tuples whose counters are larger than tmp . The pseudocode for subroutine COMPUTECOUNTJOIN is shown below.

Algorithm COMPUTECOUNTJOIN

Input: new tuple k from window S_i and a join order $(S_1, \dots, S_{i-1}, S_{i+1}, \dots, S_n)$.
 $tmp = \arg \max_{u \in S_1, u.ts \leq k.ts}$
 $\forall u \in S_1$ and $u.cnt \geq tmp.cnt - C_1$
 If $k.attr \theta u.attr$
 $\dots \setminus \setminus$ loop through S_2 up to S_{i-2}
 $tmp = \arg \max_{v \in S_{i-1}, v.ts \leq k.ts}$
 $\forall v \in S_{i-1}$ and $v.cnt \geq tmp.cnt - C_{i-1}$
 If $k.attr \theta v.attr$
 $tmp = \arg \max_{w \in S_{i+1}, w.ts \leq k.ts}$
 $\forall w \in S_{i+1}$ and $w.cnt \geq tmp.cnt - C_{i+1}$
 If $k.attr \theta w.attr$
 $\dots \setminus \setminus$ loop through S_{i+2} up to S_{n-1}
 $tmp = \arg \max_{x \in S_n, x.ts \leq k.ts}$
 $\forall x \in S_n$ and $x.cnt \geq tmp.cnt - C_n$
 If $k.attr \theta x.attr$
 Return $k \circ \dots \circ x$

4 Cost Analysis

4.1 Insertion and Expiration

All NLJ-based algorithms incur a constant insertion cost per tuple: a new tuple is simply appended to its window. Hash-based algorithms, in addition, need to compute the hash function and insert tuples into the hash table, so their insertion costs are slightly higher. During every invalidation procedure, there are on average $\tau \lambda_i$ stale tuples in the i^{th} window, therefore the average number of tuples to be invalidated per unit time is $\sum_i \lambda_i$, which is independent of the expiration interval τ . Note that the actual insertion and expiration costs are implementation-dependent, e.g. if the sliding

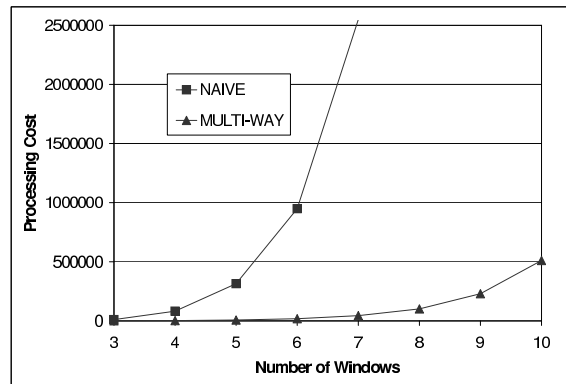


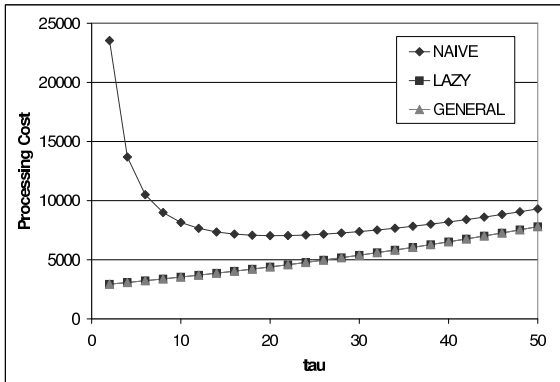
Figure 2: Scalability analysis of our (eager) multi-way join and the naive multi-way join. Derivations of the cost formulas may be found in Appendix A of the extended version of this paper [16].

window is implemented as a singly linked list with the most recent tuple at the head, removing tuples from the tail of the list may take linear time. In general, if invalidation is performed too frequently, some sliding windows may not contain any stale tuples, but we will still pay for the cost of accessing these windows. The situation is similar in hash joins: if there are many hash buckets and the invalidation frequency is high, some buckets may either be empty or may not contain any tuples to be expired. Hence, very frequent expiration may be too costly, especially in hash joins. We will validate this hypothesis empirically in Section 6.

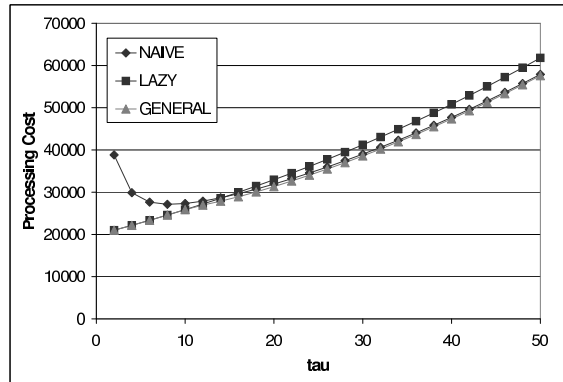
4.2 Join Processing Cost

To represent the cost of join processing, we use a per-unit-time cost model, developed by Kang et al. [20], and count the number of arithmetic operations on tuple attributes. When estimating join sizes, we make standard assumptions regarding containment of value sets and uniform distribution of attribute values. To clarify the former, suppose that an attribute $attr$ appears in several streams and has values a_1, a_2, \dots, a_v . Containment of value sets states that each window must choose its values from the front of the list a_1, a_2, \dots, a_i and have all values in this prefix. For simplicity, the subsequent discussion assumes that equi-joins on a common attribute are performed.

Given equivalent orderings, algorithm EAGER MULTI-WAY NLJ, abbreviated MULTI-WAY, outperforms algorithm NAIVE MULTI-WAY NLJ, abbreviated NAIVE, by reducing the amount of work done by the inner for-loops. For example, suppose that all windows have time size of one hundred, all streams arrive at a rate of one tuple per unit time, and each window has fifty distinct values. Figure 2 shows that our algorithm scales considerably better with the number of sliding windows than the naive technique.



(a)



(b)

Figure 3: Performance comparison of our lazy algorithms versus the naive lazy algorithm. Derivations of the cost formulas for each algorithm may be found in Appendix A of the extended version of this paper [16].

If expiration and re-execution are done every τ seconds, algorithm LAZY MULTI-WAY NLJ may not be optimal, as illustrated by the following examples. First, suppose that we have four windows with parameters as in the previous example. Figure 3(a) shows the performance of the best ordering of NAIVE MULTI-WAY NLJ (abbreviated NAIVE), LAZY MULTI-WAY NLJ (abbreviated LAZY), and GENERAL LAZY MULTI-WAY NLJ (abbreviated GENERAL). We let τ range from one to fifty, i.e. up to fifty percent of the window size. In this scenario, both LAZY and GENERAL have the same cost because the optimal ordering is always to order newly arrived tuples first. The cost steadily increases as τ increases because lazy expiration increases the average size of each window. The NAIVE strategy is more expensive than our improved algorithms, especially if τ is small, in which case the strategy of processing newly arrived tuples in the outer-most for-loop greatly reduces the amount of work done by the inner for-loops.

We now increase the arrival rate of the fourth stream to ten tuples per unit time, keeping its distinct value count at fifty. S_4 now has ten times as many tuples as before and ten times as many duplicates as the other windows. The performance of the best ordering of our algorithms is shown in Figure 3(b). As expected, GENERAL has the lowest cost of all the algorithms compared here regardless of the value of τ —this is always the case because the join ordering space considered by GENERAL is a superset of those of the others. Furthermore, LAZY works as well as GENERAL for small values of τ because the optimal ordering always places on top of the plan the source window of new tuples. However, as τ grows, there are many new S_4 -tuples generated between re-executions and placing S_4 on top of the plan causes the inner loops to do a large amount of work. In fact, the optimal ordering for processing S_4 -tuples when τ is large

is to order S_4 last—this is why LAZY is outperformed by the other algorithms in this situation. Conversely, NAIVE is the most expensive for small values of τ , but becomes nearly as good as GENERAL (and better than LAZY) when τ is large, provided that S_4 is ordered last. This is because the cost in this case is mostly made up of the cost of processing S_4 -tuples, which arrive at a faster rate than other tuples. Therefore, it is better to order S_4 last at all times, even if it means that other tuples will be processed in sub-optimal order (NAIVE), than it is to order S_4 first, even if tuples from other windows are processed in optimal order (LAZY).

4.3 Summary of Cost Analysis

EAGER MULTI-WAY NLJ always outperforms NAIVE MULTI-WAY NLJ under eager re-evaluation. Moreover, GENERAL LAZY MULTI-WAY NLJ is never worse than any of the other algorithms described here. We do not analytically compare the performance of hash joins with NLJs because while the former are clearly more efficient in terms of processing time, their expiration procedures may be more expensive. Since we do not know the cost of expiring a tuple relative to accessing a tuple during join processing (this depends on the implementation), we leave this issue for experimental analysis in Section 6.

5 Join Ordering

5.1 Eager Evaluation

We now investigate the effect of join ordering on the processing cost. All cost values given in this section are in units of attribute comparisons per unit time. We ignore the cost of inserting and expiring tuples, which is not affected by the join ordering. We also continue to assume that no materialization of intermediate results

Stream 1	$\lambda_1 = 10, T_1 = 100, v_1 = 500$
Stream 2	$\lambda_2 = 1, T_2 = 100, v_2 = 50$
Stream 3	$\lambda_3 = 1, T_3 = 200, v_3 = 40$
Stream 4	$\lambda_4 = 3, T_4 = 100, v_4 = 5$

Table 5: Stream parameters in the initial heuristic example.

may take place, leaving $n!$ possible extreme right-deep join trees from which to choose a plan (or $n \cdot n!$ for algorithm GENERAL LAZY MULTI-WAY NLJ), where n is the number of sliding windows to be joined.

We begin with a simple example of eager re-execution. Suppose that each window has the same number of distinct values. It is sensible to globally order the joins in ascending order of the window sizes (in tuples) $\lambda_i T_i$, or average hash bucket sizes $\lambda_i \frac{T_i}{b_i}$. That is, the window to which a new tuple arrives is processed in the outer-most for-loop, followed by the smallest remaining window and so on. This strategy minimizes the work done by the inner loops—by placing a small window in the outer for-loop, fewer tuples are passed down to the inner for-loops.

In the general case of eager re-execution, a sensible (but not always optimal) heuristic is to assemble the joins in descending order of binary join selectivities, leaving as little work as possible for the inner loops—we define a join predicate p_1 to be more selective than another, p_2 , if p_1 produces a smaller result set than p_2 . We calculate the intermediate result size of a join of two sliding windows S_m and S_n as $|S_m \bowtie S_n| = \frac{\lambda_m T_m \lambda_n T_n}{\max(v_m, v_n)}$. Consider four streams with parameters as shown in Table 5 and suppose that hash indices are not available. The global order chosen by this heuristic is either S_1, S_2, S_3, S_4 , or S_2, S_1, S_3, S_4 , depending on the tie-breaking procedure. The former is the optimal plan, which costs 16000 (an example derivation of this number may be found in Appendix B of the extended version of this paper [16]), while the latter is fifth-best at 19600. For comparison, the worst plan’s cost is nearly 90000. In this example, the cheapest plans are those with S_1 ordered first in the global order, suggesting an augmentation of our heuristic—explore plans where fast streams are ordered at or near the top of the plan.

To test the augmented heuristic, consider four streams with parameters given in Table 6. In the absence of hash tables, the best global order is S_2, S_1, S_3, S_4 , which costs 80400. The initial heuristic misses the optimal ordering and chooses either S_2, S_3, S_1, S_4 or S_2, S_3, S_4, S_1 , whose costs are approximately 123000 and 248000 respectively. Note that Stream 1 is faster than the others, so we try moving it up the order to get S_2, S_1, S_3, S_4 , which is the optimal global ordering. Interestingly, moving the fast stream all the way up to get S_1, S_2, S_3, S_4 is worse as it costs 120000. For comparison, the worst plan costs nearly

Stream 1	$\lambda_1 = 100, T_1 = 100, v_1 = 200$
Stream 2	$\lambda_2 = 1, T_2 = 100, v_2 = 200$
Stream 3	$\lambda_3 = 1, T_3 = 100, v_3 = 20$
Stream 4	$\lambda_4 = 3, T_4 = 100, v_4 = 2$

Table 6: Stream parameters in the augmented heuristic example.

650000.

If more than one stream is significantly faster than the others, as in the parameters shown in Table 7, the augmented heuristic still works well. Ordered by ascending join selectivity, we try S_3, S_4, S_1, S_2 , whose cost is 49542. We note that the two fast streams are ordered last, so we try to move them up. Moving up S_1 gives S_3, S_1, S_4, S_2 for a cost of 47977, which is the optimal ordering in this scenario, and moving up both S_1 and S_2 gives S_3, S_1, S_2, S_4 for a cost of 51954. Each combination considered by our heuristic costs less than the average cost per unit time over all orderings, which in this case is 63362. Again, moving the fast streams all the way to the top to get S_1, S_2, S_3, S_4 or S_2, S_1, S_3, S_4 is not recommended as these two plans have above-average costs of 68200 and 79000 respectively.

Stream 1	$\lambda_1 = 11, T_1 = 100, v_1 = 200$
Stream 2	$\lambda_2 = 10, T_2 = 100, v_2 = 100$
Stream 3	$\lambda_3 = 1, T_3 = 100, v_3 = 65$
Stream 4	$\lambda_4 = 1, T_4 = 100, v_4 = 20$

Table 7: Stream parameters in the example with two fast streams and two slow streams.

In summary, a reasonable heuristic for eager re-evaluation of the multi-way NLJ is to initially order the joins in descending order of their selectivities. If one or more streams are faster than the others, it is also beneficial to consider orderings where the fast streams are moved up the join order (but not all the way to the top). The number of orderings considered is on the order of n^f where f is the number of “fast” streams.

5.2 Lazy Re-evaluation

Recall from Figure 3 that algorithm LAZY MULTI-WAY NLJ is as efficient as GENERAL LAZY MULTI-WAY NLJ when τ is small (how small τ must be for this to be true depends on the query). If this is the case, we may use the same ordering heuristics as algorithm LAZY MULTI-WAY NLJ is a straightforward extension of its eager version. Otherwise, algorithm GENERAL LAZY MULTI-WAY NLJ is more efficient if a good join ordering is chosen. Recall that the general NLJ allows arbitrary join orderings depending on the origin of the new tuples that are being processed, i.e. newly arrived tuples are not necessarily processed in the outer for-loop. However, these local orderings are

independent of one another, as evidenced by the additivity of the processing cost derived in Appendix A of [16]. Therefore we may re-apply the existing heuristic and independently optimize each local join order O_i for $1 \leq i \leq n$, where n is the number of sliding windows.

5.3 Join Ordering in the Multi-Way Hash Join

If each hash table has the same number of buckets, the ordering problem is the same as in the NLJ. This is because the hash-join so configured operates in a nested-loop fashion like the NLJ, except that at each loop, only one hash bucket is scanned instead of the entire window. Since the number of hash buckets is the same for each window, the cost savings are equal at each for-loop. Our heuristics can also be used if the hash tables have various sizes so long as we discount each window’s time size T_i by the number of buckets to get the average bucket size $\frac{T_i}{b_i}$. This gives the average number of tuples in S_i that will be scanned when performing the join.

5.4 Join Ordering in Other Scenarios

We conclude the discussion on join ordering with examples of how our general principle of minimizing the work done by the inner for-loops may be applied in other scenarios, including those in which some of our simplifying assumptions do not hold.

- **Hybrid Hash-NLJ:** A simple heuristic is to place all the windows that contain hash indices in the inner for-loops since those repeat the most often. At this point, the index selection problem in the context of our hybrid operator comes into play: given the constraints on available memory, on which windows should we build hash indices?
- **Expensive predicates:** Those may need to be ordered near the top of the join tree in order to minimize the work done by the inner for-loops.
- **Joins on different attributes:** In this case, we may no longer arbitrarily re-order the join tree. However, it may still be more efficient to place the window from which the new tuples came at the outer-most for-loop. Thereafter, a simple greedy heuristic could be at each step to choose the most selective join predicate allowed by the join graph.
- **Fluctuating stream arrival rates:** If feasible, we re-execute the ordering heuristic whenever stream rates change beyond a given threshold. Otherwise, it may be prudent to place near the top of the plan all those streams which are expected to fluctuate widely—recall that our heuristics often select plans where fast streams are near the top.

6 Experimental Results

6.1 Experimental Setting

In this section, we validate our join ordering heuristics, compare the performance of our join operators, and investigate the trade-offs associated with re-evaluation and expiration frequencies. We have built simple prototypes of our algorithms using Sun Microsystems JDK 1.3.1, running on a Windows PC with a 1.2 GHz AMD processor and 256 megabytes of RAM. We have implemented sliding windows and hash buckets as singly-linked lists with the most recent tuple at the tail (i.e. hash tables are implemented as one-dimensional arrays of linked lists). Thus, insertion adds a tuple at the end of a list or bucket, while expiration advances the pointer to the first element. All hash functions are simple modular divisions by the number of hash buckets. For simplicity, each tuple consists of two integer attributes: a system-assigned timestamp ts and a common join attribute $attr$, and the join predicate is hard-coded to be an equality comparison. Eager and lazy query re-evaluation, as well as eager and lazy tuple expiration, are supported. Note that the expiration procedure does not delete old tuples as this is under control of Java’s garbage collection mechanism. We have experimented with Java’s `System.gc()` method, which acts as a suggestion for Java to perform garbage collection, and found that our simulation slows down when this method is called frequently. We do not discuss this further as it is an implementation-specific issue.

Our tuple generation procedure is as follows. We run a continuous for-loop, inside which one tuple per iteration is generated from a random stream i with probability equal to $\frac{\lambda_i}{\sum_{j=1}^n \lambda_j}$. The tuple is given a timestamp equal to the current loop index and an attribute value chosen uniformly at random from the set $\{1, 2, \dots, v\}$, where v is the number of distinct values in the new tuple’s source stream. This procedure simulates relative stream arrival rates and guarantees containment of value sets. We repeat the loop x times, measure the time t taken to execute the simulation, and report $\frac{x}{t}$, which is the maximum input rate that can be supported by the system. This is the cumulative rate over all inputs, measured in tuples per second. In order to eliminate transient effects occurring when the windows are initially non-full, the simulator first generates enough random tuples to fill the windows. We repeat each experiment ten times and report the average.

6.2 Validation of Cost Model and Join Ordering Heuristic

We begin by executing a simple query with parameters as shown in Table 5 and compare maximum supported input rates for various orderings. We execute the multi-way NLJ and the multi way hash join (with

Algorithm	Max. rate of best plan	Max. rate of worst plan
Eager NLJ	1614	333
Lazy NLJ, $\tau = 5$	1446	296
Lazy NLJ, $\tau = 10$	1332	274
Eager hash	11540	2524
Lazy hash, $\tau = 5$	8420	2041
Lazy hash, $\tau = 10$	7947	1848

Table 8: Experimental validation of our cost model and join ordering heuristic.

five buckets per hash table) under eager re-evaluation and expiration, and lazy re-evaluation and expiration with two choices of τ : five and ten. Results are shown in Table 8. Recall from Section 5.1 that the worst plan for this query costs roughly five times as much to process as the optimal plan in the case of algorithm MULTI-WAY NLJ. As seen in Table 8, when insertion and expiration costs have been included, this difference diminishes somewhat. We have verified that the best plan in all cases is the plan predicted by our heuristic. Note that in this particular example, the hash join outperforms the NLJ and that increasing the re-evaluation interval hurts performance; we elaborate on these points in the next two experiments.

6.3 NLJ vs. Hash Join

We now compare the performance of our NLJ versus our hash join (with bucket sizes 2, 10, and 30, abbreviated HASH2, HASH10, and HASH30 respectively). We run a four-window query in which each window’s time size is 100 units, streams arrive with rates of one tuple per unit time, and each window has fifty distinct values. Figure 4(a) shows a graph of the maximum input rate versus the re-execution (and expiration) interval τ . We will discuss the effect of expiration frequency later on; for now, we note that the relative performance of various join configurations is as expected: the NLJ is the slowest and the hash join with the largest hash table is the fastest. The greatest relative improvement occurs between the NLJ and the hash join with two hash buckets per window. Thereafter, the marginal improvement of adding more hash buckets decreases.

6.4 Effect of Query Re-Evaluation and Expiration Frequencies on Processing Time

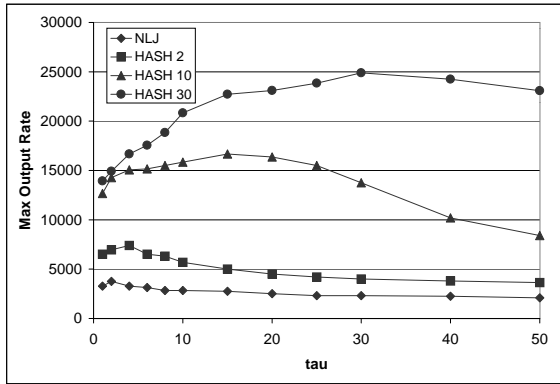
Returning to Figure 4(a), we now explain the effect of re-evaluation and expiration frequencies on the overall performance of our algorithms. Eager expiration incurs the cost of updating linked list pointers upon

every arrival of a new tuple. In contrast, lazy expiration performs fewer of these operations, but allows the windows to grow between updates, causing longer join evaluation times. Beginning with the NLJ, we see that short expiration intervals are preferred as the cost of advancing pointers is lower than the cost of scanning larger windows in the processing stage. However, very frequent expiration and re-evaluation are inefficient. The same pattern is seen in hash joins, except that the optimal frequency increases as the number of buckets increases. For instance, the two-bucket hash join performs best when expiration and execution are done every four time units (given the parameters in this example), for the ten-bucket join the optimum is roughly fifteen time units, and for the 30-bucket join, the optimum is approximately 30 time units. As discussed in Section 4.1, this is the case because if the number of hash buckets is large and the invalidation frequency is high, then many buckets are empty or do not contain any tuples to be expired. However, we still pay for the cost of accessing the buckets and searching for old tuples. This explains why the eager version of the hash join shown in Table 8 achieves a much faster output rate than the lazy version—in that example, there were only five buckets per hash table.

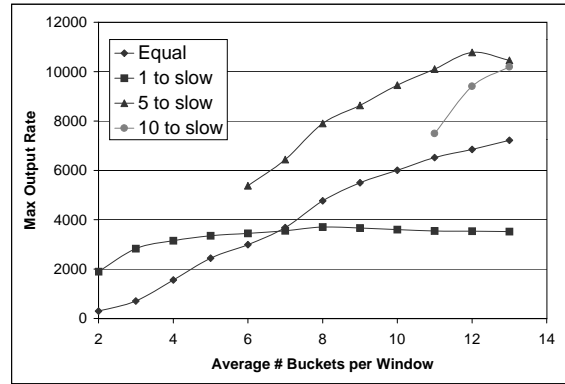
We also note that the maximum output rate drops off gradually when τ is large. Again, this is because the sliding windows are large (they have many old tuples), therefore the join spends more time scanning the windows. However, when τ is fairly small, this increase in the join processing cost is overshadowed by the savings in tuple expiration cost when using large hash tables, as seen in the increase in the optimal expiration frequency as the number of hash buckets increases.

6.5 Varying Hash Table Sizes

Thus far, we have only evaluated hash joins where each hash table has the same number of buckets. We now consider a strategy that builds larger hash tables for large windows (large in terms of time size or stream arrival rate). We execute a join of four windows, 100 time units large and with fifty distinct values. Three streams arrive with rates of one tuple per unit time, while the fourth stream is fifty times faster. We re-execute the query and expire old tuples every five time units to provide a reasonably high refresh frequency. Results are shown in Figure 4(b) for four scenarios: equal hash table sizes (Equal), allocating one hash bucket each to slow streams (1 to slow—in this case, slow streams do not have hash indices), allocating five buckets each to slow streams (5 to slow), and allocating ten buckets each to slow streams (10 to slow). The horizontal axis measures the average number of hash buckets per window. For example, if this number is 12, then in the “Equal” strategy, every hash table has 12 buckets, while in the “5 to slow” strategy, the three hash tables corresponding to slow streams have five



(a)



(b)

Figure 4: Performance comparison of our algorithms with respect to a) increasing the re-evaluation and expiration interval and b) building large hash tables on fast streams.

buckets each and the hash table over the fast stream has $4 \cdot 12 - 3 \cdot 5 = 33$ hash buckets.

In the “Equal” strategy, the maximum input rate increases gradually as the number of buckets increases, reaching over 7000 tuples per unit time when each hash table has 13 buckets. In “1 to slow”, the best rate is obtained when the average number of buckets is eight, i.e. when the largest hash table has 29 buckets. Increasing the largest hash table further hurts performance because of the large expiration cost at this (fairly small) value of τ . “5 to slow” works very well and reaches a maximum rate of over 10000 tuples per unit time when the average hash table size is 12 buckets (i.e. when the largest hash table has 33 buckets). Finally, “10 to slow” does not work as well as “5 to slow” for the hash table sizes shown on the graph, but continues to improve and surpasses “5 to slow” if more hash buckets are available (not shown on the graph). In summary, assuming that counting the average number of buckets per hash table is a fair cost metric, allocating more hash buckets to a large window improves performance. However, the largest hash table must not be too large because large hash tables do not perform well when τ is small.

7 Conclusions and Future Work

We have presented and analyzed incremental, multi-way join algorithms for sliding windows over data streams. Based on a per-unit-time cost model, we developed a join ordering heuristic that finds a good join order without iterating over the entire search space. Experiments with a Java implementation of our operators have shown that various system parameters such as stream arrival rates, tuple expiration policies, and hash table sizes greatly affect query processing efficiency. In particular, the optimal tuple expiration frequency depends on the choice of algorithm and on the

sizes of the hash tables—the larger the hash table, the higher the optimal frequency. In general, hash-based joins performed much better than NLJs and we have also discovered that allocating more hash buckets to larger windows is a promising strategy.

Our research goal is to develop a sliding window query processor that is functional, efficient, and scalable. In terms of functionality, we intend to design efficient algorithms for other query operators over sliding windows besides joins (e.g. windowed sort and windowed top- k list). Improvements in efficiency include low-overhead indices for indexing window contents, and exploiting stream properties such as near-sortedness and foreign key relationships. Finally, scalability may be improved by means of grouped evaluation of similar queries, indexing query predicates, storing materialized views, and returning approximate answers if exact answers are too expensive to compute.

References

- [1] S. Acharya, P. Gibbons, V. Poosala, S. Ramaswamy. Join synopsis for approximate query answering. In *Proc. ACM Int. Conf. on Management of Data*, 1999, pp. 275–286.
- [2] A. Arasu, B. Babcock, S. Babu, J. McAlister, J. Widom. Characterizing Memory Requirements for Queries over Continuous Data Streams. In *Proc. ACM SIGACT-SIGMOD Symp. on Principles of Database Systems*, 2002, pp. 221–232.
- [3] A. Arasu, S. Babu, J. Widom. An Abstract Semantics and Concrete Language for Continuous Queries over Streams and Relations. Stanford University Technical Report 2002-57, November 2002.
- [4] B. Babcock, S. Babu, M. Datar, R. Motwani, J. Widom. Models and Issues in Data Streams. In *Proc. ACM Symp. on Principles of Database Systems*, 2002, pp. 1–16.

- [5] B. Babcock, M. Datar, R. Motwani. Sampling from a Moving Window over Streaming Data. In *Proc. 13th SIAM-ACM Symp. on Discrete Algorithms*, 2002, pp. 633–634.
- [6] S. Babu, J. Widom. Exploiting k-Constraints to Reduce Memory Overhead in Continuous Queries over Data Streams. Stanford University Technical Report 2002-52, November 2002.
- [7] D. Carney, U. Cetintemel, M. Cherniack, C. Convey, S. Lee, G. Seidman, M. Stonebraker, N. Tatbul, S. Zdonik. Monitoring streams – A New Class of Data Management Applications. In *Proc. 28th Int. Conf. on Very Large Data Bases*, 2002, pp. 215–226.
- [8] S. Chandrasekaran, O. Cooper, A. Deshpande, M. Franklin, J. Hellerstein, W. Hong, S. Krishnamurthy, S. Madden, V. Raman, F. Reiss, M. Shah. TelegraphCQ: Continuous Dataflow Processing for an Uncertain World. In *Proc. 1st Biennial Conf. on Innovative Data Syst. Res. (CIDR)*, 2003, pp. 269–280.
- [9] J. Chen, D. DeWitt, F. Tian, Y. Wang. NiagaraCQ: A Scalable Continuous Query System for Internet Databases. In *Proc. ACM Int. Conf. on Management of Data*, 2000, pp. 379–390.
- [10] C. Cortes, K. Fisher, D. Pregibon, A. Rogers, F. Smith. Hancock: A Language for Extracting Signatures from Data Streams. In *Proc. 6th ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining*, 2000, pp. 9–17.
- [11] M. Datar, A. Gionis, P. Indyk, R. Motwani. Maintaining Stream Statistics over Sliding Windows. In *Proc. 13th SIAM-ACM Symp. on Discrete Algorithms*, 2002, pp. 635–644.
- [12] A. Dobra, M. Garofalakis, J. Gehrke, R. Rastogi. Processing Complex Aggregate Queries over Data Streams. In *Proc. ACM Int. Conf. on Management of Data*, 2002, pp. 61–72.
- [13] A. C. Gilbert, Y. Kotidis, S. Muthukrishnan, M. Strauss. QuickSAND: Quick Summary and Analysis of Network Data. DIMACS Technical Report 2001-43, Dec. 2001.
- [14] A. C. Gilbert, Y. Kotidis, S. Muthukrishnan, M. Strauss. Surfing Wavelets on Streams: One-Pass Summaries for Approximate Aggregate Queries. In *Proc. 27th Int. Conf. on Very Large Data Bases*, 2001, pp. 79–88.
- [15] L. Golab, M. T. Özsu. Issues in Data Stream Management. In *SIGMOD Record*, 32(2), June 2003, to appear.
- [16] L. Golab, M. T. Özsu. Processing Sliding Window Multi-Joins in Continuous Queries over Data Streams. University of Waterloo Technical Report, February 2003. Available at <http://db.uwaterloo.ca/~ddbms/publications/stream/multijoins.pdf>.
- [17] S. Guha, N. Koudas, K. Shim. Data-Streams and Histograms. In *Proc. 33rd ACM Symp. on Theory of Computing*, 2001, pp. 471–475.
- [18] P. Haas, J. Hellerstein. Ripple Joins for Online Aggregation. In *Proc. ACM Int. Conf. on Management of Data*, 1999, pp. 287–298.
- [19] M. Hammad, M. Franklin, W. Aref, A. Elmagarmid. Scheduling for Shared Window Joins over Data Streams. In *Proc. 29th Int. Conf. on Very Large Data Bases*, 2003, to appear.
- [20] J. Kang, J. Naughton, S. Viglas. Evaluating Window Joins over Unbounded Streams. In *Proc. 19th Int. Conf. on Data Engineering*, 2003, to appear.
- [21] L. Liu, C. Pu, W. Tang. Continual Queries for Internet-Scale Event-Driven Information Delivery. In *IEEE Trans. on Knowledge and Data Eng.*, 11(4): 610–628, 1999.
- [22] S. Madden, M. Shah, J. Hellerstein, V. Raman. Continuously Adaptive Continuous Queries Over Streams. In *Proc. ACM Int. Conf. on Management of Data*, 2002, pp. 49–60.
- [23] R. Motwani, J. Widom, A. Arasu, B. Babcock, S. Babu, M. Datar, G. Manku, C. Olston, J. Rosenstein, R. Varma. Query Processing, Approximation, and Resource Management in a Data Stream Management System. In *Proc. 1st Biennial Conf. on Innovative Data Syst. Res. (CIDR)*, 2003, pp. 245–256.
- [24] M. Sullivan, A. Heybey. Tribeca: A System for Managing Large Databases of Network Traffic. In *Proc. USENIX Annual Technical Conference*, 1998.
- [25] D. Terry, D. Goldberg, D. Nichols, B. Oki. Continuous Queries over Append-Only Databases. In *Proc. ACM Int. Conf. on Management of Data*, 1992, pp. 321–330.
- [26] P. Tucker, D. Maier, T. Sheard, L. Fegarar. Exploiting Punctuation Semantics in Continuous Data Streams. In *IEEE Trans. on Knowledge and Data Eng.*, 15(3): 555–568, 2003.
- [27] T. Urhan, M. Franklin. XJoin: A Reactively-Scheduled Pipelined Join Operator. In *IEEE Data Engineering Bulletin*, 23(2):27–33, 2000.
- [28] S. Viglas and J. Naughton. Rate-Based Query Optimization for Streaming Information Sources. In *Proc. ACM Int. Conf. on Management of Data*, 2002, pp. 37–48.
- [29] S. Viglas, J. Naughton, J. Burger. Maximizing the Output Rate of Multi-Way Join Queries over Streaming Information Sources. In *Proc. 29th Int. Conf. on Very Large Data Bases*, 2003, to appear.
- [30] A. Wilschut, P. Apers. Dataflow query execution in a parallel main-memory environment. In *Proc. 1st Int. Conf. on Parallel and Distributed Information Syst.*, 1991, pp. 68–77.
- [31] K.-Y. Whang, R. Krishnamurthy. Query Optimization in a Memory-Resident Domain Relational Calculus Database System. In *ACM Trans. on Database Syst.*, 15(1):67–95, 1990.
- [32] Y. Yao, J. Gehrke. Query Processing in Sensor Networks. In *Proc. 1st Biennial Conf. on Innovative Data Syst. Res. (CIDR)*, 2003, pp. 233–244.
- [33] Y. Zhu, D. Shasha. StatStream: Statistical Monitoring of Thousands of Data Streams in Real Time. In *Proc. 28th Int. Conf. on Very Large Data Bases*, 2002, pp. 358–369.