

Efficient Algorithms for Large-Scale Temporal Aggregation

Bongki Moon, Ines Fernando Vega Lopez, and Vijaykumar Immanuel

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. Third, the parallel algorithm designed on a shared-nothing architecture achieves scalable performance by delivering nearly linear scale-up and speed-up, even at the presence of data skew. The contributions made in this paper are particularly important because the rate of increase in database size and response time requirements has out-paced advancements in processor and mass storage technology.

Index Terms—Temporal databases, temporal aggregation, scalable query processing, data partitioning, balanced tree algorithm, merge-sort algorithm, temporal query processing, aggregate queries.

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 [11], [22]. In fact, there are several temporal query languages supporting temporal aggregation [20]. 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 values are computed over these groups. In general, temporal grouping is done by two types of partitioning [20]: *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 [13]. 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. Fig. 1a 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 Fig. 1b. 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 [5]. Although several sequential and parallel algorithms have been developed for computing temporal aggregates [10], [12], [13], [14], [20], [23], [25], they suffer from serious limitations such as the restriction on the size of aggregation by available memory and the 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:

- B. Moon and I. Fernando Vega Lopez are with the Department of Computer Science, University of Arizona, Tucson, AZ 85721-0077. E-mail: {bkmooon, ifvega}@cs.arizona.edu.
- V. Immanuel is with Compaq Computer Corporation, 19333 Vallico Pkwy., Cupertino, CA 95014. E-mail: vijay.immanuel@compaq.com.

Manuscript received 16 June 2000; revised 16 May 2001; accepted 22 June 2001.

For information on obtaining reprints of this article, please send e-mail to: tkde@computer.org, and reference IEEECS Log Number 112298.

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)

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)

Fig. 1. Sample database and its temporal aggregation. (a) Input database tuples. (b) Temporal aggregation results.

- 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.
- A parallel algorithm has been developed for shared-nothing architectures for large-scale aggregations. This solution achieves scalable performance by delivering nearly linear scale-up and speed-up, even at the presence of data skew.

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 sequential and parallel solutions. In Section 7, we discuss further extensions to the proposed algorithms. Finally, Section 8 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 [8]. 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 [9].

The first approach for implementing temporal aggregation was proposed by Tuma [23] 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 [13] for temporal aggregation based on instant grouping of tuples. The algorithms are called *aggregation tree* and its variant *k-ordered aggregation tree*, because 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 in depth first search. Kline [14] proposed to use 2-3 tree, which is a balanced tree, to compute temporal aggregates. The leaf nodes of the tree store the time intervals of the aggregate results. Like the aggregation tree, this approach requires only one database scan.

Kim et al. proposed a *point-based aggregation tree (PA-tree)* [12], which stores timestamps instead of intervals in an AVL tree. In addition to timestamps, each node in the *PA-tree* stores either a single aggregate value for count, sum, and average aggregation, or a list of *value-length* pairs for min and max aggregation. Computing count, sum, and average aggregates is performed by doing an in-order traversal of the tree and updating aggregate values by the amount indicated on each encountered node. Min and max aggregates are computed by merging the lists of pairs associated to each tree node in a similar way to the skyline problem [16].

Yang and Widom introduced the *SB-tree* [24] for incrementally computing temporal aggregates using a materialized view approach. This is a disk-based approach that computes temporal aggregates over a base relation which may gradually change by insertion and deletion. The *SB-tree* contains a hierarchy of intervals associated with partially computed aggregates. Aggregation over a given temporal interval is evaluated by performing a depth first

search on the tree and accumulating the partial aggregate values along the way. Note that the SB-tree cannot compute min and max aggregates when deletion operations are allowed. A *multiversion SB-tree* was proposed to deal with temporal aggregates coupled with nontemporal key-value range predicates [26]. The multiversion SB-tree is essentially a series of SB-trees, one for each time instant.

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 this reason, the worst-case time to create an aggregation tree is $\mathcal{O}(N^2)$ for N tuples sorted in time. An 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 $\mathcal{O}(N^2)$.¹

Apparently, the aggregation tree, the most efficient among the aforementioned algorithms, suffers from poor scale-up performance due to the $\mathcal{O}(N^2)$ worst-case running time and memory requirement. Recently, there have been some research efforts to develop parallel algorithms for computing temporal aggregates for large-scale databases. Ye and Keane proposed two approaches to parallelize the aggregation tree algorithm on a shared-memory architecture [25]. Gendrano et al. have also developed several new parallel algorithms [10] for computing temporal aggregates, specifically on a shared-nothing architecture, by parallelizing the aggregation tree algorithm. Gendrano et al. showed promising scale-up performance of the parallel algorithms through extensive empirical studies under various conditions. Nonetheless, all the aforementioned parallel algorithms inherit the same limitations from the aggregation tree algorithm, as the

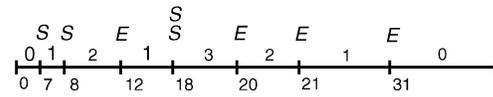


Fig. 2. Example of *count* aggregation by sorting timestamps and tags.

parallel algorithms were developed by parallelizing the aggregation tree. In particular, the size of the database those parallel algorithms can handle will be limited by the aggregate memory of participating processors.

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 [13]. 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 columns. 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*, and *avg* in one group, and *max* and *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 *merge-sort aggregation* algorithm, which is similar to the classical merge-sort algorithm [15]. 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 *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 Fig. 2 for

1. If the 2-3 tree approach is used, its running time will be $\mathcal{O}(N \log N)$, because 2-3 tree is a balanced tree. However, its main limitation lies on the requirement that a database be initially sorted by start timestamps. It has been shown that, for a randomly ordered database, the aggregation tree performs better than the 2-3 approach [14]. This is due to the preprocessing cost required by the 2-3 tree approach to sort the database.

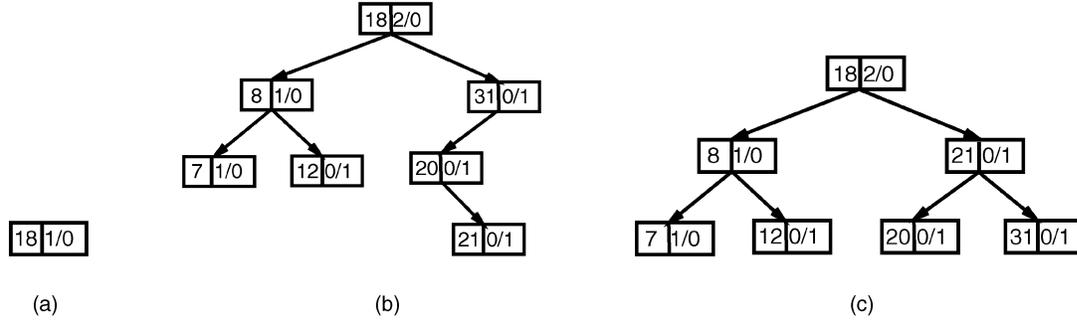


Fig. 3. Example of balanced tree construction.

a sorted list of timestamps and tags for the sample database given in Fig. 1a.

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 Fig. 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 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.² Additionally, a color tag is stored in each node, as we use the *red-black* insertion algorithm [6] to keep the tree balanced dynamically.

Fig. 3 shows the process of building a balanced tree for the sample `Employees` table in Fig. 1a. 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 Fig. 3a. Figs. 3b and 3c 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 decremented by the end-counter value of the node. The proposed balanced tree algorithm is summarized in Fig. 4.

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 Fig. 3b, 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.

2. For 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;
for each tuple  $t$  in a table do begin
  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  $\mathcal{T}$ ;
  endif
  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  $\mathcal{T}$ ;
  endif
end
set  $count \leftarrow 0$ ;
for each node  $n$  in  $\mathcal{T}$  traversed by in-order do begin
   $count += n.no\_starts$ ;
  output  $n.ts$  and  $count$ ;
   $count -= n.no\_ends$ ;
end
end Algorithm

```

Fig. 4. Balanced tree algorithm for `count` aggregation.

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 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 nontrivial.

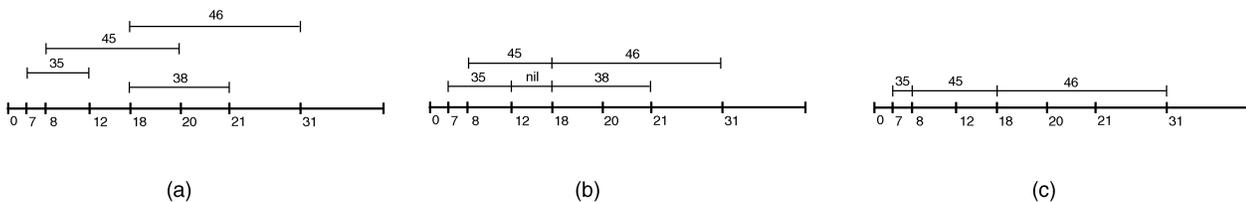
Instead, we propose a *bottom-up* aggregation approach, which we call a *merge-sort aggregation* algorithm. Like the classical merge-sort algorithm based on the divide-and-conquer strategy, the merge-sort aggregation algorithm computes a larger (intermediate) aggregate result by merging two smaller (intermediate) 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, $m_i = nil$ 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 $M_1 = \{t.attribute_value\}$.

Fig. 5 illustrates the process of merging the tuples of the sample `Employees` table in Fig. 1a. The sample tuples are described as four line segments in Fig. 5a. 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 $(\{7, 12, 18, 21\}, \{35000, nil, 38000\})$. The result of the first step is shown in Fig. 5b. 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 Fig. 5c.

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 Fig. 5, 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

Fig. 5. Example of merging for `max` aggregation. (a) Input records, (b) after the first merging step, and (c) after the second merging step.

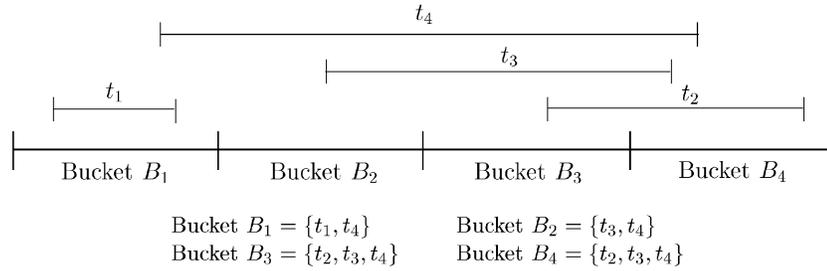


Fig. 6. Time-line partitioning and assignment of tuples into buckets.

memory required for intermediate results is likely to be smaller than the size of an input database. Although the bottom-up aggregation is also applicable to 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 requires only a constant number (say, two or three) of database scans, due to a 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. The idea of the hash join algorithm is to hash two joining relations on the join attribute, using the same hash function. Then, it is assured that tuples of one relation in a bucket can join only with tuples of the other relation in the same bucket. Thus, once both relations are partitioned, the join operation can be performed by reading the relations just once, provided that enough memory is available to keep all the tuples of one relation in a bucket in memory. Assuming uniform distribution of data, it has been shown that the hash join algorithm requires three database scans if the number of buffer pages is larger than the square root of the number of disk pages in a smaller relation [7].

Although 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 buckets, 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 Fig. 6. 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 Fig. 6, 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 Fig. 6, 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.

For example, in Fig. 7, 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 Fig. 7 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 updated. 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

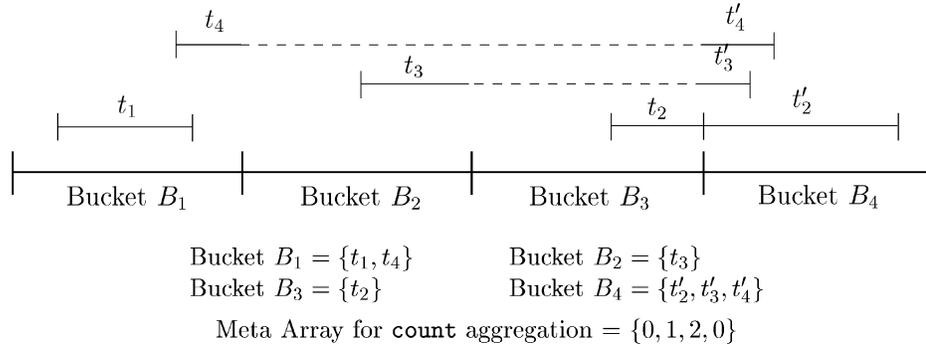


Fig. 7. Meta array and reduced data replication.

updated. The resulting data partitioning and meta array are illustrated in Fig. 7. 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. Fig. 8a 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. 8b shows the final aggregation results. The dotted vertical bars in the figure represent the merged bucket boundaries. Fig. 9 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 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

Parallel processing for database applications typically involves partitioning of data, followed by allocation of the partitions to a set of processors. Then, the processors perform operations on the partitioned data in parallel, achieving speed-up in query processing times. Among the various architectures that have been proposed for parallel database systems, a *shared-nothing* architecture [21] has made it an attractive choice for large-scale database applications due to its high potential for scalability. By scalability, we mean the capability of delivering an increase

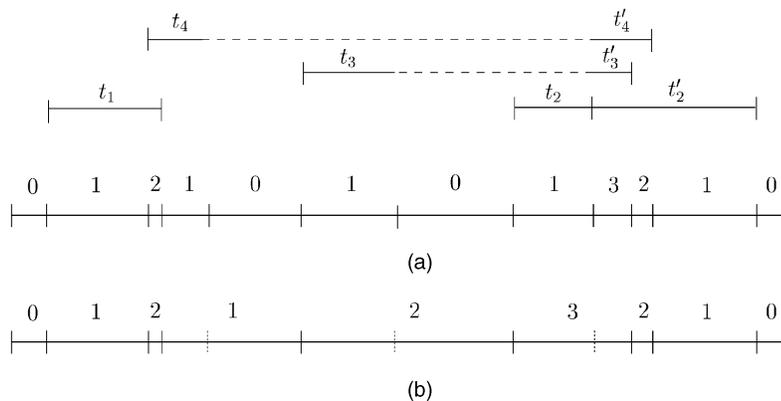


Fig. 8. Steps of the aggregation based on data partitioning and meta array. (a) After aggregated bucket by bucket. (b) After combined with a meta array.

Algorithm 2 *Temporal Bucketization*

```

set  $\mathcal{I}_B \leftarrow$  time interval for each bucket  $((\mathcal{T}_{max} - \mathcal{T}_{min})/\mathcal{N}_B)$ ;
for each tuple  $\mathbf{t}$  in a table do begin
  set  $start\_bucket \leftarrow (\mathbf{t}.start\_time - \mathcal{T}_{min})/\mathcal{I}_B$ ;
  set  $end\_bucket \leftarrow (\mathbf{t}.end\_time - \mathcal{T}_{min})/\mathcal{I}_B$ ;
  insert  $\mathbf{t}$  into a bucket  $B_{start\_bucket}$ ;
  if ( $start\_bucket \neq end\_bucket$ ) insert  $\mathbf{t}'$  into a bucket  $B_{end\_bucket}$ ;
  for ( $i=start\_bucket+1$  to  $end\_bucket-1$ ) do update  $meta\_array[i]$ ;
end
for ( $i=0$  to  $\mathcal{N}_B - 1$ ) do begin
  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;
end
end Algorithm

```

Fig. 9. Bucket algorithm based on temporal data partitioning.

in performance proportional to an increase in the number of participating processors.

In a shared-nothing architecture, each processor owns local memory and secondary storage units, and communicates with each other by message passing. Initial data placement can be either centralized or distributed across multiple processors. For most of the parallel database operations, however, some of the data may have to be redistributed among processors that actually participate in the operations. We assume that resulting aggregates remain in local storage units of the participating processors without collecting the results on a special coordinator processor. Then, the resulting aggregates can be used as intermediate data for the next phase of parallel query processing.

As was pointed out in Section 2, most of the previous attempts to develop scalable methods for computing large-scale temporal aggregates were based on parallelizing the aggregation tree algorithm. For this reason, 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.

In this section, we propose a new parallel temporal aggregation algorithm based on the bucket algorithm (Algorithm 2) presented in the previous section. 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 the processor is again partitioned into \mathcal{N}_B disjoint intervals. However, distributing the buckets is not enough to compute correct aggregate results, because the construction of meta arrays must also be processed in parallel in an efficient way.

We propose to use a *local meta array* and a *global meta array* on each processor for tuples whose life spans overlap time-lines of multiple local buckets and multiple processors, respectively. Specifically, if the life span of a tuple t overlaps the k th bucket ($B_{P_i,k}$) of processor P_i through the l th bucket ($B_{P_j,l}$) of processor P_j , the tuple t will be replicated only in $B_{P_i,k}$ and $B_{P_j,l}$. Then, a local

meta array of P_i is used to aggregate the information that the tuple t 's life span overlaps the intermediate buckets $B_{P_i,k+1}, \dots, B_{P_i,N_B}$, and so is a local meta array of P_j for $B_{P_j,1}, \dots, B_{P_j,l-1}$. Finally, a global meta array is updated on a processor that owns the tuple t to inform the intermediate processors P_{i+1}, \dots, P_{j-1} of the existence of the tuple t overlapped with their time-lines. The size of a global meta array is equal to the number of processors. The i th element of a global meta array stores an aggregate value (e.g., count) for the i th processor. Local meta arrays are identical with the ones used for the sequential bucket algorithm. Each processor computes its own global and local meta arrays independently.

In Fig. 10, for example, suppose that tuples t_1, \dots, t_4 are initially stored on processor P_0 , and four processors P_0, \dots, P_3 participate in a count aggregation. Since the time interval of t_3 spans over three remote processors P_1, P_2 , and P_3 , t_3 is split into two segments, t_3 and t'_3 , which are then sent to the processors P_1 and P_3 , respectively. Then, the third element of the global meta array of P_0 is incremented by one. In a similar way, t_4 is assigned to P_0 's local bucket B_{03} , and t'_4 is sent to processor P_3 ; the second and third elements of P_0 's global meta array are incremented by one. Fig. 10 shows the resulting data distribution across P_0 's local buckets, data shipping to other processors, and P_0 's local and global meta arrays. Note that, if data is initially distributed across multiple processors, there may be some tuples sent from other processors to P_0 , but they are not shown in Fig. 10.

The proposed parallel aggregation algorithm is summarized in Fig. 11. In the algorithm description, 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

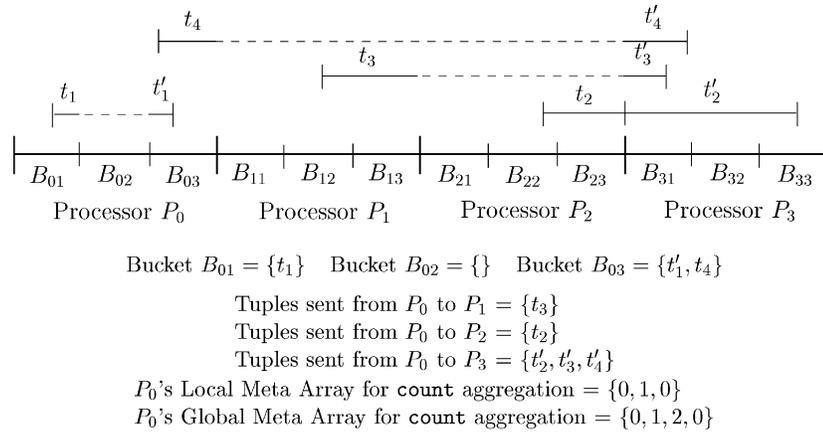


Fig. 10. Data distribution and meta arrays for P_0 's local data.

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 processor.

5.1 Handling Data Skew

When describing the sequential and parallel bucket algorithms, we have assumed that the time-line of a temporal database is partitioned into \mathcal{N}_B (or $\mathcal{N}_B \times \mathcal{P}$) disjoint buckets of an equal length. Apparently, this simple data partitioning scheme will not work for a database with skewed distribution of temporal attribute values, suffering from load imbalance. An interesting question that arises with respect to applying those bucket algorithms is how we ensure that 1) each processor receives about the same number of tuples to achieve load balance among processors, and 2) each bucket receives about the same number of tuples to minimize the local computational requirement (refer to the discussion in Section 6.3).

We address the issue of data partitioning for skewed data by determining the size of each bucket *adaptively*, utilizing the selectivity estimation mechanism provided by most database systems. For example, if frequency distribution of temporal attribute values is provided in an equi-depth histogram [17], then, the partitioned time-line of a processor P_i will be determined as an interval $[h[i \times \mathcal{N}_H/\mathcal{P}], h[(i+1) \times \mathcal{N}_H/\mathcal{P}]]$, where $h[0..\mathcal{N}_H]$ is a set of temporal attribute values of the equi-depth histogram. The boundaries for local buckets can be computed independently by each processor in a similar fashion.

6 EMPIRICAL EVALUATION

In this section, we evaluate the proposed algorithms empirically and compare with the previous work. We chose a count temporal aggregation and carried 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 sequential and parallel bucket algorithms.

6.1 Experimental Settings

Testing and benchmarks were performed on a cluster of 64 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. The workstations are connected by a 100 Mbps switched Ethernet network. The switch can handle an aggregate bandwidth of 2.4 Gbps in an all-to-all type communication. For message passing between the Pentium workstations, we used the LAM implementation of the MPI communication standard [3]. With the LAM message passing package on the Pentium cluster, we observed an average communication latency of 790 microseconds and an average transfer rate of about 5 Mbytes/second. Note that this is relatively high latency and low transfer rate compared with parallel computers equipped with high-performance switches such as IBM SP-2 parallel systems.³

For both sequential and parallel implementations, the same buffer size of 4 Kbytes was used for disk IO and message passing. Nonblocking message passing primitives were used in an attempt to minimize communication overhead by allowing interprocessor communication to be overlapped with local computation and disk IO. Throughout the experiments, we measured elapsed times including disk access time and communication overhead. 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 [13]. Each database has a time-line of one million temporal instants. We considered two basic life spans for tuples:

3. On an SP-2 system with a proprietary MPI implementation `mpif`, we observed an average communication latency of 55 microseconds and an average transfer rate of about 35 Mbytes/second.

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  $(0 \leq \text{this\_proc} < \mathcal{P})$ ;

for each tuple t in a local partition or from a remote processor do begin
  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;
end
Globally combine the global_meta_array wrt. an aggregate operator op;
for (i=0 to  $\mathcal{N}_B - 1$ ) do begin
  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 Algo-
gorithm 2;
end
end Algorithm

```

Fig. 11. Parallel bucket algorithm based on temporal data partitioning.

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 nontemporal 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.

Fig. 12a compares the balanced tree and aggregation tree algorithms for count aggregations. Fig. 12b compares the merge-sort and aggregation tree algorithms for max aggregations. The proposed balanced tree and merge-sort aggregation algorithms consistently performed about twice as fast as 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 Figs. 12c and 12d, 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 algorithms still performed two to three times faster than the *k-ordered aggregation tree* algorithm.

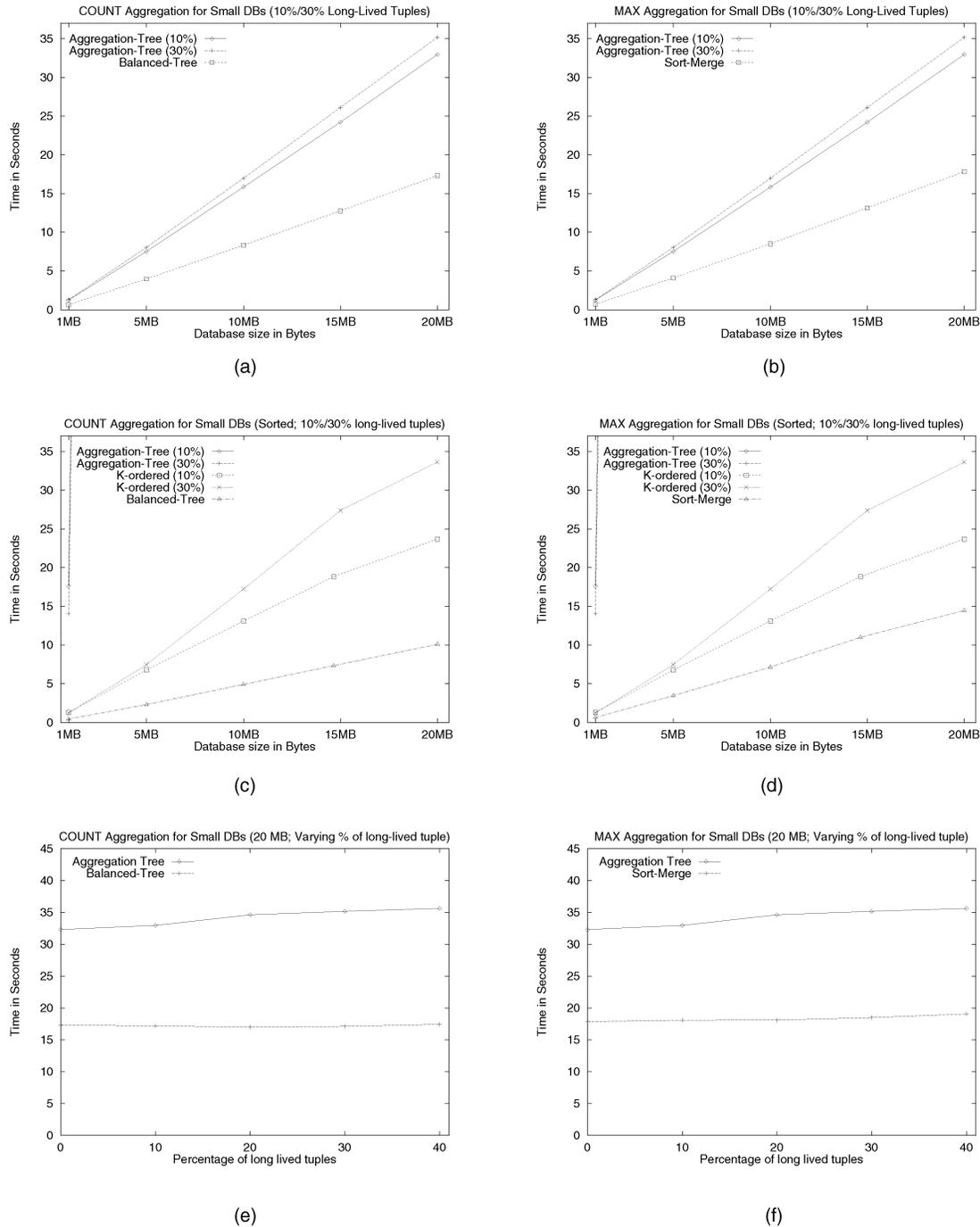


Fig. 12. Aggregation time for small-scale databases. (a) count aggregation for unsorted data, (b) max aggregation for unsorted data, (c) count aggregation for sorted data, (d) max aggregation for sorted data, (e) count aggregation with varying long-lived, and (f) max aggregation with varying long-lived.

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. The performance of the proposed algorithms was not affected by the percentage of long-lived tuples, as Figs. 12e and 12f show the processing times measured on databases (20 MB) with a varying percentage of long-lived tuples.

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. This can be explained by the fact that both algorithms have a running time of $\mathcal{O}(N \log N)$ and minimal memory requirements. In fact, the choice of a small scale aggregation algorithm is orthogonal to data partitioning. Thus, any of the algorithms proposed in this paper, and

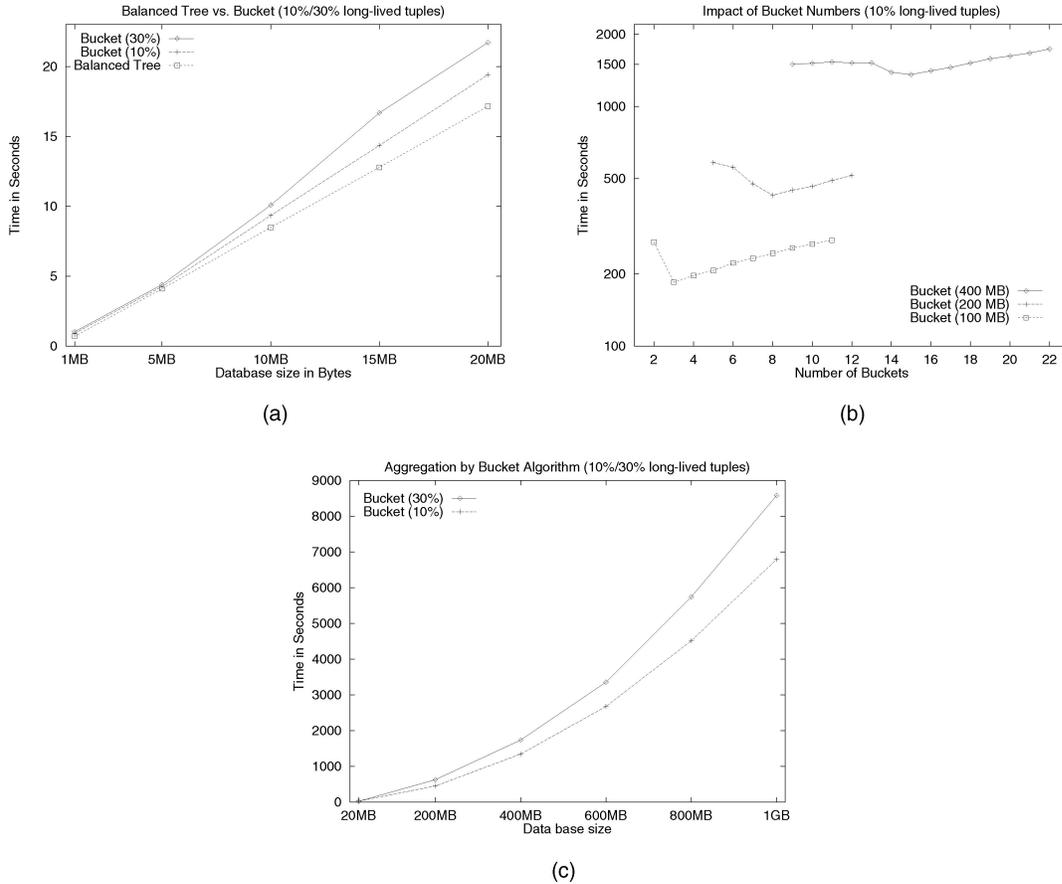


Fig. 13. Aggregation time for large-scale databases. (a) Bucket versus balanced tree, (b) impact of the number of buckets, and (c) scaleup of Bucket algorithm.

existing algorithms such as PA-tree [12] can be used to aggregate each partition.⁴ For the rest of this section, we present experimental results only for count aggregations.

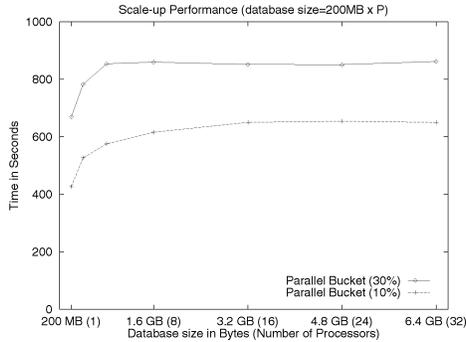
The second set of experiments were carried out to evaluate the *bucket* algorithm proposed in Section 4. First, we performed aggregations with and without data partitioning for small databases, so that we could measure the overhead of data partitioning. The balanced tree algorithm was used to compute count aggregates. In Fig. 13a, 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, however, the bucket algorithm still outperformed the aggregation tree algorithm significantly. (Compare Fig. 12a and Fig. 13a.)

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.

From the experiments, we have noticed that performance of the bucket algorithm is affected by the number of buckets used for data partitioning. More interestingly, there seemed to exist local optimum values, which were determined by database sizes. For example, in Fig. 13b, three, eight, and 15 were the optimal bucket numbers for a database of 100 MBytes, 200 MBytes, and 400 MBytes with 10 percent of long-lived tuples, respectively. Our conjecture is that this is caused by two opposite performance effects from data partitioning. First, since the computational complexity of the balanced tree algorithm is higher than linear ($\mathcal{O}(N \log N)$), the overall computational complexity will be reduced by data partitioning. Specifically, the cost of a balanced tree construction is reduced from $\mathcal{O}(N \log N)$ down to $\mathcal{O}(N \log N - N \log \mathcal{N}_B)$, where N is the number of tuples and \mathcal{N}_B is the number of buckets. Second, the more buckets are used for data partitioning, the more tuples are likely to be replicated, which will in turn increase the cost of disk accesses. We acknowledge that this issue should be addressed more carefully. Refer to Section 7.3 for more discussions.

4. The PA-tree algorithm is expected to perform as efficiently as the proposed balanced tree algorithm for count aggregations. For min and max aggregations, however, the PA-Tree incurs additional overhead to keep track of tuples associated with each node in the tree. Note that no such overhead is necessary for the proposed merge-sort algorithm.



(a)

\mathcal{P}	DB/ \mathcal{P}	Partitioning	Aggregation
	MBytes	Sec. (%)	Sec. (%)
1	200	71.0 (16.7)	355.9 (83.3)
2	200	127.8 (24.3)	399.3 (75.7)
4	200	176.3 (30.6)	399.2 (69.4)
8	200	207.2 (33.6)	408.7 (66.4)
16	200	196.2 (30.2)	454.2 (69.8)
24	200	183.5 (28.0)	470.9 (72.0)
32	200	185.4 (28.5)	464.7 (71.5)

(b)

Fig. 14. Scale-up performance of parallel aggregation. (a) Scale-up performance and (b) separate measurements (10 percent long-lived).

Fig. 13c shows processing times of the bucket algorithm for databases of size from 20 MBytes up to 1 GBytes. 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 MByte database), none of the small-scale aggregation algorithms could be used for this experiment. The results shown in Fig. 13c 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.⁵ This clearly motivates the need of scalable solutions such as the parallel bucket algorithm we proposed in Section 5.

6.4 Parallel Algorithm for Large-Scale Aggregation

The third set of experiments were designed to evaluate the scalability of the parallel bucket algorithm proposed in Section 5. For all the experiments presented in this section, input databases were distributed across participating processors by round-robin partitioning on a nontemporal attribute. By choosing such a nontemporal partitioning scheme for initial data placement, we can effectively eliminate any potential advantage that the parallel bucket algorithm can exploit for better performance. On the other hand, range partitioning on a temporal attribute would be the most favorable data placement for the parallel bucket algorithm, because the number of tuples to be shipped to remote processors could be minimized and thereby reducing communication overhead.

For the scale-up performance measurements, we fixed the size of a database partition on each processor to 10 million tuples (i.e., 200 MBytes), in a way that the entire database would grow proportionally as the number of processors increased. While the number of processors was varied from 1 to 32, the number of local buckets was fixed at eight. Thus, the total number of buckets used for data redistribution was $8 \times \mathcal{P}$, where \mathcal{P} was the number of

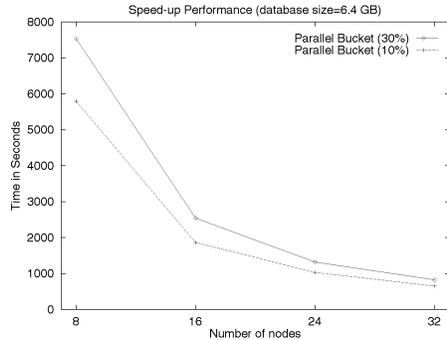
participating processors. The number of local buckets was determined from the previous experiments (see Fig. 13b) based on the local partition size. Note that we used the sequential bucket algorithm for the case $\mathcal{P} = 1$.

In Fig. 14a, the scale-up plots were fairly close to a horizontal line, which indicated a nearly linear scale-up performance with respect to the increasing number of processors. This was corroborated by the fact that the time spent on data partitioning remained quite static when the number of processors was no less than eight. See Fig. 14b for measurements (from the case of 10 percent long-lived tuples) separated into two processing stages. As the number of processors was increased from one to two, data partitioning time was increased by about 80 percent due mainly to additional cost for message passing between processors. In contrast, the time spent on aggregation was increased only by 12 percent due to increased data replication. As the number of processors increased, however, the increase of overhead leveled off and became essentially flat above the four processor case, and thereby allowing nearly linear scale-up performance.

For the speed-up performance measurements, we fixed the size of an entire database to 320 million tuples (i.e., 6.4 GBytes) and determined the size of a database partition based on the number of participating processors. That is, the size of a local partition on a single processor was 6.4 GBytes/ \mathcal{P} . Due to a limited disk space on each processor, we started experiments from eight processors and increased the number of processors up to 32, changing the size of local database partitions accordingly from 800 MBytes to 200 MBytes. The resulting speed-up performance of the parallel bucket algorithm was shown in Fig. 15a.

As a matter of fact, it was surprising that a super-linear speed-up was observed whenever the number of processors increased. From the separate measurements in Fig. 15b (from the case of 10 percent long-lived tuples), such a super-linear speed-up was largely attributed to the performance gain from local aggregation, which grew much faster than linearly as the number of processors increased. Note that the number of buckets used for data redistribution increases proportionally to the number of processors. Thus, we conjecture that the overall aggregation cost was reduced by computing many smaller aggregations rather than computing a few larger aggregations.

5. The faster than linear growth rate is due mainly to data replication. When we increased the database size, the timeline remained identical. This necessarily resulted in the timeline partitioned to shorter intervals covered by more buckets. Consequently, more tuples were expected to be replicated.



(a)

\mathcal{P}	DB/ \mathcal{P}	Partitioning	Aggregation
	MBytes	Sec. (%)	Sec. (%)
8	800	801.0 (13.8)	4995.2 (86.2)
16	400	373.3 (20.0)	1497.3 (80.0)
24	300	248.4 (24.1)	782.2 (75.9)
32	200	185.4 (28.5)	464.7 (71.5)

(b)

Fig. 15. Speed-up performance of parallel aggregation. (a) Speed-up performance and (b) separate measurements (10 percent long-lived).

6.5 Handling Data Skew

To evaluate how the parallel bucket algorithm deals with data skew, we generated synthetic data with Gaussian distribution of temporal attribute values. We assumed that equi-depth histograms were provided as a selectivity estimation mechanism. Fig. 16a illustrates the distribution of start and end times of a local database on each participating processor. We compared the scale-up performance of the parallel bucket algorithm with respect to two different data partitioning policies: 1) *fixed-length partitioning* and 2) *adaptive partitioning*, as described in Section 5.1. In Fig. 16b, the adaptive partitioning significantly improved the performance for skewed data (with 30 percent long-lived tuples). Actually, it took slightly less time than processing a database with uniform distribution using fixed-length partitioning. (Compare with Fig. 14a.) These results clearly demonstrate the effectiveness of the adaptive partitioning and scalable performance of the parallel bucket algorithm, even at the presence of data skew.

7 DISCUSSIONS FOR FURTHER EXTENSIONS

In this section, we discuss application of the proposed temporal aggregation algorithms to queries with GROUP BY clauses. We also discuss further optimization for

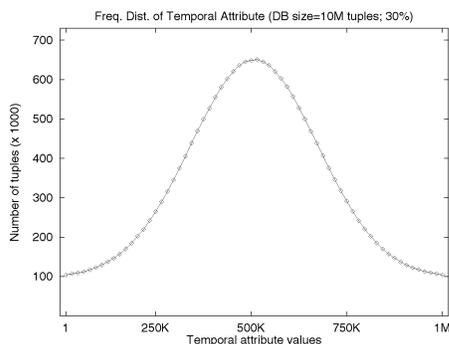
discrete temporal domains and guidelines for selecting the number of buckets for the bucket algorithm.

7.1 Aggregate Queries with GROUP BY Clause

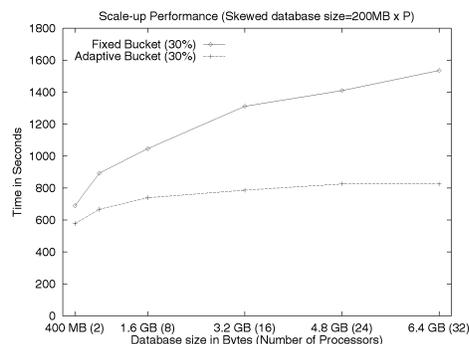
Aggregation queries are often combined with GROUP BY clauses. Thus, it is important to process an aggregation query with a GROUP BY clause efficiently. In traditional relational databases, two approaches based on *sorting* and *hashing* are commonly used. The *sorting* approach orders the tuples in an input relation by GROUP BY attribute. Aggregation is then performed by accessing the tuples in the sorted order. Aggregation for a group completes when a new grouping value is encountered. Then, aggregation for the next value will continue.

The *hashing* approach, on the other hand, requires building a hash table with entries of the form $\langle \text{grouping value}, \text{aggregate value} \rangle$. As a relation is scanned, each tuple is fetched into its corresponding hash entry according to the value of the grouping attribute and the *aggregate value* is updated. A limitation of this approach is that the number of groups should be small enough to fit in memory. If the number of groups outgrows the memory, a relation can be partitioned based on the grouping attribute values and each partition can be aggregated independently [4], [18], [19].

We can apply either of the approaches to temporal aggregation. Once all the tuples in an input relation are sorted or hashed into buckets by grouping attribute values,



(a)



(b)

Fig. 16. Scale-up performance for skewed databases. (a) Data distribution (30 percent long-lived tuples) and (b) fixed-length versus adaptive partitioning.

it is fairly straightforward to aggregate each group by using any of the temporal aggregation algorithms proposed in this paper. The choice of algorithm will be determined by the type of aggregation and the size of each group.

7.2 Fast Aggregation for Discrete Temporal Domain

For temporal attributes defined on discrete domains, certain types of temporal aggregates can be computed in linear time provided that there is enough memory available. For example, suppose that temporal attributes are defined on a time line of one million distinct chronons. Then, a *count* aggregate can be computed in linear time if the memory is large enough to hold two million counters. Two counters are needed for each chronon to keep track of the number of tuples starting and ending at each chronon. These counters are stored in a dense array so that each counter can be directly accessed by the timestamp value of a tuple. Our algorithm proceeds by scanning an input database, updating counters for each tuple. Then, the array of counters is scanned to compute *count* aggregates, in a similar way the balanced tree algorithm does. Note that it is not the database size, but the number of distinct chronons of a temporal domain that determines the memory requirements. We can apply this approach to *sum* and *average* by allocating another word for each chronon to store an aggregate value.

While this approach is very efficient for *count*, *sum*, and *average* aggregates, it is not practical to use this approach for *max* and *min* aggregates. The *max* and *min* aggregates require checking and updating aggregate values associated with all the chronons overlapped with the life span of a tuple. For example, consider a *max* aggregate. Let $\text{maxaggr}[]$ denote an array that stores the max aggregate values. When a tuple t is fetched from a database, its value is compared with all aggregate values stored in the array segment $\text{maxaggr}[t.\text{start}:t.\text{end}]$, where $t.\text{start}$ and $t.\text{end}$ are the start and end timestamps of the tuple t . Specifically, for any i such that $t.\text{start} \leq i \leq t.\text{end}$, if $\text{maxaggr}[i]$ is smaller than the value of t , then, $\text{maxaggr}[i]$ is set to the value of t . This implies that the running time of this approach will be $\mathcal{O}(N^2)$ in the worst case, where N is the number of tuples. A similar approach was proposed to use a hash table for keeping track of the temporal aggregates [14]. Unlike our approach, each tuple should be inserted separately once for each chronon that overlaps the timestamp of the tuple. This may lead to an excessive amount of data replication.

7.3 Optimal Number of Buckets

Knowing that the time complexity of the in-memory aggregation algorithms (i.e., balanced tree and merge-sort) is $\mathcal{O}(N \log N)$, one can argue that the performance of the bucket algorithm can always be improved by increasing the number of buckets. For example, the computational cost of an in-memory aggregation can be reduced from $\mathcal{O}(N \log N)$ down to $\mathcal{O}(N \log N - N \log \mathcal{N}_B)$, if \mathcal{N}_B buckets are used. However, the more buckets are used for data partitioning, the shorter the time interval each bucket represents becomes. This implies that tuples are more likely to be replicated and the overall cost of aggregate computation and disk accesses will increase.

As yet, we have not found any analytic model or mechanism that can be used to obtain an optimal number of buckets. Instead, we have observed two factors that limit the number of buckets. These two factors are essentially characteristics of an input database, and can be used to help find a reasonable choice for the number of buckets if the database characteristics are known a priori.

The first factor is *data distribution*, which may affect the way a time line of an input database is partitioned and assigned to buckets. If tuples are evenly distributed over the time line, then the buckets will represent time intervals of an equal length and will store approximately the same number of tuples. On the other hand, if tuples are unevenly distributed, then the time intervals in the dense regions will be shorter than those in the sparse regions. These shortened time intervals will result in more tuples that are intersected by the boundaries of time intervals and replicated by the bucket algorithm, which will in turn increase the cost of aggregate computation and disk accesses. Thus, it is recommended to err on the side of a smaller number of buckets for an input database with skewed distribution.

The second factor is the *percentage of long-lived tuples*. Any tuple whose time interval is longer than that of a single bucket can be considered a long-lived tuple, because the tuple will necessarily be replicated by the bucket algorithm. The more long-lived tuples exist in the database, the more bucket algorithm costs for aggregate computation and disk accesses. Since large intervals of buckets will reduce the fraction of replicated (or long-lived) tuples, it is recommended to err on the side of a smaller number of buckets for an input database with a large population of long-lived tuples.

8 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 new sequential and parallel bucket algorithms based on novel data partitioning schemes. These algorithms 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. In particular, with the adaptive data partitioning scheme and the local and global meta arrays for partitioned data, we have demonstrated that the parallel bucket algorithm achieves scalable performance for large-scale databases by delivering nearly linear scale-up and speed-up, even at the presence of data skew.

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 number 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 versus nontemporal 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.

ACKNOWLEDGMENTS

This work was sponsored in part by the US National Science Foundation CAREER Award (IIS-9876037) and Research Infrastructure programs EIA-9500991 and EIA-0080123. It was also supported by the Consejo Nacional de Ciencia y Tecnología, scholarship 117476. The authors assume all responsibility for the contents of the paper. The authors would like to thank the anonymous reviewers for all the constructive comments and suggestions. They helped them improve the quality of this paper.

REFERENCES

- [1] M. Barnett, S. Gupta, D. Payne, L. Shuler, R. van de Geijn, and J. Watts, "Interprocessor Collective Communication Library (InterCom)," *Proc. Scalable High Performance Computing Conf.*, pp. 357-364, May 1994.
- [2] Jon Louis Bentley, "Algorithms for Klee's Rectangle Problems," technical report unpublished, Computer Science Dept., Carnegie Mellon Univ., 1977.
- [3] Ohio Supercomputer Center, LAM/MPI Parallel Computing, <http://www.osc.edu/lam.html>, 1998.
- [4] S. Chaudhuri and K. Shim, "Including Group by in Query Optimization," *Proc. 20th Very Large Database Conf.*, pp. 354-366, Sept. 1994.
- [5] S. Chaudhuri and U. Dayal, "An Overview of Data Warehousing and OLAP Technology," *SIGMOD Record*, vol. 26, no. 1, Mar. 1997.
- [6] T.H. Cormen, C.E. Leiserson, and R.L. Rivest, *Introduction to Algorithms*. McGraw Hill, 1990.
- [7] D.J. DeWitt, R.H. Katz, F. Olken, L.D. Shapiro, M.R. Stonebraker, D. Wood, "Implementation Techniques for Main Memory Database Systems," *Proc. 1984 ACM-SIGMOD Conf.*, pp. 1-8, June 1984.
- [8] R. Epstein, "Techniques for Processing of Aggregates in Relational Database Systems," Technical Report UCB/ERL M7918, Univ. of California, Berkeley, Feb. 1979.
- [9] J.C. Freytag and N. Goodman, "Translating Aggregate Queries Into Iterative Programs," *Proc. 12th Very Large Database Conf.*, pp. 138-146, Aug. 1986.
- [10] J. Alvin, G. Gendrano, B.C. Huang, J.M. Rodrigue, B. Moon, and R.T. Snodgrass, "Parallel Algorithms for Computing Temporal Aggregates," *Proc. 15th Int'l Conf. Data Eng.*, Mar. 1999.
- [11] C.S. Jensen and R.T. Snodgrass, "Semantics of Time-Varying Information," *Information Systems*, vol. 21, no. 4, pp. 311-352, 1996.
- [12] J.S. Kim, S.T. Kang, and M.H. Kim, "On Temporal Aggregate Processing Based on Time Points," *Information Processing Letters*, vol. 71, no. 5-6, Sept. 1999.
- [13] N. Kline and R.T. Snodgrass, "Computing Temporal Aggregates," *Proc. 11th Int'l Conf. Data Eng.*, pp. 222-231, Mar. 1995.
- [14] R.N. Kline, *Aggregation in Temporal Databases*, PhD thesis, Univ. Arizona, Tucson, May 1999.
- [15] D.E. Knuth, "Sorting and Searching," *The Art of Computer Programming*, vol. 3, Mass.: Addison-Wesley, 1973.
- [16] U. Manber, *Introduction to Algorithms: A Creative Approach*. Addison-Wesley, 1989.
- [17] G. Piatesky-Shapiro and C. Connel, "Accurate Estimation of the Number of Tuples Satisfying a Condition," *Proc. 1984 ACM-SIGMOD Conf.*, pp. 256-276, June 1984.
- [18] R. Ramakrishnan, *Database Management Systems*. McGraw-Hill, 1998.
- [19] A. Silberschatz, H.F. Kort, and S. Sudarshan, *Database System Concepts*. McGraw-Hill, third ed., 1999.
- [20] R.T. Snodgrass, S. Gomez, and E. Mackenzie, "Aggregates in the Temporal Query Language TQuel," *IEEE Trans. Knowledge and Data Eng.*, vol. 5, no. 5, pp. 826-842, Oct. 1993.
- [21] M. Stonebraker, "The Case for Shared Nothing," *A Quarterly Bull. of the IEEE Computer Soc. Technical Committee on Database Eng.*, vol. 9, no. 1, pp. 4-9, Mar. 1986.
- [22] A. Tansel et al. *Temporal Databases: Theory, Design, and Implementation*. Database Systems and Applications Series, Benjamin/Cummings, 1993.
- [23] P.A. Tuma, "Implementing Historical Aggregates in TempIS," Master's thesis, Wayne State Univ., Nov. 1992.
- [24] J. Yang and J. Widom, "Incremental Computation and Maintenance of Temporal Aggregates," *Proc. 17th Int'l Conf. Data Eng.*, Apr. 2001.
- [25] X. Ye and J.A. Keane, "Processing Temporal Aggregates in Parallel," *Proc. IEEE Int'l Conf. Systems, Man, and Cybernetics*, pp. 1373-1378, Oct. 1997.
- [26] D. Zhang, A. Markowetz, V. Tsotras, D. Gunopulos, and B. Seeger, "Efficient Computation of Temporal Aggregates with Range Predicates," *Proc. 20th ACM SIGACT-SIGMOD-SIGART Symp. Principles of Database Systems*, May 2001.



Bongki Moon received the PhD degree in computer science from the University of Maryland, College Park, in 1996, and the MS and BS degrees in computer engineering from Seoul National University, Korea, in 1985 and 1983, respectively. He is an assistant professor in the Department of Computer Science, University of Arizona. His current research interests include XML indexing and query processing, high-performance spatial and temporal databases, scalable Web servers, data mining and warehousing, and parallel and distributed processing. He was a member of the research staff at Communication Systems Division, Samsung Electronics Corp., Korea, from 1985 to 1990. He received the US National Science Foundation CAREER Award.



Ines Fernando Vega Lopez received the MS degree in computer science from the University of Arizona in 1999, and the BE degree in computer systems from the Instituto Tecnológico y de Estudios Superiores de Monterrey, Sinaloa, Mexico, in 1994. He is a PhD student at the Department of Computer Science, University of Arizona. His current research interests include multimedia and temporal databases, as well as data mining and warehousing.



Vijaykumar Immanuel received the MS degree in computer science from the University of Arizona, Tucson, in 1998, and the BTech degree in computer science from the Indian Institute of Technology, Madras, India in 1996. He is a Software Engineer at Compaq Computer Corp. His current areas of interests include storage systems, networks, and fault-tolerance.

► For more information on this or any computing topic, please visit our Digital Library at <http://computer.org/publications/dlib>.