

Scalable Algorithms for Large Temporal Aggregation *

Bongki Moon[†] Inés Fernando Vega López[†] Vijaykumar Immanuel[‡]

[†]*Dept. of Computer Science
University of Arizona
Tucson, AZ 85721*

{*bkmoon, ifvega*} @cs.arizona.edu

[‡]*Compaq Computer Corporation
19333 Vallco Pkwy
Cupertino, CA 95014*

vijay.immanuel@compaq.com

Abstract

The ability to model time-varying natures is essential to many database applications such as data warehousing and mining. However, the temporal aspects provide many unique characteristics and challenges for query processing and optimization. Among the challenges is computing temporal aggregates, which is complicated by having to compute temporal grouping. In this paper, we introduce a variety of temporal aggregation algorithms that overcome major drawbacks of previous work. First, for small-scale aggregations, both the worst-case and average-case processing time have been improved significantly. Second, for large-scale aggregations, the proposed algorithms can deal with a database that is substantially larger than the size of available memory.

1. Introduction

Database applications often need to capture the time-varying nature of an enterprise they model. The importance of such need has been recognized by several database research groups, and temporal database models and query languages have been developed and reported in the literature [8, 13]. In fact, there are several temporal query languages supporting temporal aggregation [12]. However, temporal data and queries provide many unique characteristics and challenges for query processing and optimization. Among the challenges is computing temporal aggregates, which is complicated by having to compute *temporal grouping*.

In temporal databases, temporal grouping is a process where the time-line is partitioned over time and tuples are grouped over these partitions. Then, aggregate val-

ues are computed over these groups. In general, temporal grouping is done by two types of partitioning [12]: *span grouping* and *instant grouping*. Span grouping is based on a defined length in time, such as week or month, and is independent of temporal attribute values of database tuples. On the other hand, instant grouping depends on the data stored. Any pair of consecutive instants create a time interval, over which the aggregate value remains constant. Such intervals are called constant intervals. Aggregations based on span and instant groupings are called *span aggregation* and *instant aggregation*, respectively. In this paper, we focus on computing instant aggregates, which we believe is the most common and challenging temporal aggregation.

Computing instant aggregates is expensive because it is necessary to know which tuples overlap each instant, and simply considering each tuple in order in a sorted-by-time relation will not be sufficient due to the varying interval lengths [9]. For example, computing the time-varying maximum salary of employees involves computing the temporal extent of each maximum value, which requires determining the tuples that overlap each temporal instant. Figure 1(a) shows a sample *Employees* table with two temporal attributes, which represent the beginning and ending of the valid-times of individual tuples. The resulting instant aggregation of the maximum salary (along with the number of employees) is given in the table in Figure 1(b). Note that while multiple values are returned, the aggregation results in a single scalar value at each point in time, with the period over which the aggregate value remains constant collected into a single tuple. One could also envision an instant aggregate function, which would evaluate a time-varying maximum salary for each department.

This temporal aggregation can be processed in a sequential or parallel fashion. The parallel processing technology becomes even more attractive, as the size of data-intensive applications grows as evidenced in OLAP and data warehousing environments [3]. Although several sequential and parallel algorithms have been devel-

*This work was sponsored in part by National Science Foundation CAREER Award (IIS-9876037) and Research Infrastructure program EIA-9500991. It was also supported by Consejo Nacional de Ciencia y Tecnología, scholarship 117476. The authors assume all responsibility for the contents of the paper.

Name	Salary	Dept	Begin	End
Richard	46,000	Accounting	18	31
Karen	45,000	Shipping	8	20
Nathan	35,000	Marketing	7	12
Nathan	38,000	Accounting	18	21

(a) Input Database Tuples

Count	Max	Begin	End
1	35,000	7	8
2	45,000	8	12
1	45,000	12	18
3	46,000	18	20
2	46,000	20	21
1	46,000	21	31

(b) Temporal Aggregation Results

Figure 1. Sample Database and Its Temporal Aggregation

oped for computing temporal aggregates [7, 9, 12, 14, 15], they suffer from serious limitations such as the size of aggregation restricted by available memory and requirement of a priori knowledge about the orderedness of an input database.

In this paper, we propose a variety of temporal aggregation algorithms that overcome major drawbacks of previous work. The proposed solutions provide the following benefits over the state of the art:

- Two new algorithms proposed for small-scale aggregations do not require a priori knowledge about an input database, and they have improved both the worst-case and average-case processing time significantly.
- Another new algorithm proposed for large-scale aggregations relies on a novel data partitioning scheme, so that it can deal with a database substantially larger than the size of available memory.

It should be noted that the problem of computing temporal aggregates is different from the relational aggregation that can often be seen in the data warehousing environment. While data items in the data warehousing environment are envisioned as points in their data domain, we deal with temporal data associated with time intervals of arbitrary lengths.

The rest of this paper is organized as follows. Section 2 surveys the background and related work on computing temporal aggregates. Major limitations of previous work are also discussed in the section. In Sections 3, 4, and 5, we present the improved algorithms for small-scale aggregations, and scalable solutions for large-scale aggregations based on data partitioning and parallel processing techniques. Section 6 presents the

results of experimental evaluation of the proposed solutions. Finally, Section 7 summarizes the contributions of this paper and gives an outlook to future work.

2. Background and Previous Work

There are two types of aggregate computations in conventional relational database systems: scalar aggregates and aggregate functions. Scalar aggregates are operations such as `count`, `sum`, `avg`, `max`, and `min` that produce a single value over an entire relation, while aggregate functions first partition a relation based on some attribute value and then compute scalar aggregates independently on the individual partitions.

A scalar aggregate is composed of an aggregate expression and an optional qualification. A simple two-step algorithm was proposed by Epstein for evaluating scalar aggregates [5]. To handle many scalar aggregates in a query, the algorithm computes each of them separately and stores each result in a singleton relation, referring to that singleton relation when evaluating the rest of the query. A different approach employing program transformation methods was proposed to systematically generate efficient iterative programs for aggregate queries [6].

The first approach for implementing temporal aggregation was proposed by Tuma [14] and was based on an extension of Epstein’s algorithm. In this approach, the constant intervals are determined first, then the aggregate is evaluated using the Epstein’s technique. Since the two steps are separate and the first one must be completed before the second one, a database must be read twice.

More recent algorithms were proposed by Kline and Snodgrass [9] for temporal aggregation based on instant grouping of tuples. The algorithms are called *aggregation tree* and its variant *k-ordered aggregation tree*, as they build a tree while scanning a database. Both algorithms are fast and require minimal I/O overhead, as they need to scan the database only once to build a tree in memory. Then, the resulting tree stores enough information to compute temporal aggregates by traversing it using depth first search.

2.1. Limitations of Previous Methods

It should be noted that the order of tuples inserted into the aggregation tree affects its performance, though not its result. If the tuples are sorted via the start time and inserted in that order, the aggregation tree would look more like a linked list, causing insertions to be slower than insertions into a balanced binary tree. For the reason, the worst case time to create an aggregation tree is $\mathcal{O}(N^2)$ for N tuples sorted in time. Even more serious limitation of the aggregation tree approach is that the

entire tree must be kept in memory. Since the size of an aggregation tree is proportional to the number of distinct timestamps (both start times and end times), the size of the database the aggregation tree algorithm can deal with tends to be limited by the size of available memory and the number of distinct timestamps of tuples.

To circumvent this problem, a variant of the aggregation tree, called k -ordered aggregation tree, was proposed by the same authors. The k -ordered aggregation tree takes advantage of the k -orderedness of tuples to enable garbage collection of tree nodes, so that the memory requirements can be reduced significantly. However, the k -ordered aggregation tree approach assumes that the tuples in a table be ordered within a certain degree. Specifically, each tuple is at most k positions from its position in a totally ordered version of the table. This requirement is difficult to be met in a real database system. Without a priori knowledge about a given table, the k -orderedness is expensive to measure, as it requires an external sort of the table. The worst case running time of the k -ordered aggregation tree algorithm is still $\mathcal{O}(N^2)$.

3. Improved Algorithms for Small-Scale Aggregation

In this section, we present two new algorithms for computing temporal aggregates, as alternatives to the aggregation tree algorithm [9]. The aggregation tree is a binary tree, which is similar to the segment tree by Bentley [2]. The segment tree is a static structure, which can be balanced for a given set of abscissae. However, there is no guarantee that the aggregation tree is always balanced, because the aggregation tree is dynamically constructed as the tuples in a database are being scanned and inserted into the tree. Thus, the structure of the resulting aggregation tree depends on the order of tuples inserted. This fact may cause the worst case running time of $\mathcal{O}(N^2)$ for a database of N tuples, particularly when the tuples are ordered by their timestamp values. Such a quadratic complexity may be impractically costly for many database applications.

As will be seen in this section, we have observed that the five most common aggregation operators can be categorized into two groups, namely, `count`, `sum`, `avg` in one group, and `max`, `min` in the other. For the latter group, there is more demand to keep track of attribute values of tuples. This observation has led us to develop a different algorithm for each of the two groups of aggregation operators. The solution to the first group of operators, which we call a *balanced tree* algorithm, will be presented in Section 3.1. The main idea of this algorithm is that the tree can be balanced dynamically as tuples are being inserted, by *giving up the notion of maintaining intervals* in the tree nodes. The solution to the second group is called a *merge-sort aggregation*

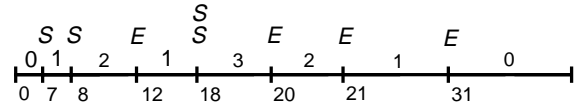


Figure 2. Example of `count` Aggregation by Sorting Timestamps and Tags

algorithm, which is similar to the classical merge-sort algorithm [10]. This algorithm will be presented in Section 3.2. In this section, we assume that the memory is large enough to store the entire data structures required by each aggregation algorithm. In the rest of this paper, we use the `count` and `max` as the representatives of the two groups of operators, respectively.

3.1. Balanced Tree Algorithm for `count` Aggregation

A relatively simple approach based on timestamp sorting can provide an efficient solution for the `count` aggregation. This approach starts with loading the entire tuples in memory. Then, the timestamp values are extracted from the tuples, and each timestamp is associated with a tag, which indicates whether the timestamp is a start time or an end time of a tuple. These timestamps and tags are then sorted in an increasing order of the timestamp values. See Figure 2 for a sorted list of timestamps and tags for a sample database given in Figure 1(a).

The `count` aggregate is computed by scanning the sorted timestamps and tags in an increasing order. Getting started with a counter initialized to zero, the counter is incremented by one when a START tag is encountered, and it is decremented by one when an END tag is encountered. When more than one tags are associated with a timestamp, the counter is incremented by the number of START tags or decremented by the number of END tags. For example, in Figure 2, when the timestamp value 18 is encountered, the counter is incremented by two from 1 to 3 because there are two START tags associated with the timestamp. Apparently, the worst case processing time of this approach is $\mathcal{O}(N \log N)$, where N is the number of tuples in an input database.

In real world temporal databases, it may be the case that many tuples share the same timestamp values for their start times and end times. Nonetheless, this timestamp-sort approach requires the same amount of memory and processing time regardless of the repeated timestamp values. Thus, we propose a *balanced tree* algorithm to further optimize its performance for such databases with repeated timestamp values.

The motivation behind the balanced tree algorithm is

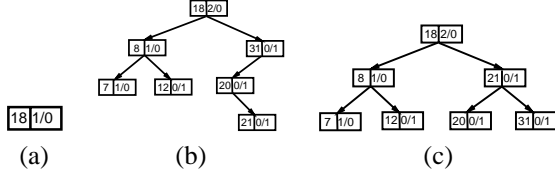


Figure 3. Example of Balanced Tree Construction

that the sorted list of timestamps can be built even without loading an entire database into memory at once. Instead, the timestamps can be sorted *incrementally* by inserting them into a balanced tree, as the tuples of an input database are being scanned. Each node of a balanced tree stores a timestamp, either a start time or an end time, but need not store a START/END tag. Instead of the tag, each node stores two counters: one storing the number of tuples starting at the timestamp and the other storing the number of tuples ending at the timestamp.¹ Additionally, a color tag is stored in each node, as we use the *red-black* insertion algorithm [4] to keep the tree balanced dynamically.

Figure 3 shows the process of building a balanced tree for the sample `Employees` table in Figure 1(a). In the figure, we only show timestamps and counters, which are relevant to temporal aggregate computation. When the start time 18 of the first record is inserted into an empty tree, a new node is created for the timestamp, and then its start-counter and end-counter are set to one and zero, respectively. The resulting tree having a single node is shown in Figure 3(a). Figures 3(b) and (c) illustrate snapshots of the tree before and after the tree is balanced by the red-black insertion algorithm. We do not elaborate on the red-black insertion because it is not the focus of this paper.

The balanced tree algorithm proceeds in two steps, first by creating the tree and then by traversing the tree. Whenever a tuple is read from an input database, the balanced tree is probed to see whether the start and end times of the tuple are already in the tree. If the start (or end) timestamp is not found in the tree, then a new node is created and inserted into the tree. Otherwise, the start time (or end time) counter of a node that contains the timestamp is incremented by one without inserting a new node. Once the balanced tree has been built, the algorithm computes aggregate values while performing an in-order traversal of the tree. Specifically, whenever a tree node is visited, the `count` aggregate value is incremented by the start-counter value of the node and

¹For sum aggregation, each node stores two variables: one storing the attribute value sum of the tuples starting at the timestamp and the other storing the attribute value sum of the tuples ending at the timestamp.

Algorithm 1: Balanced Tree

```

set  $\mathcal{T} \leftarrow$  an empty balanced tree;
foreach tuple  $t$  in a table do
  if ( $t.start\_time = n.ts$  for any node  $n$  in  $\mathcal{T}$ ) then
     $n.no\_starts++$ ;
  else
    insert a new node  $n'$  (with  $n'.ts = t.start\_time$ ) into  $\mathcal{T}$ ;
  if ( $t.end\_time = n.ts$  for any node  $n$  in  $\mathcal{T}$ ) then
     $n.no\_ends++$ ;
  else
    insert a new node  $n'$  (with  $n'.ts = t.end\_time$ ) into  $\mathcal{T}$ ;

set  $count \leftarrow 0$ ;
foreach node  $n$  in  $\mathcal{T}$  traversed by in-order do
   $count += n.no\_starts$ ;
  output  $n.ts$  and  $count$ ;
   $count -= n.no\_ends$ ;

```

decremented by the end-counter value of the node. The proposed balanced tree algorithm is summarized in Algorithm 1.

By eliminating redundant timestamp values from the tree, the balanced tree algorithm reduces the memory requirements and tree traversal time substantially especially for a database with a small percentage of unique timestamps. The balanced tree stores information needed for temporal grouping and aggregation both in internal nodes and leaf nodes. Thus, the balanced tree algorithm uses only half the nodes required by the aggregation tree algorithm, which stores constant intervals only in leaf nodes.

3.2. Merge-Sort Algorithm for `max` Aggregation

While the balanced tree algorithm is simple and efficient for `count` aggregations, it cannot be used for `max` aggregations. Since a balanced tree stores only unique timestamps and associated counters for `count` aggregation, it is not possible to keep track of all the tuples that are alive at a given time instant with the information available in the tree. For example, in Figure 3(b), the root node shows that there exist two tuples whose start times are 18. However, the tree does not convey any information about the life spans of the tuples (*i.e.*, the exact end times of the two specific tuples). Unlike `count` aggregations, it is impossible to compute `max` aggregations without knowing the exact life spans of tuples in a database.

One can modify the balanced tree algorithm to compute `max` aggregates, by allowing repeated timestamp values in a tree and using additional data structures such as dual heaps while traversing the tree. The dual heaps store the attribute values (on which the `max` aggregation is performed) of live tuples and dead tuples, separately. While traversing the tree, the `max` aggregate can be computed by comparing two maximum values in both the heaps and popping matched maximum values from the heaps. In fact, the dual heaps are used to keep track

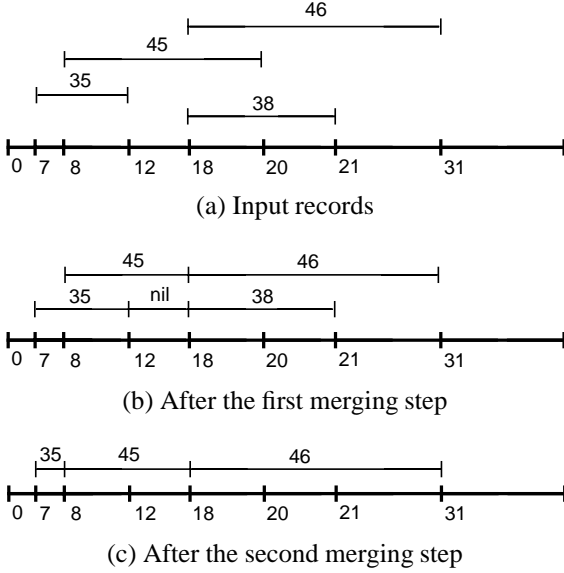


Figure 4. Example of Merging for \max Aggregation

of the life spans of tuples that are required to compute the \max aggregate. However, with this modification, we will lose all the benefits of using the balanced tree algorithm, because the tree will need exactly two nodes per each tuple (*i.e.*, no reduction in memory requirements due to repeated timestamps) and additional overhead for processing the heaps will be non-trivial.

Instead, we propose a *bottom-up* aggregation approach, which we call a *merge-sort aggregation* algorithm. Like the classical merge-sort algorithm based on the divide-and-conquer strategy, the merge-sort aggregation algorithm computes a larger (intermediate) aggregate result by merging two smaller (intermediate) aggregate results. The algorithm starts with merging tuples in pairs at the bottom and terminates when a final aggregate result is obtained at the top.

Formally, an intermediate aggregate can be defined as (T_k, M_k) , where $T_k = \{t_0, t_1, \dots, t_k\}$ and $M_k = \{m_1, m_2, \dots, m_k\}$ for an integer $k \geq 1$. T_k is a set of $k + 1$ unique timestamps in an increasing order ($t_0 < t_1 < \dots < t_k$). M_k is a set of k attribute values, where m_i ($1 \leq i \leq k$) is a maximum attribute value associated with a time interval $[t_{i-1}, t_i]$ if there exist at least one live tuple in $[t_{i-1}, t_i]$. Otherwise, $m_i = \text{nil}$ for an empty interval. No two consecutive values in M_k are equal (*i.e.*, $m_i \neq m_{i+1}$ for any i ($1 \leq i \leq k - 1$)). Each tuple t in an input database can be considered as a (T_1, M_1) with $T_1 = \{t.start_time, t.end_time\}$ and $M_1 = \{t.attribute_value\}$.

Figure 4 illustrates the process of merging the tuples of the sample `EMPLOYEES` table in Figure 1(a). The sample tuples are described as four line seg-

ments in Figure 4(a). In the first step, the first two tuples in the `EMPLOYEES` table are merged into an intermediate result $(\{8, 18, 31\}, \{45000, 46000\})$; the last two tuples are merged into an intermediate result $(\{7, 12, 18, 21\}, \{35000, \text{nil}, 38000\})$. The result of the first step is shown in Figure 4(b). In the second step, the two intermediate results are merged together into the final aggregate result $(\{7, 8, 18, 31\}, \{35000, 45000, 46000\})$, as shown in Figure 4(c).

As an input database of N tuples is scanned, the merge-sort aggregation algorithm generates $\lceil N/2 \rceil$ first-step intermediate aggregates in memory. Then, the algorithm recursively merges the intermediate results until a final aggregate result is obtained. Thus, the worst case processing time of the algorithm is $\mathcal{O}(N \log N)$. As is shown in Figure 4, the size of an intermediate result (T_k, M_k) may be smaller than the tuples themselves covered by (T_k, M_k) , because two consecutive intervals can be merged into a single interval if they share the same aggregate value (*i.e.*, maximum in the example). Thus, the amount of additional memory required for intermediate results is likely to be smaller than the size of an input database. Nonetheless, for `COUNT` aggregations, the balanced tree will remain as the algorithm of choice. This is because the balanced tree algorithm will keep the memory requirement (*i.e.*, the number of tree nodes) down to the minimum by building a balanced tree incrementally and by removing repeated timestamps, and thereby minimizing its processing time.

4. Bucket Algorithm for Large-Scale Aggregation

In addition to the algorithms for small-scale aggregations proposed in the previous section, another major component of the work proposed in this paper is to develop new techniques for computing temporal aggregates under the constraint of limited buffer space. Then, the size of databases we can deal with is not limited by the size of available memory. Additionally, it is crucial that temporal aggregation require only a constant number (say, two or three) of database scans, due to potentially huge amount of temporal data. It will be prohibitively costly for a large-scale database, if the number of required database scans is not limited and is rather proportional to the size of database. For this reason, we do not consider as an acceptable solution any method that requires more than a small constant number of database scans.

In this section, we propose a new algorithm based on partitioning database tuples into several buckets, which has been used for many important database operations such as the relational hash join algorithm. Although

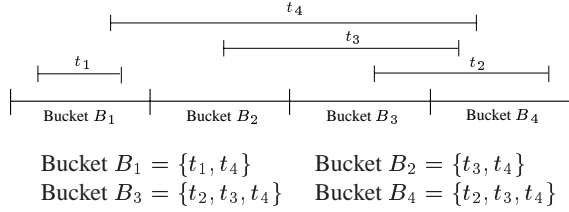


Figure 5. Time-line partitioning and assignment of tuples into buckets.

the idea of data partitioning appears promising for relational hash join operation, it cannot be applied directly to temporal aggregation. Tuples associated with time intervals are not readily partitioned into temporally disjoint equivalence classes (*e.g.*, hash buckets), because the time intervals of tuples may be of any length. Some tuples may overlap with the intervals of more than one bucket, and such tuples must be checked with tuples in all the overlapping buckets. That is, there is no guarantee that temporal aggregates can be computed by reading the buckets only a constant number of times.

To circumvent this problem, one can allow assignment of a data object into multiple buckets by replicating it. This approach can be best described by an example given in Figure 5. The time-line of a given temporal database is partitioned into \mathcal{N}_B disjoint intervals, where \mathcal{N}_B is the number of buckets. If a tuple's life span is contained in the interval of a bucket, the tuple is assigned to the bucket. For example, in Figure 5, tuple t_1 will be assigned to bucket B_1 as t_1 's life span is properly contained in that of bucket B_1 . On the other hand, if a tuple's life span overlaps two or more intervals (say, k intervals), the tuple's life span is split into k pieces and these pieces may be assigned to k buckets. (It turns out that splitting a tuple into several does not impact the result of the aggregation.) In Figure 5, the life spans of tuples t_2 , t_3 and t_4 overlap with 2, 3 and 4 buckets, respectively. Thus, tuple t_2 will be assigned to buckets B_3 and B_4 , t_3 to buckets B_2 , B_3 and B_4 , and t_4 to buckets B_1 , B_2 , B_3 and B_4 .

This process entails replicating tuples and may lead to considerable duplication of data, especially for long-lived tuples. To minimize duplication of tuples, we propose to assign each tuple solely to the buckets where the tuple's start and end timestamps lie. Suppose the life span of a tuple t overlaps buckets B_i, B_{i+1}, \dots, B_j ($0 \leq i < j < \mathcal{N}_B$). Then, the tuple t will be replicated only in the buckets B_i and B_j , but the intermediate buckets will not store the tuple t . Instead, a *meta array* is used to aggregate the information that the tuple t 's life span overlaps the intermediate buckets B_{i+1}, \dots, B_{j-1} . The size of a meta array is equal to the number of buckets. The i -th element of a meta array stores an aggregate value (*e.g.*, count) for the i -th bucket.

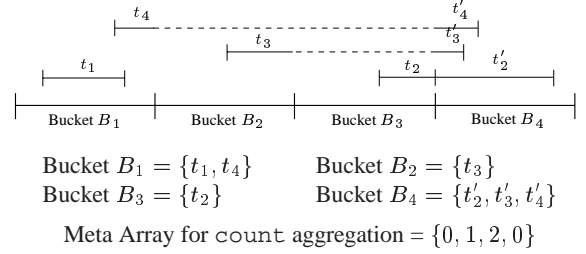


Figure 6. Meta Array and Reduced Data Replication

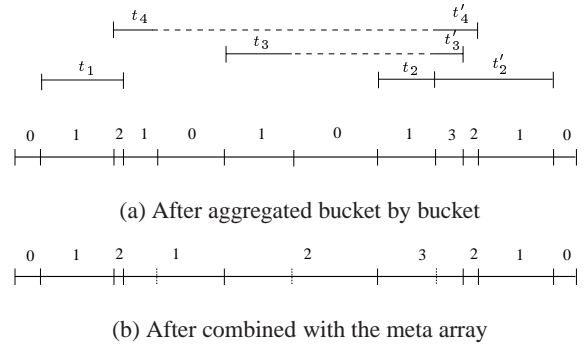


Figure 7. Steps of the aggregation based on data partitioning and meta array

For example, in Figure 6, the time interval of tuple t_3 spans over three buckets B_2 , B_3 and B_4 . Thus, t_3 is split into two segments (*i.e.*, t_3 and t'_3) with adjusted time intervals so that each segment can be properly contained in the interval of its corresponding bucket. (Solid lines in Figure 6 represent adjusted time intervals of split tuples.) Then, t_3 and t'_3 are assigned to two buckets B_2 and B_4 , respectively; the third element of the meta array is incremented by one. In a similar way, t_4 and t'_4 are assigned to two buckets B_1 and B_4 respectively, and the second and third elements of the meta array are incremented by one. The resulting data partitioning and meta array are illustrated in Figure 6. Note that neither the first nor the last element of the meta array stores a valid aggregate value, as no tuple can have a life span longer than the time-line of an entire database.

Once all the tuples are scanned and partitioned into buckets and a meta array is created, the temporal aggregate operation can be performed on each bucket independently. Figure 7(a) shows the partial results of the aggregation performed on each bucket. Then, each aggregate value stored in the meta array is combined with the aggregation results from each corresponding bucket (*e.g.*, simply by adding counts for count aggregation). Lastly, the final aggregation results can be obtained by merging each pair of adjacent buckets at their boundaries if the two adjacent aggregate values are equal. Fig-

Algorithm 2: Temporal Bucketization

```
set  $\mathcal{I}_B \leftarrow$  time interval for each bucket  $((\mathcal{T}_{max} - \mathcal{T}_{min})/\mathcal{N}_B)$ ;  
foreach tuple  $t$  in a table do  
  set start_bucket  $\leftarrow (t.start\_time - \mathcal{T}_{min})/\mathcal{I}_B$ ;  
  set end_bucket  $\leftarrow (t.end\_time - \mathcal{T}_{min})/\mathcal{I}_B$ ;  
  insert  $t$  into a bucket  $B_{start\_bucket}$ ;  
  if (start_bucket  $\neq$  end_bucket) then  
    insert  $t'$  into a bucket  $B_{end\_bucket}$ ;  
  for (i=start_bucket+1 to end_bucket-1) do  
    update meta_array[i];  
for (i=0 to  $\mathcal{N}_B - 1$ ) do  
  perform temporal aggregation on the bucket  $B_i$ ;  
  combine the scalar value of meta_array[i] to the bucket  $B_i$ ;  
  merge the bucket boundary with  $B_{i-1}$  as needed;
```

ure 7(b) shows the final aggregation results. The dotted vertical bars in the figure represent the merged bucket boundaries. Algorithm 2 outlines the proposed temporal aggregation algorithm based on data partitioning. In the algorithm description, it is assumed that the entire time-line of a table is partitioned into \mathcal{N}_B disjoint intervals of an equal length, each of which is associated with a bucket. Note that any small-scale aggregation algorithm proposed in the previous section can be used to aggregate each individual bucket.

Provided that the meta array is small enough to fit in memory and sufficient memory is available to hold all the tuples in a bucket, the temporal aggregate operation can be performed by reading each bucket just once. Thus, in total, this approach requires three database accesses (*i.e.*, two reads and one write) to compute temporal aggregates. Considering the data replication for the tuples overlapped with multiple buckets, the database access requirement of this approach is likely to increase to some extent depending on various factors such as the life spans of tuples and the number of buckets used. Even in the worst case, however, the size of a given table can increase only up to twice its original size by replicating each tuple in the table into two buckets. Thus, the database access requirement of this approach is still bounded to a small constant number of scans. We will show the performance impact of data replication in Section 6.

5. Parallel Bucket Algorithm

Previous attempts [7, 15] to develop scalable methods for computing large-scale temporal aggregates were based on parallelizing the aggregation tree algorithm. Consequently, those approaches inherit all the limitations the aggregation tree algorithm has. Specifically, these approaches will suffer from $\mathcal{O}(N^2)$ worst-case running time and tight limitations on a database size they can deal with. As a solution for this worst-case running time, we propose a new parallel temporal aggregation algorithm based on the bucket algorithm (Algorithm 2) presented above.

Algorithm 3: Parallel Temporal Bucketization

```
set  $\mathcal{P} \leftarrow$  number of participating processors;  
set  $\mathcal{I}_B \leftarrow$  time interval for each bucket  $((\mathcal{T}_{max} - \mathcal{T}_{min})/(\mathcal{N}_B \times \mathcal{P}))$ ;  
set this_proc  $\leftarrow$  a local processor id ( $0 \leq \text{this\_proc} < \mathcal{P}$ );  
foreach tuple  $t$  in a local partition or from a remote processor do  
  set start_proc  $\leftarrow (t.start\_time - \mathcal{T}_{min})/(\mathcal{I}_B \times \mathcal{P})$ ;  
  set end_proc  $\leftarrow (t.end\_time - \mathcal{T}_{min})/(\mathcal{I}_B \times \mathcal{P})$ ;  
  if (start_proc  $\neq$  this_proc) then  
    send  $t$  to a processor  $P_{start\_proc}$ ;  
  if (end_proc  $\neq$  this_proc) then  
    send  $t'$  to a processor  $P_{end\_proc}$ ;  
  for (i=start_proc+1 to end_proc-1) do  
    update global_meta_array[i];  
  insert  $t$  into one or two local buckets as in Algorithm 2;  
  update local_meta_array as in Algorithm 2;
```

```
Globally combine the global_meta_array wrt. an aggregate operator  $op$ ;  
for (i=0 to  $\mathcal{N}_B - 1$ ) do  
  local_meta_array[i]  $\leftarrow$   
   $op(\text{local\_meta\_array}[i], \text{global\_meta\_array}[\text{this\_proc}])$ ;  
  perform temporal aggregation on the bucket  $B_i$   
  with local_meta_array[i] as in Algorithm 2;
```

It is relatively straightforward to parallelize the bucket algorithm by distributing buckets across participating processors. The time-line of a given temporal database is partitioned into \mathcal{P} disjoint intervals, where \mathcal{P} is the number of processors. Then, on each processor, the time-line of its local database is again partitioned into \mathcal{N}_B disjoint intervals, where \mathcal{N}_B is the number of local buckets. We propose to use a *global* meta array and a *local* meta array on each processor. The global meta array keeps track of tuples whose life spans are extended over the time-lines of more than one processors. The local meta array, on the other hand, keeps track of tuples whose life spans are extended over more than one local buckets.

The proposed parallel aggregation algorithm is summarized in Algorithm 3. As mentioned above, it is assumed that the entire time-line of a table is partitioned into $\mathcal{N}_B \times \mathcal{P}$ disjoint intervals of an equal length, each of which is associated with a bucket, and the buckets are distributed across \mathcal{P} processors by range partitioning so that each processor is assigned \mathcal{N}_B consecutive buckets. This range partitioning scheme obviously minimizes the size of a global meta array in a way that only one array element is required per each processor. Since each processor computes a global meta array independently only for its local data, all the \mathcal{P} processors need to communicate each other to compute a final global meta array for an entire database with respect to a given operator op . The operator op is determined by a kind of aggregate operation. For example, op will be an *addition* operator for a count aggregation and a *maximum* operator for a max aggregation. Such collective communication for computing a final global meta array can be implemented efficiently on most parallel computers and networks of workstations [1]. Thus the overhead for combining global meta arrays is expected to be negligible because the volume of communication is only \mathcal{P} words per pro-

cessor. A complete description of the parallel bucket algorithm and its empirical evaluation are given in [11].

6. Empirical Evaluation

In this section, we evaluate the proposed algorithms empirically and compare with the previous work. We chose the `count` and `max` temporal aggregates to carry out experiments under various operational conditions that may affect the performance of the algorithms. In particular, we focus on the performance gain by the proposed algorithms for small-scale aggregations, and the scalability of the bucket algorithm.

6.1. Experimental Settings

Testing and benchmarks were performed on Intel Pentium workstations with 200 MHz clock rate. Each workstation has 128 MBytes of memory and 2 or 4 GBytes of disk storage with Ultra-wide SCSI interface, and runs on Linux kernel version 2.0.30. Throughout the experiments, we measured elapsed times including disk access time. For accurate measurement, we averaged elapsed times from multiple runs after eliminating extreme cases. Additionally, we avoided the system cache effects for disk accesses by loading irrelevant data into the entire memory between consecutive runs of our experiments.

We generated synthetic data in the same way as in [9]. Each database has a time-line of one million temporal instants. We considered two basic life spans for tuples: short-lived and long-lived. The life span of a short-lived tuple was determined randomly between one and 1,000 instants; the life span of a long-lived tuple was determined randomly between 200,000 and 800,000 instants, namely, between 20 and 80 percent of the time-line of a database. In most of our experiments, the population of long-lived tuples was fixed at 10 percent or 30 percent. The start times of tuples were uniformly distributed over the time-line of a database. Each tuple was 20 bytes including two temporal attributes (start time and end time) and other non-temporal attributes as well. Synthetically generated databases used in our experiments were not sorted by any temporal attribute unless stated otherwise.

6.2. Small-Scale Aggregation

The first set of experiments were carried out on relatively small databases between 1 MBytes and 20 MBytes so that all the required data structures can fit in available memory. Recall that the algorithms proposed in Section 3 as well as the aggregation tree algorithm and its variation require that the entire data structures be kept in memory. In this section, we used the *balanced tree* algorithm for `count` aggregations, and the *merge-sort aggregation* algorithm for `max` aggregations.

Figure 8(a) compares the balanced tree and aggregation tree algorithms for `count` aggregations; Figure 8(b) compares the merge-sort and aggregation tree algorithms for `max` aggregations. The proposed balanced tree and merge-sort aggregation algorithms consistently performed about twice faster than the aggregation tree algorithm for `count` and `max` aggregations, respectively. While the aggregation tree took more time to aggregate a database with higher percentage of long-lived tuples, the processing times of the two proposed algorithms remained constant for different percentage of long-lived tuples. Note that the performance of the aggregation tree algorithm remains unchanged for `count` and `max` aggregations, since the algorithm works essentially in the same way for both the aggregations.

In Figures 8(c) and 8(d), the tuples in input databases were sorted by their start time, where we expected the worst-case performance from the aggregation tree algorithm. The processing times of the aggregation tree were several orders of magnitude slower than the two proposed algorithms, and were plotted as almost vertical lines in the figures. Thus, we compared with the k -ordered aggregation tree algorithm (with $k = 1$) instead. The proposed algorithm still performed two to three times faster than the k -ordered aggregation tree algorithm.

In summary, the proposed algorithms outperformed the aggregation tree and k -ordered aggregation tree consistently by a significant margin. The k -ordered aggregation tree requires a priori knowledge about the orderedness of databases, whereas the proposed algorithms do not.

6.3. Bucket Algorithm for Large-Scale Aggregation

Despite the fact that the balanced tree and merge-sort aggregation algorithms were designed for two different groups of aggregate operations, both algorithms showed almost identical performance behaviors in the previous experiments. Thus, for the rest of this section, we present experimental results only for `count` aggregations.

The second set of experiments were carried out to evaluate the *bucket* algorithm proposed in Section 4. First, we performed aggregations with and without data partitioning for small databases, so that we could measure the overhead of data partitioning. The balanced tree algorithm was used to compute `count` aggregates. In Figure 9(a), we used 64 buckets irrespective of database sizes, which was large enough to demonstrate the overhead of data partitioning. Compared with the balanced tree algorithm without data partitioning, we observed about 10 to 30 percent increase in processing time of the bucket algorithm. Despite the additional overhead,

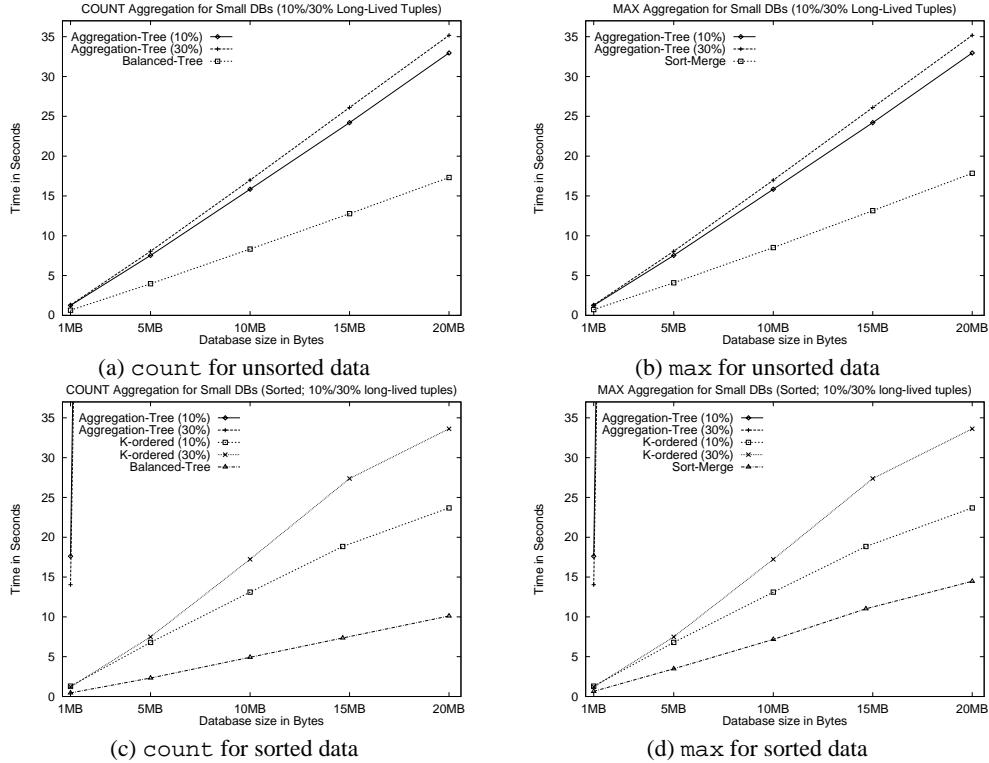


Figure 8. Aggregation time for small-scale databases

however, the bucket algorithm still outperformed the aggregation tree algorithm significantly. (Compare Figure 8(a) and Figure 9(a).)

For small databases, the amount of overhead of data partitioning was expected to be smaller than what it should be for large databases, because all the buckets might remain in memory even after they were written to disk. So, for the next step of aggregating individual buckets, the cached buckets would be used instead of the disk copies. Also note that performance of the bucket algorithm is affected by the percentage of long-lived tuples. The reason appears quite obvious because long-lived tuples are more likely to be replicated than short-lived tuples, leading to increased computation time and disk access time.

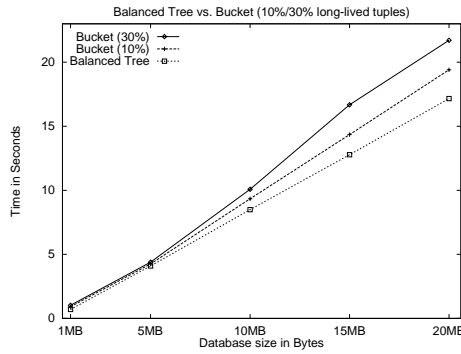
Figure 9(b) shows processing times of the bucket algorithm for databases of size from 20 MBytes up to 1 G-Bytes. The number of buckets used for data partitioning was 2, 8, 16, 24, 32 and 40 for 20 MBytes, 200 MBytes, 400 MBytes, 600 MBytes, 800 MBytes and 1 GBytes databases, respectively. Since each of these databases is too large to fit in memory (with an exception of a 20 M-Byte database), none of the small-scale aggregation algorithms could be used for this experiment. The results shown in Figure 9(b) demonstrate that the proposed bucket algorithm can compute temporal aggregates for databases substantially larger than the size of available

memory. However, it should be noted that the processing time of the algorithm grows faster than linearly as the size of a database increases. This clearly motivates the need of scalable solutions such as the parallel bucket algorithm we proposed in Section 5. Readers are referred to [11] for complete description of the empirical evaluation of the parallel bucket algorithm.

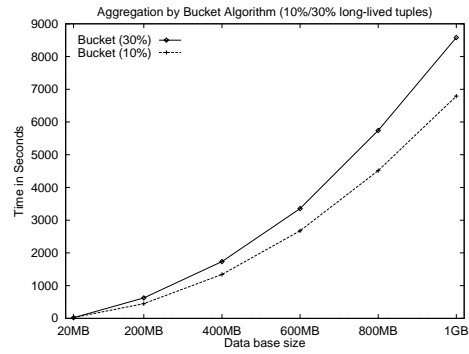
7. Conclusions and Future Work

We have developed new algorithms for computing temporal aggregates. The proposed algorithms provide significant benefits over the current state of the art in different ways. The balanced tree and merge-sort aggregation algorithms have improved the worst-case and average-case processing time significantly for small databases that fit in memory. We have also developed a new sequential bucket algorithm based on novel data partitioning schemes. This algorithm can be used to compute temporal aggregates for databases that are substantially larger than the size of available memory, by processing data partitions in a sequential or parallel fashion.

From our experiments, we have observed that there are a few factors that affect the performance. They include the percentage of long-lived tuples and the num-



(a) Bucket vs. Balanced Tree



(b) Scaleup of Bucket algorithm

Figure 9. Aggregation time for large-scale databases

ber of buckets used for data partitioning. Although the proposed algorithms outperformed previous approaches consistently irrespective of such conditions, we believe it is worth elaborating further on the issues. Additionally, we plan to study performance impacts of such factors as initial data placement (*e.g.*, temporal partitioning vs. non-temporal partitioning) and data reduction by aggregation.

We also plan to extend the data partitioning approach to spatio-temporal databases, which requires computing aggregates for data objects with two or more dimensional extents. Unlike the temporal aggregation, we expect that the process of data partitioning and generating meta arrays will be more sophisticated.

References

- [1] M. Barnett, S. Gupta, D. Payne, L. Shuler R. van de Geijn, and J. Watts. Interprocessor collective communication library (InterCom). pages 357–364, Knoxville, TN, May 1994.
- [2] Jon Louis Bentley. Algorithms for Klee’s rectangle problems. Technical Report unpublished, Pittsburgh, PA, 1977.
- [3] Surajit Chaudhuri and Umeshwar Dayal. An overview of data warehousing and OLAP technology. *SIGMOD Record*, 26(1), March 1997.
- [4] Thomas H. Cormen, Charles E. Leiserson, and Ronald L. Rivest. *Introduction to Algorithms*. McGraw Hill, Cambridge, Mass., 1990.
- [5] Robert Epstein. Techniques for processing of aggregates in relational database systems. Technical Report UCB/ERL M7918, University of California, Berkeley, CA, February 1979.
- [6] Johann C. Freytag and Nathan Goodman. Translating aggregate queries into iterative programs. In *Proceedings of the 12th VLDB Conference*, pages 138–146, Kyoto, Japan, August 1986.
- [7] Jose Alvin G. Gendrano, Bruce C. Huang, Jim M. Rodrigue, Bongki Moon, and Richard T. Snodgrass. Parallel algorithms for computing temporal aggregates. In *Proceedings of the 15th Inter. Conference on Data Engineering*, Sydney, Australia, March 1999.
- [8] Christian S. Jensen and Richard T. Snodgrass. Semantics of time-varying information. 21(4):311–352, 1996.
- [9] Nick Kline and Richard T. Snodgrass. Computing temporal aggregates. In *Proceedings of the 11th Inter. Conference on Data Engineering*, pages 222–231, Taipei, Taiwan, March 1995.
- [10] Donald E. Knuth. *Sorting and Searching*, volume 3 of *The Art of Computer Programming*. Addison-Wesley Publishing Company, Reading, Mass., 1973.
- [11] Bongki Moon, Ines Fernando Vega Lopez, and Vijaykumar Immanuel. Scalable algorithms for large-scale temporal aggregation. Technical Report TR 98-11, Tucson, AZ 85721, November 1998. <http://www.cs.arizona.edu/research/reports.html>.
- [12] R. T. Snodgrass, S. Gomez, and E. Mackenzie. Aggregates in the temporal query language TQuel. *IEEE Transactions on Knowledge and Data Engineering*, 5(5):826–842, October 1993.
- [13] A. Tansel et al. *Temporal Databases: Theory, Design and Implementation*. Database Systems and Applications Series. Benjamin/Cummings, Redwood City, CA, 1993.
- [14] Paul A. Tuma. Implementing historical aggregates in TempIS. Master’s thesis, Wayne State University, Detroit, Michigan, November 1992.
- [15] Xinfeng Ye and John A. Keane. Processing temporal aggregates in parallel. In *IEEE Inter. Conf. on Systems, Man, and Cybernetics*, pages 1373–1378, Orlando, FL, October 1997.