

PARSIMONY: An Infrastructure for Parallel Multidimensional Analysis and Data Mining

Sanjay Goil¹

Sun Microsystems Inc., 901 San Antonio Road, Palo Alto, California 94303

E-mail: sgoil@eng.sun.com

and

Alok Choudhary²

Northwestern University, 2145 Sheridan Road, Evanston, Illinois 60208

E-mail: choudhar@ece.nwu.edu

Received March 1, 1999; accepted December 31, 1999

Multidimensional analysis and online analytical processing (OLAP) operations require summary information on multidimensional data sets. Most common are aggregate operations along one or more dimensions of numerical data values. Simultaneous calculation of multidimensional aggregates are provided by the *Data Cube* operator, used to calculate and store summary information on a number of dimensions. This is computed only partially if the number of dimensions is large. Query processing for these applications requires different views of data to gain insight and for effective decision support. Queries may either be answered from a materialized cube in the data cube or calculated on the fly.

The multidimensionality of the underlying problem can be represented both in relational and in multidimensional databases, the latter being a better fit when query performance is the criteria for judgment. Relational databases are scalable in size for OLAP and multidimensional analysis and efforts are on to make their performance acceptable. On the other hand multidimensional databases have proven to provide good performance for such queries, although they are not very scalable. In this article we address (1) scalability in multidimensional systems for OLAP and multidimensional analysis and (2) integration of data mining with the OLAP framework. We describe our

¹ This work was performed while this author was a graduate student in the ECE Department at Northwestern University.

² This work was supported in part by NSF Young Investigator Award CCR-9357840 and NSF CCR-9509143.

system PARSIMONY, parallel and scalable infrastructure for multidimensional online analytical processing, used for both OLAP and data mining. Sparsity of data sets is handled by using *chunks* to store data either as a dense block using multidimensional arrays or as sparse representation using a bit encoded sparse structure. Chunks provide a multidimensional index structure for efficient dimension oriented data accesses much the same as multidimensional arrays do. Operations within chunks and between chunks are a combination of relational and multidimensional operations depending on whether the chunk is sparse or dense. Further, we develop parallel algorithms for data mining on the multidimensional cube structure for attribute-oriented association rules and decision-tree-based classification. These take advantage of the data organization provided by the multidimensional data model.

Performance results for high dimensional data sets on a distributed memory parallel machine (IBM SP-2) show good speedup and scalability. © 2001

Academic Press

1. INTRODUCTION

Online analytical processing (OLAP) and multidimensional analysis are used for decision support systems and statistical inferencing to find interesting information from large databases. Multidimensional databases are suitable for OLAP and data mining since these applications require dimension oriented operations on data. Traditional multidimensional databases store data in multidimensional arrays on which analytical operations are performed. Multidimensional arrays are good to store dense data, but most data sets are sparse in practice for which other efficient storage schemes are required.

It is important to weigh the trade-offs involved in reducing the storage space versus the increase in access time for each sparse data structure, in comparison to multidimensional arrays. These trade-offs are dependent on many parameter, some of which are (1) number of dimensions, (2) sizes of dimensions, and (3) degree of sparsity of the data. Complex operations such as those required for OLAP can be very expensive in terms of data access time if efficient data structures are not used.

Sparse data structures such as the R-tree and its variants have been used for OLAP [27]. Range queries with a lot of unspecified dimensions are expensive because many paths have to be traversed in the tree to calculate aggregates. Chunking has been used in [32] with dense and sparse chunks. Sparse chunks store an offset-value pair for the data present. Dimensional operations on these require materializing the sparse chunk into a multidimensional array and performing array operations on it. For a high number of dimensions this might not be possible since the materialized chunk may not fit in memory. A split storage is described in [6], where the dimensions are split into sparse and dense subsets. The sparse dimensions use a sparse index structure to index into the dense blocks of data stored as multidimensional arrays. This treats the dimensions nonuniformly, requiring the user to partition the dimensions. Further, none of these address parallelism and scalability to large data sets in a high number of dimensions.

We compare the storage and operational efficiency in OLAP and multidimensional analysis of various sparse data storage schemes in [10]. A novel data structure

using bit encodings for dimension indices called bit-encoded sparse structure (BESS) is used to store sparse data in chunks, which supports fast OLAP query operations on sparse data using bit operations without the need for exploding the sparse data into a multidimensional array. This allows for high dimensionality and large dimension sizes.

In this article we present a parallel and scalable OLAP and data mining framework for large data sets. Parallel data cube construction for large data sets and a large number of dimensions using both dense and sparse storage structures is presented. Sparsity is handled by using compressed *chunks* using a bit encoded sparse structure. Data are read from a relational data warehouse which provides a set of tuples in the desired number of dimensions. Reference [13] presents data partitioning and *sort-based* loading of chunks in the *base* cube (which is a n -dimensional structure at level n of the data cube) from which the data cube is computed. OLAP queries can now be answered from the precomputations available in the data cube. Precomputed values are used in the probability calculations for association rule mining. Also, classification trees can be built by using the aggregates of class-id counts to calculate the splitting criterion for each dimension.

Figure 1 motivates the overall design of our parallel multidimensional database (MDDB) engine for OLAP and data mining. At the lowest level is a parallel machine which can be either a distributed memory message passing machine or a shared memory machine. This machine has a layer of system software that allows programming and provides communication mechanisms. We build the parallel MDDB engine using this layer since the parallel functionality is optimized by the machine vendor. Data distribution, sparse representation, and load balancing are issues that are addressed in the parallel design of the framework. A query processing layer for point queries, range queries, and hierarchical consolidation queries provide support for OLAP and data mining activities. This layered approach is modular and optimizations are done at each level to tune performance on the underlying architecture. Further details of this infrastructure can be obtained from [9].

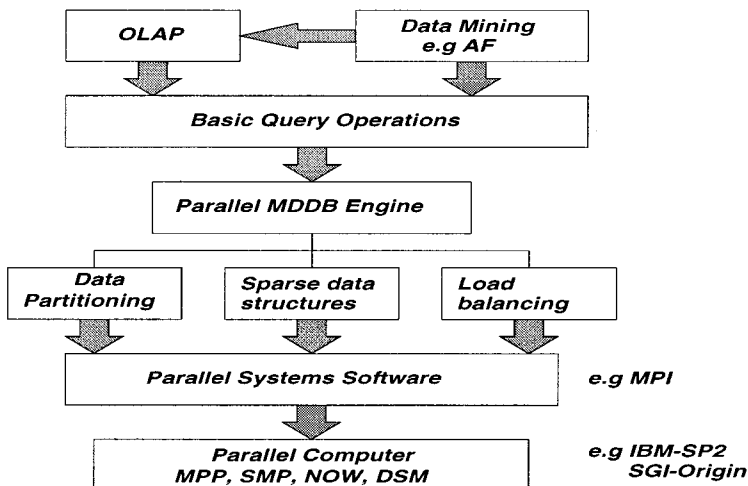


FIG. 1. Overall schematic.

We do not address the question of selective cube materialization here. Our focus in this work is on the scalable maintenance of the multidimensional data cube structure in parallel. Clearly, the selection of subcubes for materialization based on their benefit to query response times and cost-saving ratio is an important problem and has received much attention recently [7, 19, 28]. We believe these can be easily incorporated into our framework. Our techniques for multidimensional analysis can also be applied to scientific and statistical databases [22, 30].

The rest of the article is organized as follows. Section 2 describes OLAP using the data cube operator and queries on them. Section 3 presents multidimensional storage using chunks and BESS for sparse data. Section 4 presents steps in the computation of the data cube on a parallel machine and the overall design. Section 5 describes the algorithms, techniques, and optimizations in the parallel building of the simultaneous multidimensional aggregates and the factors affecting performance. Section 6 describes data mining algorithms such as association rule mining and classification on the multidimensional structure. Section 7 presents performance results for cube building and analysis for communication and I/O for it. Section 8 concludes the article.

2. OLAP AND DATA CUBES

OLAP is used to summarize, consolidate, view, apply formulae to, and synthesize data according to multiple dimensions. Traditionally, a relational approach, using relational databases (relational OLAP), has been taken to build and query such systems. A complex analytical query is cumbersome to express in SQL and it might not be efficient to execute. Alternatively, multidimensional database techniques (multidimensional OLAP) with multidimensional APIs have been applied to decision-support applications. Data are stored in multidimensional structures which is a more natural way to express the multidimensionality of the enterprise data and is more suited for analysis. A *cell* in multidimensional space represents a tuple, with the attributes of the tuple identifying the location of the tuple in the multidimensional space and the *measure* values represent the content of the cell.

Data can be organized into a data cube by-calculating all possible combinations of *Group-bys* [16]. This operation is useful for answering OLAP queries which use aggregation on different combinations of attributes (Fig. 2). For a data set with n attributes this leads to 2^n GROUP-BY calculations. A data cube treats each of the k , $0 \leq k < n$ aggregation attributes as a dimension in k -space. Figure 3 shows an example of a data cube constructed from a data set with three attributes. In relational systems, a *fact* table is constructed, using a *star schema* from the original

```

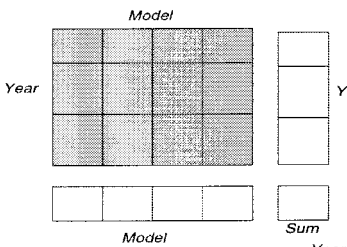
SELECT Model, Year, Color, SUM(sales) AS Sales
FROM Sales
WHERE Model IN {'Ford', 'Chevy'} AND
        Year BETWEEN 1990 AND 1992
GROUP BY CUBE (Model, Year, Color);

```

FIG. 2. Query for the data cube example.

Model	Year	Color	Sales
Chevy	1990	Blue	87
Chevy	1990	Red	5
Chevy	1990	ALL	92
Chevy	ALL	Blue	87
Chevy	ALL	Red	5
Chevy	ALL	ALL	92
Ford	1990	Blue	99
Ford	1990	Green	64
Ford	1990	ALL	163
Ford	1991	Blue	7
Ford	1991	Red	8
Ford	1991	ALL	15
Ford	ALL	Blue	106
Ford	ALL	Green	64
Ford	ALL	Red	8
ALL	1990	Blue	186
ALL	1990	Green	64
ALL	1991	Blue	7
ALL	1991	Red	8
Ford	ALL	ALL	178
ALL	1990	ALL	255
ALL	1991	ALL	15
ALL	ALL	Blue	193
ALL	ALL	Green	64
ALL	ALL	Red	13
ALL	ALL	ALL	270

Cross Tab



Data Cube

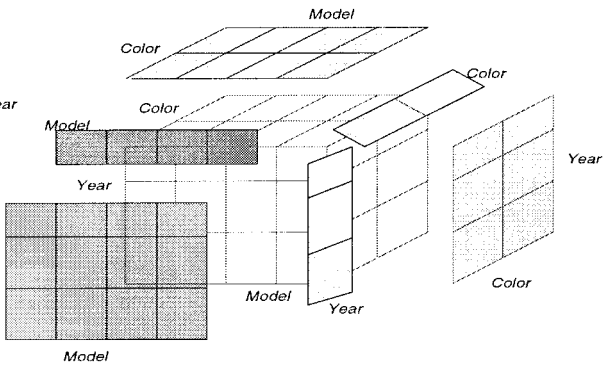


FIG. 3. Data cube example.

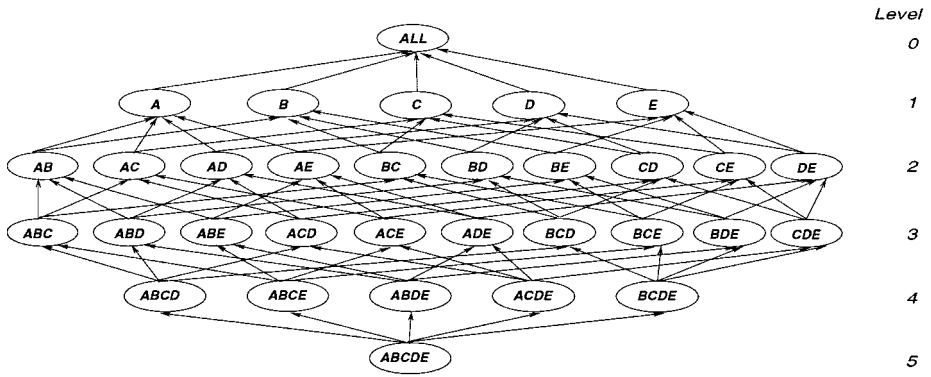


FIG. 4. Lattice for cube operator.

table. The fact table has keys of the *dimensional* tables and contains the various aggregates, the aggregated attribute represented by a special keyword *ALL*. Data access for queries is through appropriate relational operations such as **join** and **projection**. In SQL, the *Cube* operator represents the data cube calculation as illustrated in Fig. 2. A multidimensional representation uses a multicube architecture to store data organized with dimensional information in a suitable multidimensional structure. This is a more natural representation of dimensional data and leads to efficient data access operations.

Data cube operators generalize the histogram, cross-tabulation, roll-up, drill-down and subtotal constructs. A lattice structure is used to represent the various aggregate calculations. Figure 4 shows a lattice structure for the data cube with five dimensions (A, B, C, D, and E). At a level i , $0 \leq i \leq n$ of the lattice, there are $C(n, i)$ subcubes (aggregates) with exactly i dimensions, where the function C gives all the combinations having i distinct dimensions from n dimensions. A total of $\sum_{i=0}^n C(n, i) = 2^n$ subcubes are present in the data cube including the base cube. Optimizations of calculating the aggregates in the subcubes can be performed using the lattice structure, augmented by the various computation and communication costs to generate a DAG of cube orderings. The arrows are labeled by the costs of constructing an aggregate and an order of construction is chosen which minimizes the cost [1]. This is described in detail later in the section which presents details on the cost model (Section 4.3).

2.1. OLAP Queries

OLAP queries can be answered by the aggregates in the data cube. The following operations are defined on the cube:

- *Pivoting*: This is also called *rotating* and involves rotating the cube to change the dimensional orientation of a report or page on display. It may consist of swapping the two dimensions (row and column in a 2D-cube) or introducing another dimension instead of some dimension already present in the cube.

- *Slicing-dicing*: This operation involves selecting some subset of the cube. For a fixed attribute value in a given dimension, it reports all the values for all the other dimensions. It can be visualized as a *slice* of the data in a 3D-cube.

- *Roll-up*: Some dimensions have a hierarchy defined on them. Aggregations can be done at different levels of hierarchy. Going up the hierarchy to higher levels of generalization is known as roll-up. For example, aggregating the dimension up the hierarchy (*day* → *month* → *quarter*) is a roll-up operation.

- *Drill-down*: This operation traverses the hierarchy from lower to higher levels of detail. Drill-down displays detail information for each aggregated point.

The operations needed for OLAP and data mining require access to data along a particular dimension or a combination of dimensions. Usual operations are of the type of aggregation along a particular attribute, which is represented by a dimension in the multidimensional database. The set of operations required for OLAP and data mining should be efficiently supported by the sparse data structure for good performance. The operations fall under the following categories.

- Retrieval of a random cell element
- Retrieval along a dimension or a combination of dimensions
- Retrieval for values of dimensions within a range (range queries)
- Aggregation operations on dimensions
- Multidimensional aggregation (generalization–consolidation: lower to higher level in hierarchy)

Most data dimensions have hierarchies defined on them. Typical OLAP queries probe summaries of data at different levels of the hierarchy. Consolidation is a widely used method to provide roll-up and drill-down functions in OLAP systems. Each dimension in a cube can potentially have a hierarchy defined on it. This hierarchy can be applied to all the subcubes of the data cube. We do not discuss query processing any further and instead concentrate on the design of the storage and computation infrastructure for parallel data cubes.

3. DATA STORAGE: CHUNKS AND BESS

Multidimensional database technology facilitates flexible, high performance access and analysis of large volumes of complex and interrelated data [21]. It is more natural and intuitive for humans to model a multidimensional structure. A *chunk* is defined as a block of data from the multidimensional array which contains data in all dimensions. A collection of chunks defines the entire array. Figure 5a shows chunking of a three dimensional array. A chunk is stored contiguously in memory and data in each dimension are strided with the dimension sizes of the chunk. Most sparse data may not be uniformly sparse. Dense clusters of data can be stored as multidimensional arrays. Sparse data structures are needed to store the sparse portions of data. These chunks can then be either stored as dense arrays or stored using an appropriate sparse data structure as illustrated in Fig. 5b. Chunks also act as an index structure which helps in extracting data for queries and OLAP operations.

Typically, sparse structures have been used for the advantages they provide in terms of storage, but operations on data are performed on a multidimensional array which is populated from the sparse data. However, this is not always possible when

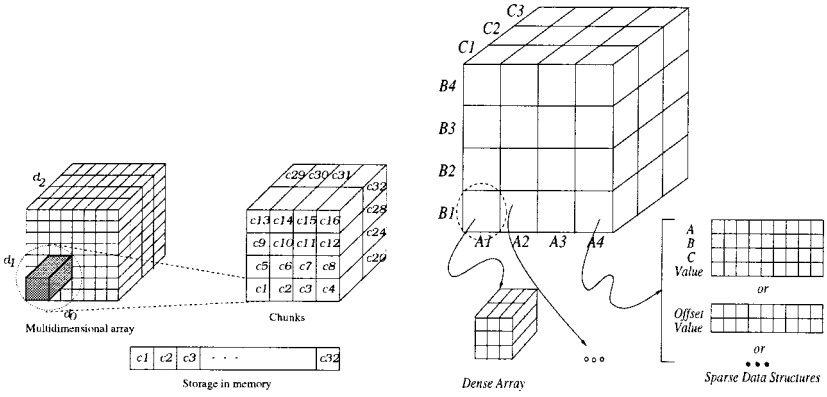


FIG. 5. Storage of data in chunks: (a) chunking for a 3D array, and (b) chunked storage for the cube.

either the dimension sizes are large or the number of dimensions is large. Since we are dealing with multidimensional structures for a large number of dimensions, we are interested in performing operations on the sparse structure itself. This is desirable to reduce I/O costs by having more data in memory to work on. This is one of the primary motivations for our BESS. For each cell present in a chunk a dimension index is encoded in $\lceil \log |d_i| \rceil$ bits for each dimension d_i of size $|d_i|$. An 8-byte encoding is used to store the BESS index along with the value at that location. A larger encoding can be used if the number of dimensions is larger than 20. A dimension index can then be extracted by a bit mask operation. Aggregation along a dimension d_i can be done by masking its dimension encoding in BESS and using a sort operation to get the duplicate resultant BESS values together. This is followed by a scan of the BESS index, aggregating values for each duplicate BESS index. For dimensional analysis, aggregation needs to be done for appropriate chunks along a dimensional plane. This will be elaborated further in a later section.

Figure 6a illustrates an example for storing a two dimensional sparse chunk of size 3×2 . Values $v1$ and $v2$ in chunk 0 are stored using a bit-encoded structure. Dimension x can be encoded in two bits since there are only three distinct

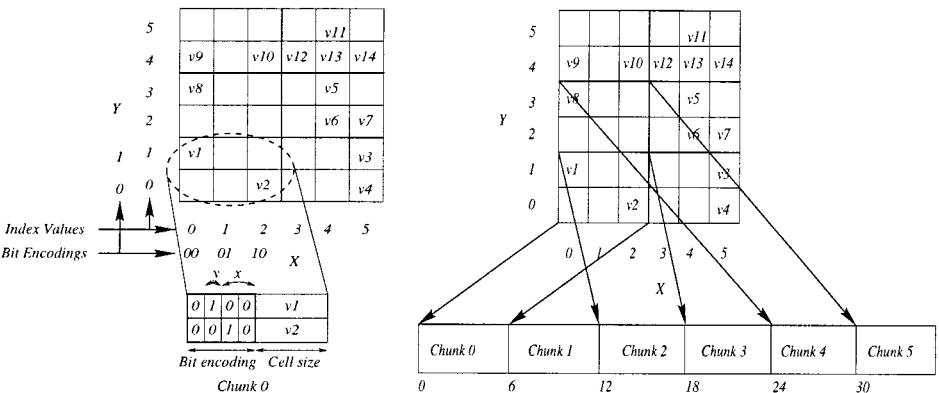


FIG. 6. BESS used in chunk storage: (a) Chunk storage using BESS, and (b) storing chunk offsets for index retrieval.

values. Dimension y needs one bit to differentiate between 0 and 1. Hence three bits are needed for an index. Higher dimensional structures will use more bits but the minimum storage is 32 bits since data are word aligned. Thus, a large number of dimensions can be stored if two integers are used.

A chunk index structure stores the chunk offsets for each chunk. Chunks are stored in memory in some order of dimensions. Without loss of generality let us assume that the order of storage is in the order X followed by Y in Fig. 6b. To reference an element in the chunk, first the chunk offset is dereferenced to get the dimension index values. Since the dimension extents are known for the chunk this can be done easily. The second step is to add the index value in the chunk to this by retrieving it from the bit encoding. This can be done by shifting bits in the bit-encoded structure and using a bit-mask to extract the bits that are used to encode each dimension.

4. OVERALL DESIGN

In this section we describe our design for a parallel and scalable data cube on coarse grained parallel machines (e.g, 4BM SP-2) or a network of workstations characterized by powerful general purpose processors (few to a few hundred) with a fast interconnection. The programming paradigm used is a high level programming language (e.g., C/C++) embedded with calls to a portable communication library (e.g., message passing interface).

In what follows, we address issues of data partitioning, parallelism, schedule construction, data cube building, chunk storage, and memory usage on this machine architecture. Figure 7 summarizes the various options available for these, especially in terms of storage of cubes, parallelism at different levels, aggregation calculation orderings, and the chunk access for aggregations. Moreover, a partial cube can be constructed if the number of dimensions is large or a specific level of the cube is needed. For example, in 2-way attribute-oriented data mining of associations, all cubes at level 2 are materialized by using the base cube and the minimum materializations of subcubes at the intermediate levels between 3 and $n - 1$ [12].

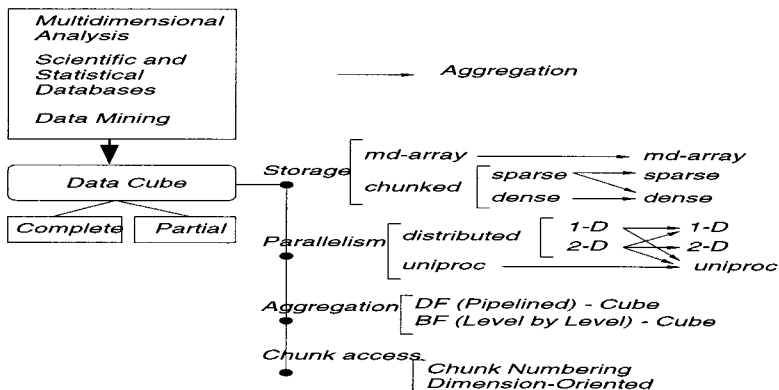


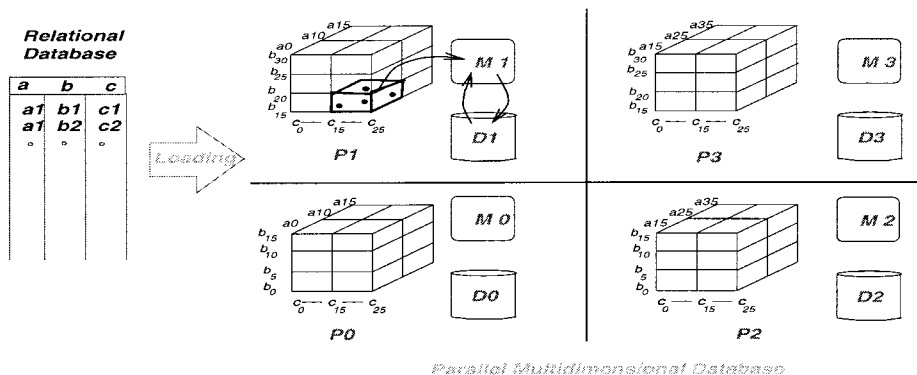
FIG. 7. Scheduling for data cube.

4.1. Data Partitioning and Parallelism

Data are partitioned on processors to distribute work equally. In addition, a partitioning scheme for multidimensional structures has to be *dimension-aware* and for dimension-oriented operations it must have some regularity in the distribution. A dimension, or a combination of dimensions, can be distributed. In order to achieve sufficient parallelism, it would be required that the product of cardinalities of the distributed dimensions be much larger than the number of processors. For example, for five-dimensional data (*ABCDE*), a 1D distribution will partition *A* and a 2D distribution will partition *AB*. We assume that dimensions are available that have cardinalities much greater than the number of processors in both cases. That is, either $|A_i| \gg p$ for some i or $|A_i| |A_j| \gg p$ for some $i, j, 0 \leq i, j \leq n - 1, n$ is the number of dimensions. Partitioning determines the communication requirements for data movement in the intermediate aggregate calculations in the data cube. We support both 1D and 2D partitions in our implementations. Let K (1 or 2) be the number of dimensions used for the partitioning. For $K=2$, a two-dimensional partitioning, p processors are divided into two groups ($k, \frac{p}{k}$), where k is the number of processors in the group to be created by the first distributed dimension, chosen to be a divisor of p . For each tuple, A_i is used to choose among the k partitions and then A_j is further used to find the correct processor among the $\frac{p}{k}$ processors in that group. k can be appropriately chosen to reflect the ratios of the dimension sizes used in the partitioning. Figure 8 shows the base cube loading for a three-dimensional cube (*ABC*) distributed on four processors in a 2D distribution, where dimensions *A* and *B* are distributed in a (2, 2) grid. Clearly, when $K=1$, a 1D partitioning is done with $k = p$.

Since 2^n cubes are being constructed, we keep them distributed as well. The distribution of these cubes depends on the cardinalities of their largest one or two dimensions. The same criteria are used here as for the base cube. However, redistribution of dimensions and chunks may be required if a dimension is partitioned anew or is repartitioned.

Table 1 shows the various distributions for aggregate calculations supported in our framework. The underlined dimensions are partitioned. Calculations are either *local* or *nonlocal*. Local calculations maintain the data distribution on each processor



Parallel Multidimensional Database

FIG. 8. Base cube loading for three-dimensional cube on four processors.

TABLE 1
Partitioning of Subcubes following Aggregation Calculations

Distribution	Local	Nonlocal	
		Dimension 1	Dimension 2
$2D \rightarrow 2D$	$\underline{ABC} \rightarrow \underline{AB}$	$\underline{ABC} \rightarrow \underline{BC}$	$\underline{ABC} \rightarrow \underline{AC}$
$2D \rightarrow 1D$		$\underline{ABC} \rightarrow \underline{BC}$	$\underline{ABC} \rightarrow \underline{AC}$
$1D \rightarrow 1D$	$\underline{ABC} \rightarrow \underline{AB}, \underline{AC}$	$\underline{ABC} \rightarrow \underline{BC}$	
$2D \rightarrow 0D$		$\underline{ABC} \rightarrow \underline{BC}$	$\underline{ABC} \rightarrow \underline{AC}$
$1D \rightarrow 0D$		$\underline{ABC} \rightarrow \underline{BC}$	$\underline{ABC} \rightarrow \underline{AC}$
$0D \rightarrow 0D$	$ABC \rightarrow AB, AD$		

and the aggregation calculation does not involve any interprocessor communication. Nonlocal calculations distribute an undistributed dimension such as in $\underline{ABC} \rightarrow \underline{AC}$, where dimension B is aggregated and C , which was previously undistributed, is distributed. Another calculation is $\underline{ABC} \rightarrow \underline{BC}$, where A is aggregated, B , the second distributed dimension, becomes the first distributed dimension, and C gets distributed as the second dimension. These can be categorized as either dimension 1 or dimension 2 being involved in the (re) distribution. Also, as illustrated in Fig. 7, the subcubes can be stored as *chunked* or as *multidimensional arrays* which are distributed or on a single processor with these distributions. The multidimensional arrays are, however, restricted to a 1D distribution since their sizes are small and 2D distribution will not provide sufficient parallelism. The data cube build scheduler does not evaluate the various possible distributions currently, instead calculating the costs based on the estimated sizes of the source and the target subcubes and uses a partitioning based on the dimension cardinalities.

4.2. Schedule Generation for Data Cube

Several optimizations can be done over the naive method of calculating each aggregate separately from the initial data [16]. *Smallest parent* computes a group-by by selecting the smallest of the previously computed group-bys from which it is possible to compute the group-by. Consider a four attribute cube ($ABCD$). Group-by AB can be calculated from $ABCD$, ABD , and ABC . Clearly the sizes of ABC and ABD are smaller than that of $ABCD$ and are better candidates. The next optimization is to compute the group-bys in an order in which the next group-by calculation can benefit from the cached results of the previous calculation. This can be extended to disk-based data cubes by reducing disk I/O and caching in main memory. For example, after computing ABC from $ABCD$ we compute AB followed by A . An important multiprocessor optimization is to *minimize interprocessor communication*. The order of computation should minimize the communication among the processors because interprocessor communication costs are typically higher than computation costs. For example, for a 1D partition, $BC \rightarrow C$ will have a higher communication cost to first aggregate along B and then divide C among the processors in comparison to $CD \rightarrow C$ where a local aggregation on each processor along D will be sufficient.

A lattice framework to represent the hierarchy of the group-bys was introduced in [19]. This is an elegant model for representing the dependencies in the calculations and also to model costs of the aggregate calculations. A scheduling algorithm can be applied to this framework substituting the appropriate costs of computation and communication. A lattice for the group-by calculations for a five-dimensional cube ($ABCDE$) is shown in Fig. 4. Each node represents an aggregate and an arrow represents a possible aggregate calculation which is also used to represent the cost of the calculation.

Calculation of the order in which the Group-bys are created depends on the cost of deriving a lower order (one with a lower number of attributes) group-by from a higher order (also called the *parent*) group-by. For example, between $ABD \rightarrow BD$ and $BCD \rightarrow BD$ one needs to select the one with the lower cost. Cost estimation of the aggregation operations can be done by establishing a cost model. Some calculations do not involve communication and are *local*; others involving communication are labeled as *nonlocal*. Details of these techniques for a parallel implementation using multidimensional arrays can be found in [11]. However, with chunking and the presence of sparse chunks the cube size cannot be taken for calculating computation and communication costs. Size estimation is required for sparse cubes to estimate computation and communication costs when dimension aggregation operations are performed. We use a simple analytical algorithm for size estimation in the presence of hierarchies presented in [31]. This is shown to perform well for uniformly distributed random data and also works well for some amount of skew. Since we need reasonable estimates to select the materialization of a subcube from a subcube at a higher level, this works well.

4.3. Cost Model

Size estimation is important to label the edges of the cube lattice with costs of various aggregation calculations. Apart from sizes, computation and communication costs have to be modeled for the various cubes and their distributions. Cubes are partitioned on processors and communication is generated for the various combinations of distributions shown in Table 1. Let $\mathcal{V}(\prod_{i=0}^{n-1} D_i)$ be the estimated value of the number of tuples in the region defined by the dimensions $|D_i|$. Let P be the number of processors, divided into a $P_x \times P_y$ grid of processors for a two-dimensional partitioning. The communication cost for a 2D to 2D calculation where the first dimension, D_0 , is aggregated, dimension D_1 is redistributed from P_y to P_x , and dimensions D_2 is distributed on P_y is given by

$$\begin{aligned} & \mathcal{V} \left(\frac{|D_0|}{P_x} \frac{|D_1|}{P_y} \prod_{i=2}^{n-1} D_i \right) T_{calc} + \frac{\mathcal{V} \left(\frac{|D_0|}{P_x} \frac{|C_1|}{P_y} \prod_{i=2}^{n-1} D_i \right)}{S_m} T_{I/O} \\ & + \left(T_s + T_w \mathcal{V} \left(\frac{|D_1| P_x}{P_y} \frac{|D_2|}{P_y} \prod_{i=3}^{n-1} D_i \right) P \right). \end{aligned}$$

Dimension D_0 , distributed on P_x processors, is aggregated locally, dimension D_1 , distributed on P_y , is redistributed to P_x processors, and dimension D_2 , which is local, now distributed on P_y processors. The computation cost is the cost of aggregating the tuples in all the chunks on the processor, which are stored as a collection of minichunks for sparse chunks. S_m is the size of a minichunk, and S_b is the size of the aggregate buffer. The I/O cost is in terms of the number of minichunks read and the number of minichunks written after the aggregation. The number of minichunk writes for dimension order traversal is estimated as a ratio of the total number of tuples in the target cube and the minichunk size S_m . We estimate all chunks as sparse chunks for schedule calculations. However, chunks above the user defined threshold density are converted to dense chunks in the implementation. Table 2 summarizes the cost estimations for the various distributions for the aggregation of the first distributed dimension is aggregated. Table 3 summarizes the cost estimation when dimension 2 is aggregated under the various distributions.

Reference [32] uses a minimum memory spanning tree (MMST) for a cube in an optimal dimension order to reduce the memory requirements to overlap the computation of the MMST subtrees. The optimal order is shown to be the increasing order of cardinalities of the dimensions. This is due to the fact that chunks are always accessed in storage order and main memory is allocated for all the target chunks that map in the region defined by the chunks that are being read. Since the

TABLE 2

Cost Model for Parallel Aggregation Calculations for Aggregating Dimension 1

Distribution	T_{calc}	Cost	
		T_{comm}	$T_{I/O}$
$2D \rightarrow 2D$	$\gamma \left(\frac{ D_0 }{P_x} \frac{ D_1 }{P_y} \prod_{i=2}^{n-1} D_i \right)$	$P_x(T_s + T_w) \frac{\gamma \left(\frac{ D_1 P_x}{P_y} \frac{ D_2 }{P_y} \prod_{i=3}^{n-1} D_i \right)}{S_{sendbuf}}$	$\frac{\gamma \left(\frac{ D_0 }{P_x} \frac{ D_1 }{P_y} \prod_{i=2}^{n-1} D_i \right)}{S_{minichunk}}$
$2D \rightarrow 1D$	$\gamma \left(\frac{ D_0 }{P_x} \frac{ D_1 }{P_y} \prod_{i=2}^{n-1} D_i \right)$	$P(T_s + T_w) \frac{\gamma \left(\frac{ D_1 P}{P_y} \prod_{i=2}^{n-1} D_i \right)}{S_{sendbuf}}$	$\frac{\gamma \left(\frac{ D_0 }{P_x} \frac{ D_1 }{P_y} \prod_{i=2}^{n-1} D_i \right)}{S_{minichunk}}$
$2D \rightarrow 0D$	$\gamma \left(\frac{ D_0 }{P_x} \frac{ D_1 }{P_y} \prod_{i=2}^{n-1} D_i \right)$	$(T_s + T_w) \frac{\gamma \left(\frac{ D_1 }{P_y} \prod_{i=2}^{n-1} D_i \right)}{S_{sendbuf}}$	$\frac{\gamma \left(\frac{ D_0 }{P_x} \frac{ D_1 }{P_y} \prod_{i=2}^{n-1} D_i \right)}{S_{minichunk}}$
$1D \rightarrow 1D$	$\gamma \left(\frac{ D_0 }{P} \prod_{i=1}^{n-1} D_i \right)$	$P(T_s + T_w) \frac{\gamma \left(\frac{ D_1 }{P} \prod_{i=2}^{n-1} D_i \right)}{S_{sendbuf}}$	$\frac{\gamma \left(\frac{ D_0 }{P} \prod_{i=1}^{n-1} D_i \right)}{S_{minichunk}}$
$1D \rightarrow 0D$	$\gamma \left(\frac{ D_0 }{P} \prod_{i=1}^{n-1} D_i \right)$	$(T_s + T_w) \frac{\gamma \left(\prod_{i=1}^{n-1} D_i \right)}{S_{sendbuf}}$	$\frac{\gamma \left(\frac{ D_0 }{P} \prod_{i=1}^{n-1} D_i \right)}{S_{minichunk}}$

TABLE 3

Cost Model for Parallel Aggregation Calculations for Aggregating Dimension 2

Dimension	Cost		
	T_{calc}	T_{comm}	$T_{I/O}$
$2D \rightarrow 2D$	$\gamma \left(\frac{ D_0 }{P_x} \frac{ D_1 }{P_y} \prod_{i=2}^{n-1} D_i \right)$	$P_x(T_s + T_w) \frac{\gamma \left(\frac{ D_1 }{P_x} \frac{ D_2 }{P_y} \prod_{i=3}^{n-1} D_i \right)}{S_{sendbuf}}$	$\gamma \left(\frac{ D_0 }{P_x} \frac{ D_1 }{P_y} \prod_{i=2}^{n-1} D_i \right) \frac{1}{S_{minichunk}}$
$2D \rightarrow 1D$	$\gamma \left(\frac{ D_0 }{P_x} \frac{ D_1 }{P_y} \prod_{i=2}^{n-1} D_i \right)$	$P(T_s + T_w) \frac{\gamma \left(\frac{ D_1 }{P_x} \prod_{i=2}^{n-1} D_i \right)}{S_{sendbuf}}$	$\gamma \left(\frac{ D_0 }{P_x} \frac{ D_1 }{P_y} \prod_{i=2}^{n-1} D_i \right) \frac{1}{S_{minichunk}}$
$2D \rightarrow 0D$	$\gamma \left(\frac{ D_0 }{P_x} \frac{ D_1 }{P_y} \prod_{i=2}^{n-1} D_i \right)$	$(T_s + T_w) \frac{\gamma \left(\frac{ D_0 }{P_x} \prod_{i=2}^{n-1} D_i \right)}{S_{sendbuf}}$	$\gamma \left(\frac{ D_0 }{P_x} \frac{ D_1 }{P_y} \prod_{i=2}^{n-1} D_i \right) \frac{1}{S_{minichunk}}$
$1D \rightarrow 0D$	$\gamma \left(\frac{ D_0 }{P} \prod_{i=1}^{n-1} D_i \right)$	$P(T_s + T_w) \frac{\gamma \left(\frac{ D_0 }{P} \prod_{i=2}^{n-1} D_i \right)}{S_{sendbuf}}$	$\gamma \left(\frac{ D_0 }{P} \prod_{i=1}^{n-1} D_i \right) \frac{1}{S_{minichunk}}$

operations are on dense chunks the sizes of the chunks are considered as the product of their dimension cardinalities. However, a true sparse representation designed for both storage and operations, such as BESS, needs to estimate the sizes of the sparse chunks after group-bys.

4.4. Data Structure Management

For large data sets the sizes of the cubes and the number of cubes will not fit in the main memory of the processors. A scalable parallel implementation will require disk space to store results of computations, often many of them with intermediate results. This is similar to a *paging*-based system which can either rely on virtual memory system of the computer or perform the paging of data structures to the needs of the application. We follow the latter approach. Figure 9 shows the data structures for our design and the ones which are paged in and out from disk into main memory on each processor.

A global cube topology is maintained for each subcube by distributing the dimension equally on each processor. A dimension of size d_i , $0 \leq i < n$ gets distributed on p processors, a processor i gets the $\lceil \frac{d_i}{p} \rceil$ portion of d_i , if $i < d_i \bmod p$, or else it gets $\lfloor \frac{d_i}{p} \rfloor$. Each processor thus can calculate what portion belongs to which processor. Furthermore, a constant chunk size is used in each dimension across subcubes. This allows for a simple calculation to find the target chunk which a chunk maps to after aggregating a dimension. However, the first distribution of the dimensions in the base cube is done using sample-based partitioning scheme which may result in an inexact partition and it is kept the same until any of the distributed dimension gets redistributed.

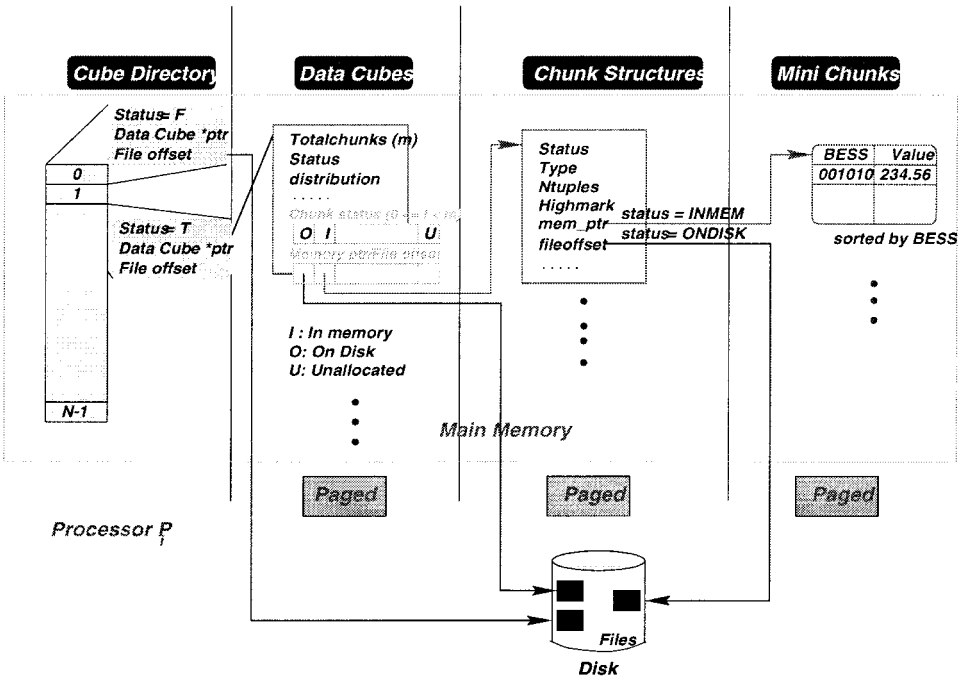


FIG. 9. Data structure management.

A *cube directory* structure is always maintained in memory for each cube at the highest level. For each cube this contains a pointer to a *data cube* structure which stores information about the cube and its chunks. It also contains a file offset to indicate the file address if the data cube structure is paged out. A status parameter indicates whether the data cube structure is in memory (INMEM) or on disk (ONDISK).

A data cube structure maintains the cube topology parameters, the number of unique values in each dimension, whether the chunk structure for the cube is in memory (*status*), a pointer to the chunk structure if it is in memory, and a file offset if it is on disk. The total number of chunks for the chunk structure of the cube is in *totalchunks*. Additionally, for each chunk of the chunk structure, a chunk status *cstatus* is maintained to keep track of chunk structure paging. The chunk address is a pointer to the chunk structure in memory which stores information for each chunk. This is when *cstatus* is set to INMEM. Otherwise, *cstatus* can either be UNALLOCATED or ONDISK. In the latter case the chunk address will be a file offset value. For a multidimensional array, the size of the array and the dimension factor in each dimension are stored for the lookup calculations involving aggregations instead of calculating them on the fly every time.

A chunk structure for a subcube can either be in its entirety or parts of it can be allocated as they are referred to. The *cstatus* field of the data cube will keep track of allocations. A chunk structure keeps track of the number of BESS + value pairs (*ntuples*) in the chunk, which are stored in *minichunks*. Whether a chunk is dense or sparse is tracked by *type*. A dense chunk has a memory pointer to a dense array whereas a sparse chunk has a memory pointer to a mini-chunk. Chunk index for

each dimension in the cube topology is encoded in an 8-byte value $cidx$. Furthermore, dimensions of the chunk are encoded in another 8-byte value $cdim$. This allows for quick access to these values instead of calculating them on the fly.

Minichunks can be unallocated (UNALLOCATED), in memory (INMEM), on disk (ONDISK), or both in memory and on disk (INMEM_ONDISK). Initially, a minichunk for a chunk is allocated memory when a value maps to the chunk (UNALLOCATED \rightarrow INMEM). When the minichunk is filled it is written to disk and its memory reused for the next minichunk (INMEM \rightarrow INMEM_ONDISK). Finally, when a minichunk is purged to disk it is deallocated (INMEM_ONDISK \rightarrow ONDISK). A chunk can thus have multiple minichunks. Hence, choosing the minichunk size is an important parameter to control the number of disk I/O operations for aggregation calculations.

On each processor, there is one file for the chunk structure of each cube, one file for the minichunks of all the sparse chunks in the cube, and one file for the multidimensional chunks of the cube (or if the cube is a multidimensional array). One file is used for the datacube structure. For 10 (20) dimensions there are 1024 (1,048,576) chunk files and 1024 (1,048,576) chunk structure files. With a limit of 2000 open files, file management is required to use the open file descriptors efficiently.

5. ALGORITHMS AND ANALYSIS

Since chunks can be either sparse or dense, we need methods to aggregate sparse chunks with sparse chunks, sparse with dense chunks, and dense with dense chunks. The case of dense chunks to sparse chunks does not arise since a dense chunk never gets converted to a sparse chunk. Also, a chunked organization may be converted into a multidimensional array. The various options are illustrated in Fig. 7. In this section we discuss the algorithms for cube aggregations and chunk mappings. Table 4 describes the notation used in the algorithm description.

TABLE 4
Notation for Algorithms

p	Parent (source) cube
c	Child (target) cube
\mathcal{D}_k	Set of dimensions for cube k
\mathcal{D}_k^d	Set of distributed dimensions for cube k
\mathcal{D}_k^{nd}	Set of non-distributed dimensions for cube k
\mathcal{E}_j^k	Processor extents for cube k in the j th dimension
\mathcal{F}_j^k	Dimension factor of chunks for cube k for the j th dimension
$\langle variable \rangle_A^j$	Value for cube A (parent or child) for the j th dimension
\mathcal{P}_k	Processor set partitioned into \mathcal{P}_k^{di} , $0 \leq i < \mathcal{D}_k^d $ for cube k
\mathcal{B}_{proc}	Buffer for processor $proc$, communication buffer if $proc \neq myid$ else it is a local buffer
\mathcal{R}_{proc_i}	Receive buffer for values from processor $proc_i$, for $0 \leq i < \mathcal{P}_k$, $i \neq myid$
dim_agg_k	Dimension being aggregated in cube k
N_k	Number of chunks for cube k
$chunks_k[i]$	Chunk structure for cube k , $0 \leq i < N_k$
$cidx_j$	chunk index in dimension j
$cdim_j$	chunk extent in dimension j

5.1. Chunk Mapping to Processors

Algorithm 5.1 describes the algorithm for chunked aggregation of a parent cube with sparse chunks. It describes the chunk mapping process and the distinction between *split* and *nonsplit* chunks and *local* and *nonlocal* mappings.

Each chunk in the source cube is processed to map its values to the target chunk. The chunk structure carries information about the chunk's dimensional offsets in *cidx*. This, along with *cdim*, the chunk extents, is used to calculate the local value in each dimension. For distributed dimensions we need to add the start of the processor range to calculate the global value. This is then used to calculate the target start and stop values. This is used to determine the destination target processor and the target chunk. The source can map to the same target chunk on the same processor, map to the same target chunk on another processor, split among chunks on the same processor, or split among chunks on different processors. These cases are illustrated in Fig. 10 for a two-dimensional source to target aggregation of chunks.

Each mapping to a target chunk on the local processor is copied to a local aggregation buffer. If this buffer is full or the current value maps to a different target chunk than the ones in the buffer, it is aggregated with the target. For a non-local calculation a send buffer is kept for each remote processor. A split chunk needs to evaluate each of the index values by decoding the BESS values and adding it to the chunk index values. A target processor needs to be evaluated for the distributed dimensions since this can potentially be different. For a split source chunk, a corrected BESS + value and target chunk-id is sent; otherwise just the BESS + value is sent. Asynchronous send is used to overlap computation and communication. Hence, before the send buffer to a processor is reused, a receive of the previous send must be completed. Asynchronous receive operations are posted from all processors and periodically checked to complete the appropriate sends. A processor receives the BESS values and the target chunk-id and does the aggregation operation.

For a conversion of a chunked source cube to a multidimensional target array, offsets are calculated. Dense chunks are similarly treated. This is explained in a later section.

5.2. Local–Nonlocal Aggregation

Figure 11 illustrates the local and nonlocal aggregation calculations for chunks. The same partitioning dimensions in the source and target subcubes result in a local aggregation calculation. For example, $\underline{ABC} \rightarrow \underline{AB}$ has both *A* and *B* partitioned in both subcubes and this results in a local aggregation. On the other hand, $\underline{ABC} \rightarrow \underline{AB}$ has only *A* partitioned in the result subcube and *B* goes from being distributed to being undistributed. This results in communication and is a nonlocal aggregation. Other cases are illustrated in Table 1.

The distribution of the subcubes is ascertained by the cardinalities of its two largest dimensions. A subcube can be 2D, 1D distributed, or on a single processor (OD). To generate the schedule of data cube calculations, the lattice structure is used, exploiting the various optimizations discussed earlier in the section on schedule

```

for i ← 1 to Np
  for j ← 1 to ndim, j ≠ dim_aggp
    local_start_indexjP = chunksp[i].cidxj * chunk_extentj
    local_stop_indexjP = local_start_indexjP + chunksp[i].cdimj
    /* Calculate the global index in dimension */
    if(j ∈ Dpd)
      global_start_indexjP = local_start_indexjP + EjP[myid]
      global_stop_indexjP = local_stop_indexjP + EjP[myid]
    if(j ∈ Dcd)
      j' = Rank(j, Dcd)
      procb = BinarySearch(global_start_indexjP, EjP[], Pcj')
      local_start_indexjC = global_start_indexjP - EjC[procb]
      proce = procb /* (e)end and (b)beginning processors */
      while(global_stop_indexjP > EjC[proce])
        proce ← proce + 1
      if(procb ≠ proce)
        local_stop_indexjC = EjC[proce]
      else /* For non-distributed dimension j in target */
        local_start_indexjC = global_start_indexjP
        chunk_start_numjC = local_start_indexjC / chunk_extentj
        chunk_stop_numjC = local_stop_indexjC / chunk_extentj
    /* Check if source chunk maps to the same processor/same chunk or same processor/different chunks or different
    processors */
    if(procb = proce and chunk_start_numjC = chunk_stop_numjC, j ∈ Dc)
      target_chunk_id = 0
      for j ← 1 to ndim, j ≠ dim_aggp
        target_chunk_id += chunk_start_numjC * FjC
        /* Calculate the BESS correction values for the target chunk */
        bess_correctionj = |local_start_indexjC - local_start_indexjP|
      if(procb ≠ myid) /* Off-processor value */
        Correct and copy BESS+value pairs into communication buffer Bprocb
        (If Bprocb is full, do an Asynchronous SEND to procb and reuse buffer)
      else /* Local processor value */
        Correct and copy BESS+value pairs into local buffer Bmyid
        (If Bmyid is full, aggregate buffer and reuse)
    else
      if(procb = proce) /* Same processor different chunks */
        for each BESS value in chunk
          target_chunk_id = 0, Initialize prev_target_chunk_id
          for j ← 1 to ndim, j ≠ dim_aggp
            target_chunk_id += {(local_start_indexjP + BESSj) / chunk_extentj} * FjC
          if(procb ≠ myid) /* Off-processor value */
            Copy (BESS+value, target_chunk_id) pairs into communication buffer Bprocb
            (If Bprocb is full, do an Asynchronous SEND to procb and reuse buffer)
          else /* Local processor value */
            if(target_chunk_id ≠ prev_target_chunk_id)
              Aggregate local buffer Bmyid with target_chunk_id of target
            else
              Copy BESS value to local buffer Bmyid
              prev_target_chunk_id = target_chunk_id
      else
        /* Different processors and hence different chunks */
        for each BESS value in chunk
          target_chunk_id = 0, Initialize prev_target_chunk_id
          for j ← 1 to ndim, j ≠ dim_aggp
            global_indexj = global_start_indexjP + BESSj
            if(j ∈ Dcd)
              proc = BinarySearch(global_indexj, EjP[], Pcj')
              local_indexj = global_indexj - EjP[proc]
            else
              local_indexj = global_indexj
            target_chunk_id += (local_indexj / chunk_extentj) * FjC
          if(proc = myid) /* Local processor value */
            if(target_chunk_id ≠ prev_target_chunk_id)
              Aggregate local buffer Bmyid with target_chunk_id of target
            else
              Copy BESS value to local buffer Bmyid
              prev_target_chunk_id = target_chunk_id
          else /* Off-processor value */
            Copy (BESS+value, target_chunk_id) pairs into communication buffer Bproc
            (If Bproc is full, do an Asynchronous SEND to proc and reuse buffer)
    Asynchronous RECEIVE from other processors in Rproci for 0 ≤ i < |Pp|, i ≠ myid, and aggregate to
    local chunks.
end

```

Algorithm 5.1. Chunked cube mapping.

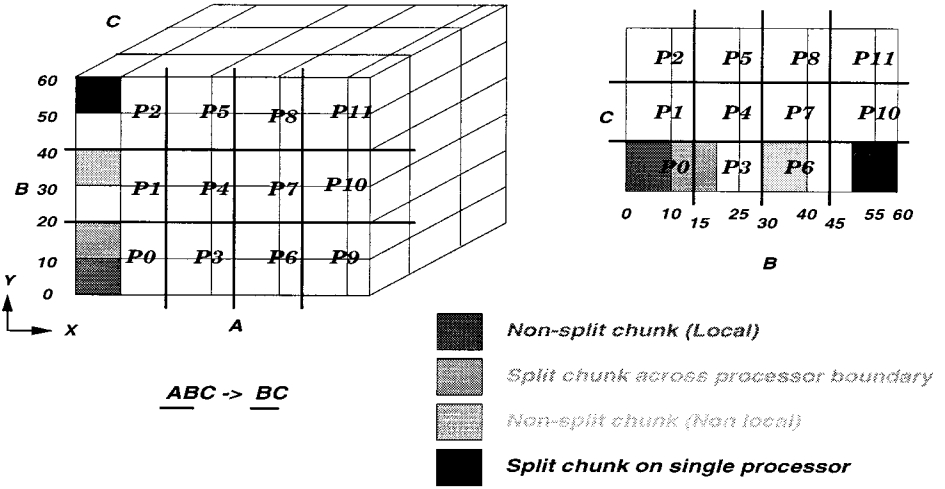


FIG. 10. Chunk mapping after aggregation for a 2D cube distribution.

generation. Computation and communication costs for aggregation are used by the scheduler by analyzing the chunk aggregation operations and the resulting partitioning. The cost analysis of the various aggregations between 1D, 2D, and uniprocessor aggregations is not included in this article.

The extents of a chunk of the aggregating cube can be contained in the extents of a chunk of the aggregated cube. In this case the BESS + value pairs are directly mapped to the target chunk, locally or nonlocally. However, the BESS index values need to be modified to encode the offsets of the new chunk. Figure 12 illustrates a case for $\underline{AB} \rightarrow \underline{A}$ where the chunks with \underline{A} extents of 10–13 on processor P0 map to the chunk with extents 9–13 on processor P1. The BESS encoding of \underline{A} needs to be incremented by 1 to correctly reflect the target BESS encoding. If the chunk is overlapping over a target chunk boundary, then each BESS value has to be extracted to determine its target chunk. This is computationally more expensive than the

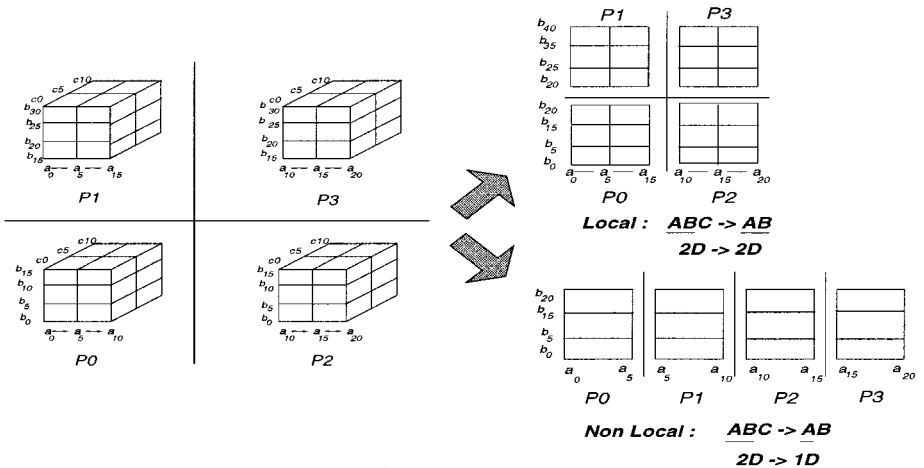


FIG. 11. Local and Nonlocal aggregation calculations.

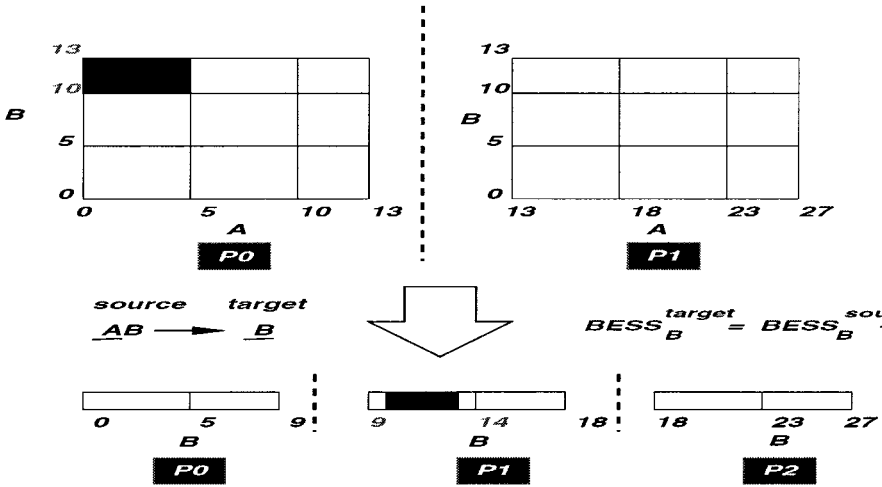


FIG. 12. Modification of BESS indices while mapping a nonsplit chunk.

direct case. It is to be noted that a two-dimensional distribution may result in more overlapped chunks than a one-dimensional distribution, because the former has more processor boundary area than the latter.

5.3. Sparse-Sparse Chunk Aggregation

Algorithm 5.2 describes the aggregation algorithm of a sparse chunk with a sparse chunk. Sparse chunks store BESS + value pairs in minichunks. Sparse to sparse aggregations involve accessing these minichunks. The BESS values are kept sorted in the minichunks to facilitate the aggregation calculations by using sort and

```

for i ← 1 to  $N_p$ 
  visited[i] = FALSE;
done = FALSE
Initialize prev_target_chunk_id
for i ← 1 to  $N_p$ 
  nextchunk = i
  if(! visited[nextchunk])
    do {
      /* Check if this is the chunk mapping to the same target */
      for each BESS value in nextchunk
        if(Buffer is not full and target_chunk_id = prev_target_chunk_id)
          for j ← 1 to ndimp
            Mask dim_aggp from BESS index and copy to memory buffer
        else
          Sort memory buffer on masked BESS indices
          Merge sorted buffer with sorted BESS values of target_chunk_id
          Empty the buffer for reuse
      /* If dimension oriented chunk access, get the next chunk along dimension dim_aggp */
      if(chunk_access = DIM_ORIENTED)
        if((nextchunk ← Get_next_chunk(dim_aggp, nextchunk)) = FALSE)
          done = TRUE
    } else
      done = TRUE
      visited[nextchunk] = TRUE
  } while (! done)
end

```

Algorithm 5.2. Sparse-sparse chunk aggregation algorithm.

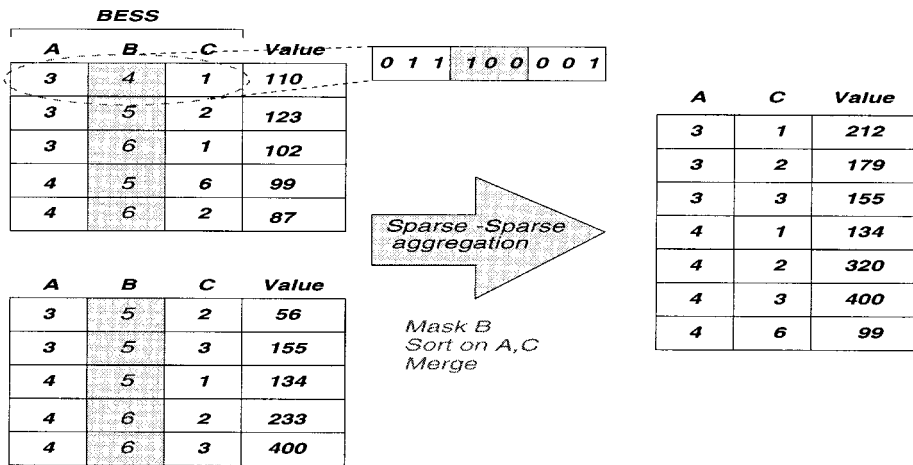


FIG. 13. Sparse-sparse chunk aggregation.

scan operations used in relational processing. Figure 13 shows an aggregation in which *B* is aggregated in *ABC* to give an *AC* subcube. Here the bit encoding for *B* is masked from the BESS values for both chunks, which means we are aggregating along dimension *B*. An integer sort is done on the remaining encoding of *A* and *C* on both the chunks. This gets the values which map to the same *A* and *C* contiguous to each other. This is followed by a merge of sorted values aggregating where the values of *A* and *C* are the same.

Buffer management for this aggregation depends on the order in which chunks to be aggregated are accessed. There are two alternative ways to access chunks. The first is *dimension oriented*, in which the chunks are accessed along the dimension to be aggregated. The other method, *chunk numbering*, accesses chunks in the order in which they are laid out in memory. Figure 14 illustrates these two choices. In the dimension oriented access, a buffer can be filled with the source chunks that map

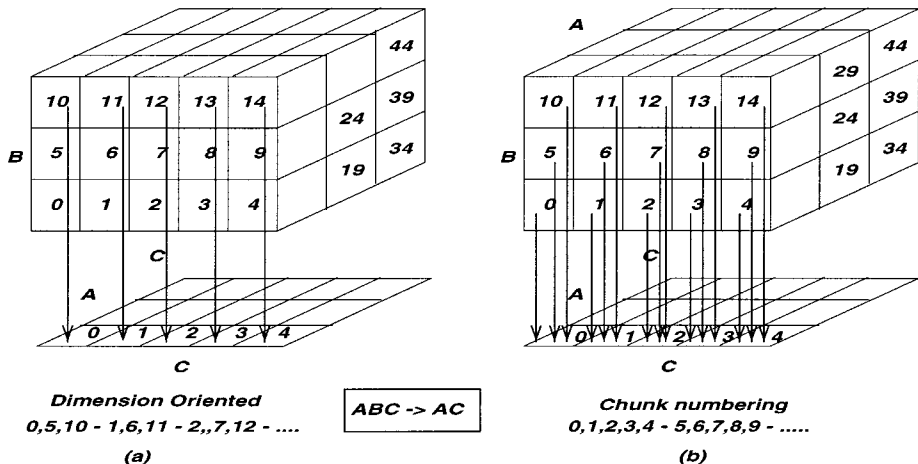


FIG. 14. Dimension oriented vs. chunk numbering access of chunks.

to a single target chunk and the sort and scan operations are done on the buffer in memory. For the chunk order access, the consecutive chunks accessed are mapped to different chunks for all dimensions except the innermost dimension. This results in a sort of each individual chunk and a merge with the sorted values calculated earlier mapping to the same chunk. This involves a disk read to get the previous result from the minichunk of the target chunk, which was written to disk if it was full.

It is seen that the dimension oriented access performs better than chunk ordering since the sort can be done more efficiently than the sort-merge operation which involves disk I/O.

5.4. Sparse-Dense Chunk Aggregation

A sparse chunk can be aggregated to a dense chunk by converting the BESS dimension encodings to a chunk offset value. The chunk structure has the values for the index of the chunk in each dimension in the variable *cidx*. This is used to calculate the offset of the chunk in the resulting cube topology by using the number of chunks in each dimension of the resulting subcube.

The BESS index values are extracted and each dimension value is multiplied with the appropriate dimension factor for the chunk stored in the chunk structure. This converts the BESS values to an offset within the chunk which is added to the chunk offset after adjusting the source chunk offset with the resulting chunks offset. The value is then aggregated to this location in the resultant dense chunk.

5.5. Chunked-Multidimensional Array Aggregation

A chunk is aggregated to a dense multidimensional array in a way similar to the sparse-dense chunk described above. The chunk offset is calculated from *cidx* and the BESS value is used to calculate the offset within the chunk. This is added to the chunk offset which is the correct offset in the resulting multidimensional array. This could be either local or nonlocal. Local values are aggregated to the correct offset. Nonlocal offsets are copied to a buffer for each processor along with the value at that location. Once the buffer is full it is sent to the appropriate processor, which receives it and performs the aggregation. Algorithm 5.3 describes the algorithm in more detail.

5.6. Cube and Chunk Level Reuse and I/O

Calculating simultaneous multi-dimensional aggregates involves calculations for multiple cubes that are either chunked or stored as multidimensional arrays. Since a large number of cubes are calculated and disk storage is used for minichunks, chunk structures, and data cube structures, strategies to reduce the number of scans for an intermediate calculation should reuse data in memory effectively. This is determined by the order of calculation of subcubes in the lattice structure. Figure 15 illustrates the order by a depth first and breadth first traversal of the lattice structure. In Fig. 15a, the order of computation is $ABC \rightarrow AB \rightarrow A$, reusing AB to calculate A before calculating AC from ABC . The result of an aggregation can be

```

for i ← 1 to  $N_p$ 
  for j ← 1 to  $ndim$ ,  $j \neq dim\_agg_p$ 
     $local\_start\_index_j^p = chunks_p[i].cid_x * chunk\_extent$ 
     $local\_stop\_index_j^p = local\_start\_index_j^p + chunks_p[i].cdim$ 
    /* Calculate the global index in dimension */
    if ( $j \in \mathcal{D}_p^d$ )
       $global\_start\_index_j^p = local\_start\_index_j^p + \mathcal{E}_j^p[myid]$ 
       $global\_stop\_index_j^p = local\_stop\_index_j^p + \mathcal{E}_j^p[myid]$ 
    if ( $j \in \mathcal{D}_c^d$ )
       $j' = Rank(j, \mathcal{D}_c^d)$ 
       $proc_j^b = BinarySearch(global\_start\_index_j^p, \mathcal{E}_j^p[], \mathcal{P}_c^{j'})$ 
       $local\_start\_index_j^c = global\_start\_index_j^p - \mathcal{E}_j^c[proc_j^b]$ 
       $proc_j^e = proc_j^b$  /* (e)end and (b)beginning processors */
      while ( $global\_stop\_index_j^p > \mathcal{E}_j^c[proc_j^e]$ )
         $proc_j^e \leftarrow proc_j^e + 1$ 
      if ( $proc_j^b \neq proc_j^e$ )
         $local\_stop\_index_j^c = \mathcal{E}_j^c[proc_j^e]$ 
      else /* For non-distributed dimension j in target */
         $local\_start\_index_j^c = global\_start\_index_j^p$ 
    /* Check if source chunk maps to the same processor or different processors */
    if ( $proc_b = proc_e$ ) /* Same processor */
      for j ← 1 to  $ndim$ ,  $j \neq dim\_agg_p$ 
         $target\_offset += local\_start\_index_j^c * \mathcal{F}_j^c$ 
      if ( $proc_b \neq myid$ ) /* Off-processor value */
        Copy  $target\_offset$ +value pairs into communication buffer  $\mathcal{B}_{proc_b}$ 
        (If  $\mathcal{B}_{proc_b}$  is full, do an Asynchronous SEND to  $proc_b$  and reuse buffer)
      else /* Local processor value */
        Copy  $target\_offset$ +value pairs into local buffer  $\mathcal{B}_{myid}$ 
        (If  $\mathcal{B}_{myid}$  is full, aggregate buffer and reuse)
    else
      for each BESS value in chunk
         $target\_offset = 0$ 
        for j ← 1 to  $ndim$ ,  $j \neq dim\_agg_p$ 
           $global\_index_j = global\_start\_index_j^p + BESS_j$ 
          if ( $j \in \mathcal{D}_c^d$ )
             $proc = BinarySearch(global\_index_j, \mathcal{E}_j^p[], \mathcal{P}_c^{j'})$ 
             $local\_index_j = global\_index_j - \mathcal{E}_j^p[proc]$ 
          else
             $local\_index_j = global\_index_j$ 
             $target\_offset += (local\_index_j / dimension\_size_j^{proc}) * \mathcal{F}_j^c$ 
          if ( $proc = myid$ ) /* Local processor value */
            if (local buffer  $\mathcal{B}_{myid}$  is full)
              Aggregate local buffer  $\mathcal{B}_{myid}$  with  $target\_chunk\_id$  of target
            else
              Copy  $target\_offset$ +value to local buffer  $\mathcal{B}_{myid}$ 
          else /* Off-processor value */
            Copy  $target\_offset$ +value to communication buffer  $\mathcal{B}_{proc}$ 
            (If  $\mathcal{B}_{proc}$  is full, do an Asynchronous SEND to  $proc$  and reuse buffer)
        Asynchronous RECEIVE from other processors in  $\mathcal{R}_{proc_i}$  for  $0 \leq i < |P_p|$ ,  $i \neq myid$ , and
        aggregate to local address offsets.
  end
end

```

Algorithm 5.3. Chunked–multidimensional array cube algorithm.

reused in this case. On the other hand, in Fig. 15b the source ABC can be reused by calculating AB , AC , BC from ABC before calculating the next level.

Only a minichunk is in memory for a particular chunk; the remaining minichunks (if any) are stored on disk. A result of an aggregation results in a disk write for the target minichunk when it gets full. A depth first traversal, which reuses the resultant chunk, will still result in disk reads for chunks which have many minichunks. This is a consequence of having only a minichunk size memory buffer for each chunk for scalability reasons. A larger minichunk size would have data of more resultant chunks in memory to be reused.

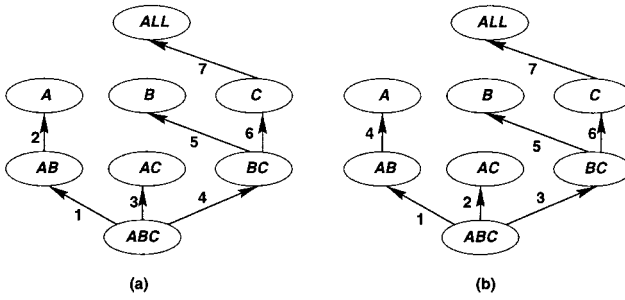


FIG. 15. (a) Depth first traversal of lattice. (b) Breadth first traversal of lattice.

In case of the breadth first traversal, the source minichunks are read for a chunk in a memory buffer and aggregated with the target. Once this is done, the source buffer is released, again to accommodate many such calculations. At a cube level reuse, if the minichunk size can fit the values of the source chunk, they will remain in memory after the $ABC \rightarrow AB$ calculation. This can be reused to calculate the other aggregates from the source. In case there are multiple minichunks for a chunk not much reuse is possible at the cube level. However, at a chunk level reuse, the memory buffer allocated for ABC 's chunks aggregation can be reused to calculate the target of AB , AC , and BC by doing the appropriate maskings. Again, this will hold true only in the case of chunk accesses by chunk number and not for dimension oriented chunk access. This is because the chunks are read in dimension order and this order is different for $ABC \rightarrow AB$ (dimension order is C), $ABC \rightarrow AC$ (dimension order is B) and $ABC \rightarrow BC$ (dimension order is A). Thus, reuse at chunk-level does not take advantage of the dimension-oriented chunk access optimizations.

5.7. Interprocessor Communication

The processor with the least load is assigned to a uniprocessor subcube for load balancing. To perform a parallel aggregation calculation every processor cooperates since each contains a part of the source and possibly the target cube. Since there are dependencies between any two levels of the lattice, i.e. between the source and the target, the latter cannot be used for any-further calculation unless it is calculated completely. This is especially important when the calculation is in parallel since other processors contribute to the aggregated values. Each processor goes through the chunks assigned to it independently and communicates with other processors by sending and receiving messages. Asynchronous sends and receives are used to overlap computation with communication and reduce the synchronization that results from waiting for messages. For example, when a processor needs to send a message in buffer B , it posts an asynchronous send and goes on to perform its processing. But before modifying B again it needs to check that the message has been sent. Hence, each processor also posts an asynchronous receive to get values from other processors. Once data are received from a particular processor, another asynchronous receive is posted for that processor. Processors exchange sentinel values to signal the end of the current communication phase. Algorithm 5.4 shows the communication phase on each processor.

```

for  $i \leftarrow 0$  to  $P$ ,  $i \neq myid$ 
  Post Asynchronous RECEIVE from processor  $i$ 
  Initialize  $send\_request[i]$  to FALSE
for  $j \leftarrow 1$  to  $N_p$  (for each chunk)
  ...Computation...
  Find target processor for chunk as  $proc$ 
  While ( $send\_request[proc] = TRUE$ )
    Wait for  $b$  to be received at  $proc$  before modifying  $b$ 
  if( $target_{proc} \neq myid$ )
    Post Asynchronous SEND to processor  $proc$ 
    Set  $send\_request[proc]$  to TRUE
  for  $i \leftarrow 0$  to  $P$ ,  $i \neq myid$ 
    If( $Test\_for\_receives[i] = TRUE$ )
      Receive and aggregate buffer
      [This will set  $send\_request[myid]$  to FALSE on processor  $i$ ]
      Post Asynchronous RECEIVE from processor  $i$ 
end

```

Algorithm 5.4. Communication in source \rightarrow target aggregation.

6. DATA MINING

Data mining can be viewed as an automated application of algorithms to detect patterns and extract knowledge from data [8]. An algorithm that enumerates patterns from or fits models to data is a data mining algorithm. Data mining is a step in the overall concept of knowledge discovery in databases (KDD). Large data sets are analyzed for searching patterns and discovering rules. Automated techniques of data mining can make OLAP more useful and easier to apply in the overall scheme of decision support systems. Data mining techniques such as *Associations, Classification, and Clustering* [8] can be used together with OLAP to discover knowledge from data. In the next few sections we present an integration of association rule mining and classification with the parallel OLAP and multidimensional analysis infrastructure presented in this article.

6.1. Association Rule Mining

Discovery of quantitative rules is associated with quantitative information from the database. The data cube represents quantitative summary information for a subset of the attributes. Attribute-oriented approaches [2, 17, 18] to data mining are data-driven and can use this summary information to discover association rules. Transaction data can be used to mine association rules by associating *support* and *confidence* for each rule. Support of a pattern A in a set S is the ratio of the number of transactions containing A and the total number of transactions in S . Confidence of a rule $A \rightarrow B$ is the probability that pattern B occurs in S when pattern A occurs in S and can be defined as the ratio of the support of AB and support of A . The rule is then described as $A \rightarrow B$ [*support, confidence*] and a *strong* association rule has a support greater than a predetermined minimum support and a confidence greater than a predetermined minimum confidence. This can also be taken as the measure of *interestingness* of the rule. Calculation of support and confidence for the rule $A \rightarrow B$ involve the aggregates from the cube AB , A , and ALL. Additionally,

dimension hierarchies can be utilized to provide multiple level data mining by progressive generalization (roll-up) and deepening (drill-down). This is useful for data mining at multiple concept levels and interesting information can potentially be obtained at different levels.

An approach to data mining called *attribute focusing* targets the end-user by using algorithms that lead the user through the analysis of data. Attribute focusing has been successfully applied in discovering interesting patterns in the NBA [2] and other applications. Earlier applications of this technique were to software process engineering [3]. Since data cubes have aggregation values on combinations of attributes already calculated, the computations of attribute focusing are greatly facilitated by data cubes. We present a parallel algorithm to calculate the interestingness function used in attribute focusing on the data cube.

Consider a three-attribute data cube with attributes A , B , and C , defining $E_3 = ABC$. For showing two-way associations, we will calculate the interestingness function between A and B , A and C , and B and C . When calculating associations between A and B , the probability of E , denoted by $P(AB)$ is the ratio of the aggregation values in the subcube AB and ALL. Similarly the independent probability of A , $P(A)$, is obtained from the values in the subcube A , dividing them by ALL. $P(B)$ is similarly calculated from B . The calculation $|P(AB) - P(A)P(B)| > \delta$, for some threshold δ , is performed in parallel. Assume that the cubes AB and A are distributed along the A dimension; no replication of A is needed. However, since B is distributed in subcube B , and B is local on each processor in AB , B needs to be replicated on all processors. AB and A cubes are distributed, but B is replicated on all processors. Figure 16 shows a sample calculation of $P(AB)$, $P(A)$, and $P(B)$ on three processors. A sample calculation is highlighted in the figure, $|0.03 - 0.22 \times 0.08| = 0.0124$, which is greater than δ values from 0.001 to 0.01, and the corresponding attribute values are associated within that threshold.

Aggregates with m -attributes and below present in the data cube can be used to calculate m -way associations by looking at subcubes for all combinations of m dimensions. Since these cubes are precomputed, calculations are very efficient. We do not discuss this further in this article. Next, we describe algorithms for classification on the multidimensional structure.

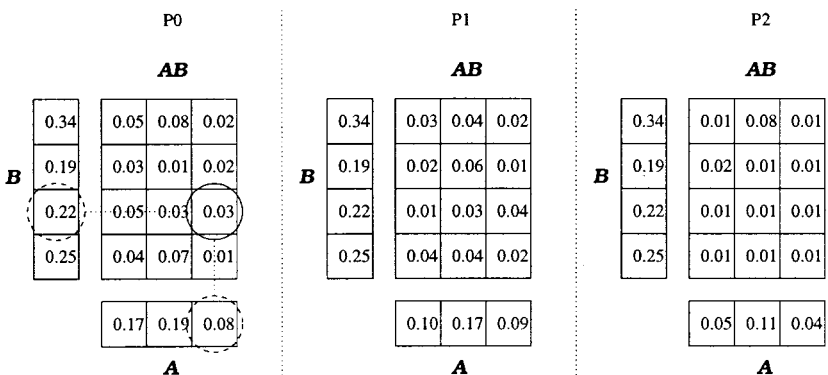


FIG. 16. Use of aggregate subcubes AB , A , and B for calculating interestingness on three processors.

6.2. Classification

A set of sample records called the *training data* set is given consisting of several attributes. Attributes can either be *continuous*, if they come from an ordered domain, or *categorical*, if they are from an unordered domain. One of the attributes is the *classifying* attribute that indicates the *class* to which the record belongs. The objective of classification is to build a model of the classifying attribute based upon the other attributes of the record.

Several classification models have been used in the past, notably neural networks [24], genetic algorithms [15], decision trees [4, 26] and some others. Among these models, the decision tree models are considered to be the most useful in the domain of data mining because they are relatively inexpensive to construct, easy to interpret, and easy to integrate with database systems. Also, for a variety of problem domains they yield comparable or better accuracy as compared to the other models [24].

A decision tree recursively partitions the training set until each partition consists entirely or dominantly of records from one class. Each *nonleaf* node of the tree contains a split point which is a test on an attribute and determines how the data are partitioned. Once the decision tree is built from the training set it can be used to classify future instances. The decision-tree-based classifiers that can handle large data sets are important because use of larger data sets improves the classification accuracy [24].

Previous work in classifying large data sets has been to use sampled data sets or multiple partitions of the data set [5, 24]. Recent work has focused on using the entire data set, in classifiers such as SLIQ [23] and SPRINT [29]. A parallel classifier in the same spirit has been developed in ScalParC [20]. Classifiers such as CART [4] and C4.5 [26] perform sorting at every node of the decision tree, which makes them expensive for large data sets since disk-based sorting has to be done. The approach of SPRINT and SLIQ is to sort the continuous attribute once in the beginning and maintain the sorted order in the subsequent splitting steps. Separate lists are kept for each attribute which maintains a record identifier for each sorted value. In the splitting phase the same records need to be assigned to a node, which may be in a different order in the different attribute lists. A hash table is used to provide a mapping between record identifiers and the node to which it belongs after the split. This mapping is then probed to split the attribute lists in a consistent manner.

Table 5 is an example training set with three attributes, Age, Car color, and Gender, and a class attribute. Figure 17a shows its classification tree. At each node the attribute that best divides the training set is chosen to split. Several splitting criteria have been used in the past to evaluate the goodness of a split. Calculating the *gini* index is commonly used [4]. This involves computing the frequency of records of each class in each of the partitions. If a parent node having n records and c possible classes is split into p partitions, the *gini* index of the i th partition is $gini_i = 1 - \sum_{j=1}^c (n_{ij}/n_i)^2$. n_i is the total number of records in partition i , of which n_{ij} records belong to class j . The *gini* index of the total split is given by $gini_{split} = \sum_{i=1}^p (n_i/n) gini_i$. The attribute with the least value of $gini_{split}$ is chosen to split the records at that node. The matrix n_{ij} is called the *count matrix*. The count matrix needs to be calculated for each evaluated split point for a continuous attribute.

TABLE 5
Training Set

Row-id	Age	Car-color	Gender	Class-id
0	10	Green	F	0
1	50	Blue	M	1
2	40	Yellow	F	0
3	30	Green	F	0
4	20	Red	M	1
5	40	Blue	M	0
6	20	Yellow	M	1

Categorical attributes have only one count matrix associated with them and hence computation of the gini index is straightforward. For the continuous attributes an appropriate splitting value has to be determined by calculating the $gini_{split}$ and choosing the one with the minimum value. If the attribute is sorted then a linear search can be made for the optimal split point by evaluating the gini index at each attribute value. The count matrix is calculated at each possible split point to evaluate the $gini_{split}$ value. The $gini$ index calculations and the node splits for the example above are given in Fig. 18. At node 0, the attribute *Gender* yields the optimal $gini_{split}$ value of 0.214. This creates a split with one partition with *M* values for gender and another with *F* values. After this split is made, two child nodes are created. The record values need to be partitioned consistently between the two nodes for the *split-attribute* and the *nonsplit attributes*. Splitting the split-attribute is straightforward by adjusting pointer values. The challenge is to split the nonsplit attributes efficiently. Existing implementations such as SPRINT and ScalParC maintain a mapping of the row-id and class-id with the values assigned to each node. The values are split physically among nodes such that the continuous attributes maintain their sorted order in each node to facilitate the sequential scan

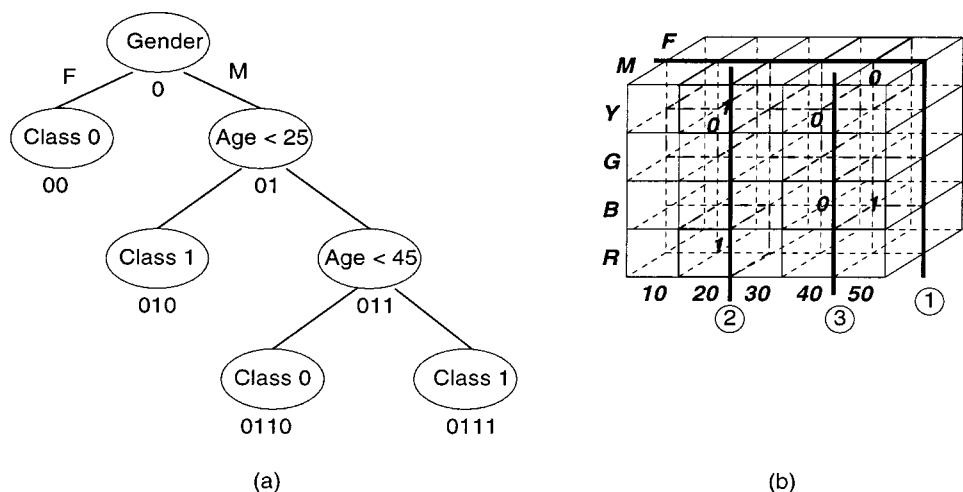


FIG. 17. (a) Classification tree for training set. (b) Classification tree embedded on the base cube.

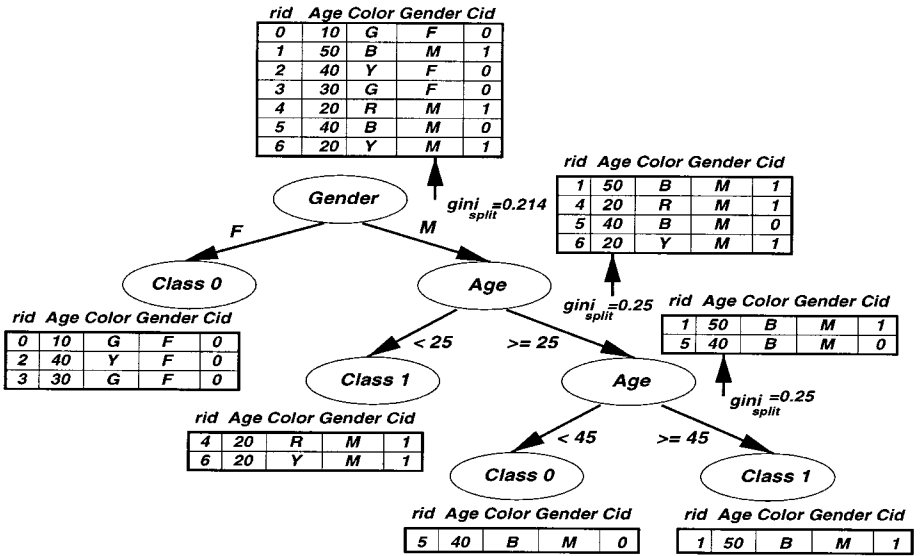


FIG. 18. Gini index calculation for the attributes and node splitting.

for the next split determination phase. A hash list maintains the mapping of record ids to nodes. The record ids in the lists for nonsplitting attributes are searched to get the node information and perform the correct split.

6.2.1. *Classification on the multidimensional cube structure.* We propose that classification trees can be built using structure imposed on data using the multidimensional data model. Gini index calculation relies on the count matrix which can be efficiently calculated using the dimensional model. Each populated cell represents a record in the array. For the base cube (which is a multidimensional representation of the records without any aggregation) the class value of the record is stored in each cell. The gini index calculation uses the count matrix which has information about the number of records in each partition belonging to each possible class.

To evaluate split points for a continuous attribute the $gini_{split}$ needs to be evaluated for each possible split point in a continuous attribute and once for a categorical attribute. This means the aggregate calculations present in each of the one-dimensional aggregates can be used if they have a number of records belonging to each class. Therefore for each aggregate we store the number of records in each class. Figure 19a gives an example training set with two dimensions, A , a continuous dimension, and B , a categorical dimension, and two class values 0 and 1.

Figure 19b is the corresponding multidimensional model. The continuous dimensions A are stored in the sorted order. The aggregates store the number of records mapping to that cell for both classes 0 and 1. To calculate the $gini_{split}$ for the continuous attribute A it is now easy to look at the A aggregate and sum the values belonging to both classes 0 and 1 on both sides of the split point under consideration to get the count matrix. Gini index calculation is done on an attribute list which in the case of a multidimensional model is a dimension. Count matrix is

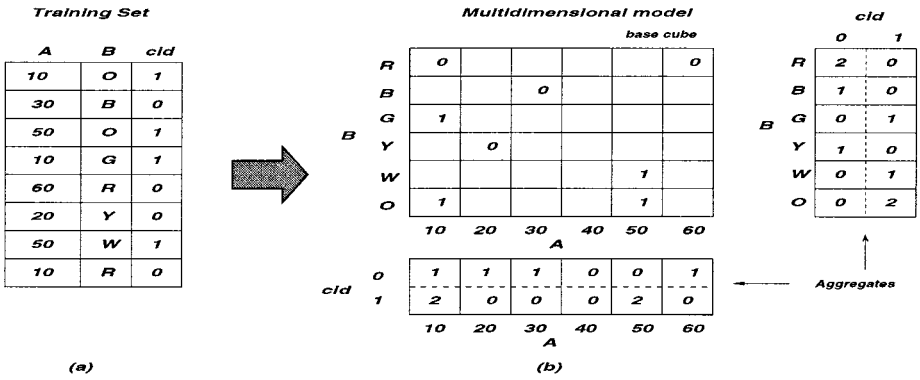


FIG. 19. (a) Training set record and (b) corresponding multidimensional model.

repeatedly calculated on the sorted attribute list which is readily available in the cube structure as a higher level one-dimensional aggregate. Each dimension is sorted in the dimensional structure as shown in Fig. 19b. For a parallel version, the count matrices are computed in parallel and each processor maintains the classification tree nodes, identifying the portion of the hypercube present locally. Each processor computes the count matrices locally for each active node of the tree. Collective communication is used to consolidate the global count matrices. Communication is $O(\mathcal{N}_n^k |d|)$ at a level k , where \mathcal{N}_n^k is the number of active nodes to classify at level k , and $|d|$ is the size of the largest dimension. This is much less than the communication time of $O(N/p)$ ($O(N)$ worst case) at each level for each attribute in the attribute-list approach discussed earlier. Further details can be found in [14].

Figure 20 illustrates the classification tree building process using the multidimensional model and the aggregates maintained at the highest level of the cube structure, one for each dimension.

6.3. Partial Cubes

Partial cube calculations are performed to calculate all the aggregate subcubes at level 1, 2, or 3 of the cube lattice as needed in data mining. This is done by maintaining the minimum count of intermediate calculations. In a full data cube, the number of aggregates is exponential. In the cube lattice there are n levels and for each level i there are $\binom{n}{i}$ aggregates at each level containing i dimensions. The total number of aggregates is thus 2^n . For a large number of dimensions, this leads to a prohibitively large number of cubes to materialize.

Next, we need to determine the number of intermediate aggregates needed to calculate all the aggregates for all levels below a specified m , $m < n$, for m -way associations and for metarules with m predicates. The number of aggregates at level 2 is $\sum_{j=1}^{n-1} \binom{n}{j}$. For level 3 it is $\sum_{j=1}^{n-2} \binom{n}{j}$. At each level k , the first $k-1$ literals in the representation are fixed and the others are chosen from the remaining. The total number of aggregates is then $\sum_{i=1}^{n-1} \sum_{j=1}^i \binom{n}{j}$. After simplification this becomes $\frac{n(n-1)(2n+2)}{12}$. For example, when $n = 10$, this number is 165, a great reduction from $2^{10} = 1024$. The savings are even more significant for $n = 20$, where this number is

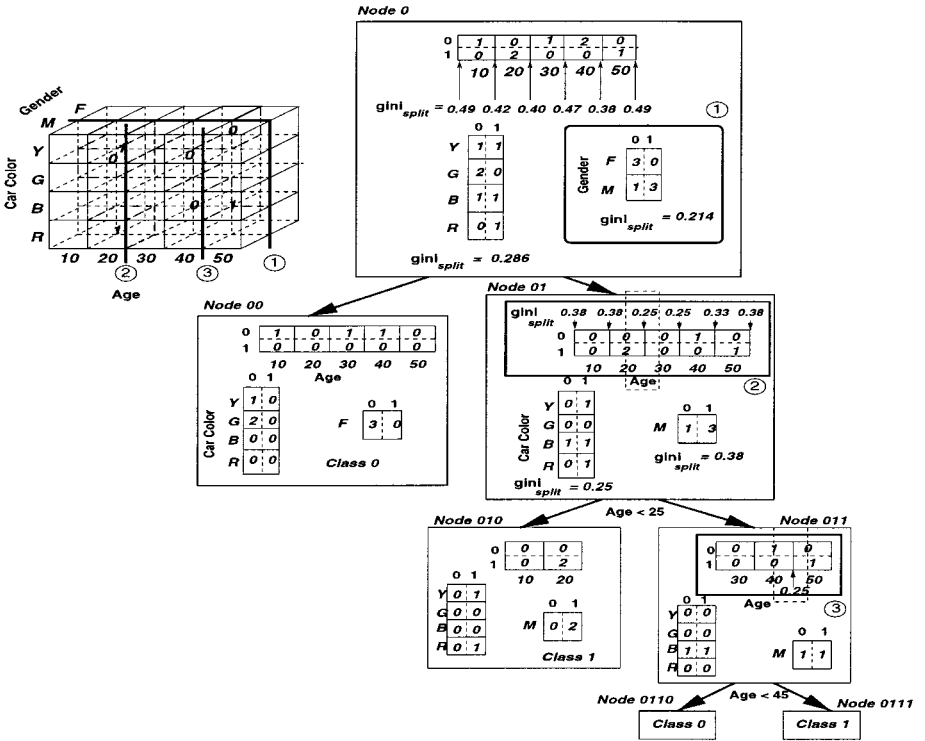


FIG. 20. Gini calculation for the attributes and node splitting with multidimensional aggregates.

1330, much smaller than $2^{20} = 1,048,576$. A similar analysis of the number of cubes needed for calculating all the one-attribute aggregates at level 1 of the cube lattice is given in [14].

7. PERFORMANCE RESULTS AND DISCUSSION

In this section we present performance results for our system on a 16-node IBM SP-2 distributed memory parallel computer available to us at Northwestern University. Assume N tuples and p processors. Initially, each processor reads $\frac{N}{p}$ tuples from a shared disk, assuming that the number of unique values is known for each attribute. These are partitioned using a sample-based partitioning algorithm (*partitioning phase*) so that the attribute (dimension) values are ordered on processors and distributed almost equally. To load the base cube, tuples are sorted (*sorting phase*) on the combined key of all the attributes so that the access to chunks is conformant to its layout in memory or disk. (Sorting in the order $A_0 \rightarrow A_1 \rightarrow A_2 \dots A_{n-1}$ is conformant to the layout of chunks where A_0 is the outermost dimension and A_{n-1} is the innermost for loading a sorted run of values.)

The base cube is loaded on each processor (Fig. 8) from these tuples locally on each processor. The subcubes of the data cube are calculated from here (*building phase*). This is followed by the analysis and data mining phase on the computed aggregates.

We choose four data sets, one each of dimensionality 3, 5, 10, and 20 to illustrate performance. Random data with a uniform distribution are currently used for the

TABLE 6
Description of Datasets and Attributes

Data set	Dim	Cardinalities (d_i)	$(\prod_i d_i)$	Tuples
Dataset I	3	1024(S), 256(N), 512(N)	2^{27}	10 M
Dataset II	5	1024(S), 16(N), 32(N), 16(N), 256(S)	2^{31}	1 & 10 M
Dataset III	10	1024(S), 16(S), 4(S), 16, 4, 4, 16, 4, 4, 32(N)	2^{37}	5 & 10 M
Dataset IV	20	16(S), 16(S), 8(S), 2, 2, 2, 2, 4, 4, 4, 4, 8, 2, 8, 8, 8, 2, 4, 1024(N)	2^{51}	1 M

Note. (N) Numeric (S) String, M, million.

performance figures. We have evaluated other types of data, skewed data sets and the OLAP council benchmark [25], which model a realistic business situation with similar results.

Table 6 illustrates the data sets for our experiments. The five-dimensional data set has chunk sizes in each dimension of 32, 4, 8, 8, 16 for a total chunk dimension product of 2^{17} . The number of chunks is 2^{14} for the base cube and is calculated for each subcube depending on the dimensions it has, again getting distributed across processors. Finally, for the ten-dimensional data set the chunk dimension sizes are 64, 4, 2, 4, 2, 2, 4, 2, 2, 4 or a total chunk dimension product of 2^{19} and total number of chunks distributed across processors as 2^{18} for the base cube. The other subcube chunk sizes and chunk structure sizes are calculated from these.

The number of subcubes in the datacube for Dataset I is $2^3 = 8$, Dataset II is $2^5 = 32$, and Dataset III is $2^{10} = 1024$. We report the results of complete data cube construction for these sets. For Dataset IV we report the results of partial data cube construction in which 1350 subcubes are calculated.

Figure 21 shows the time taken by the various phases of the data cube construction algorithm. Each phase of the cube construction process shows good speedup for all the data sets when a single dimension is partitioned in the base cube. A two-dimensional partitioning performs better than a one-dimensional partitioning

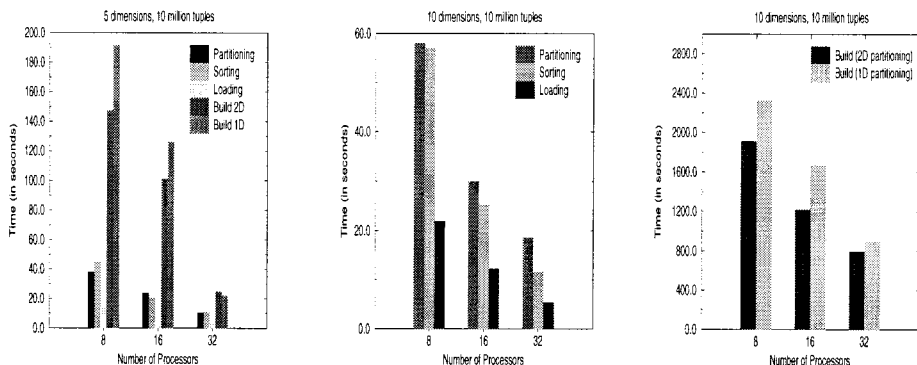


FIG. 21. Time taken by various phases of the data cube construction algorithm for 5 (32 subcubes) and 10 (1024 subcubes) dimensions.

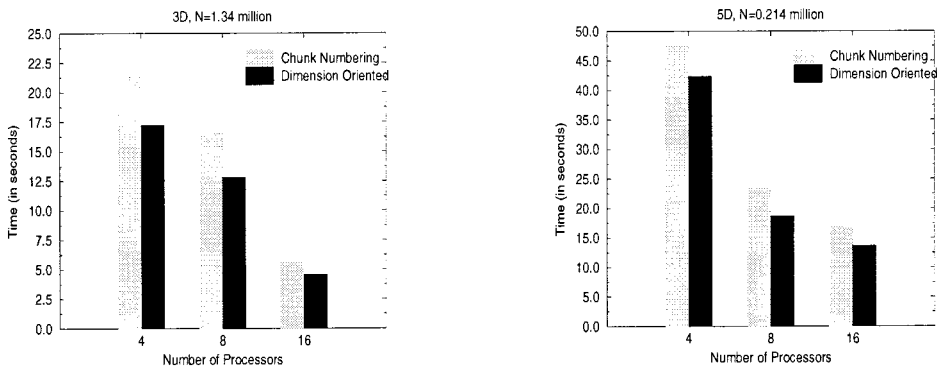
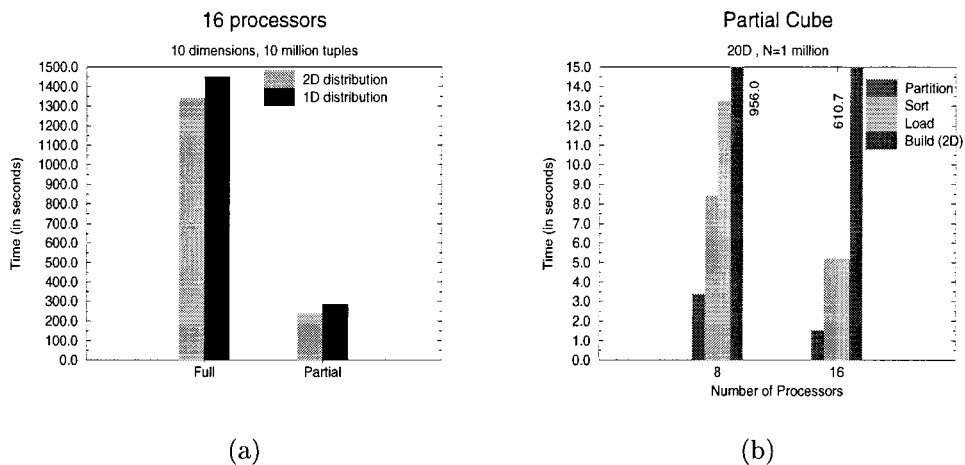


FIG. 22. Comparison of the two chunk access methods, chunk numbering and dimension oriented for a dataset with 3 and 5 dimensions.

because there are more chunks in the partitioned dimension that allows for better sparse-sparse aggregation performance.

Figure 22 shows the time taken by the two different methods of accessing chunks for aggregation calculations. Dimension oriented method is better than the access that follows chunk numbering because of better buffer use in the sparse aggregation calculations. In all cases we observe that the dimension ordering method works better than the access pattern that follows chunk numbering. This effect is more pronounced on a lower number of processor and lower dimensional data sets because the effect of amortizing sort costs among chunks mapping to the target chunks is more. This is due to the fact that there are more chunks along a given dimension to be aggregated in these cases. This effect will be more pronounced with higher data densities.

Figure 23 shows the time for full and partial cube build for Dataset III and for partial cube building for Dataset IV. Partial build times are much smaller than the



(a)

(b)

FIG. 23. Comparing full and partial cubes (a) 10D, 10 million tuples processors, and (b) partial cubes for 20D on 8 and 16 processors.

Operation	Chunksize I	Chunksize II
$ABC \rightarrow AB$	7.446	2.067
$ABC \rightarrow AC$	37.41	12.40
$ABC \rightarrow BC$	103.26	47.71
$AC \rightarrow A$	0.148	0.084
$BC \rightarrow B$	0.194	0.045
$BC \rightarrow C$	0.042	0.032
Total	143.78	62.34

(a)

	Processors		
	4	8	16
Chunksize I	143.78	44.95	30.43
Chunksize II	62.34	11.59	9.65

(b)

FIG. 24. 3D dataset with Chunksize I (32, 8, 16) and Chunksize II (64, 8, 32). (a) Aggregation time (in seconds) on four processors. (b) Speedup performance.

full cube building time. This meant that for the mining of two-way association rules the partial cube can be constructed fairly quickly. Similarly, three-way associations can also be mined efficiently. These two are the more important ones in terms of using support and confidence measures to find useful rules from large databases.

Figure 24 shows the time for calculating various subcubes for two different chunk sizes. A larger chunk size means fewer chunks and hence less chunk traversal and aggregation overhead. Hence we observe lower numbers, indicating better performance. Choice of a chunk size depends on the number of dimensions, size of the dimensions, and density of the data set. A chunk size should be maximized given these variables.

Table 7 shows the parameters that are indicators of the computation communication and I/O parameters and performance for Dataset I on four processors. With a larger chunk size we see that all these reduce, notably the number of split chunks, and the BESS merge read and write operations. These are reflected in the improved times seen in Fig. 24.

Finally, Fig. 25 shows the performance of 100 random range queries on the base cube at level n and the $n - 1$ dimension aggregates for $n = 5$ and $n = 10$. We see that the performance scales very well as the number of processors is increased. Ranges are defined on both the distributed and the non-distributed dimensions. Queries on

TABLE 7

Count of Each Parameter for 3D Dataset with Chunk Size I (32, 8, 16) and Chunk size II (64, 8, 32) on four Processors

Parameter	Chunk size I	Chunk size II
Split chunks	16,928	4,384
BESS communicated	708,489	642,724
Offsets communicated	45,014	19,174
Minichunks write	68,944	58,087
Minichunks read	309,292	221,544
Merge BESS write	99,356	45,915
Merge BESS read	102,521	47,179

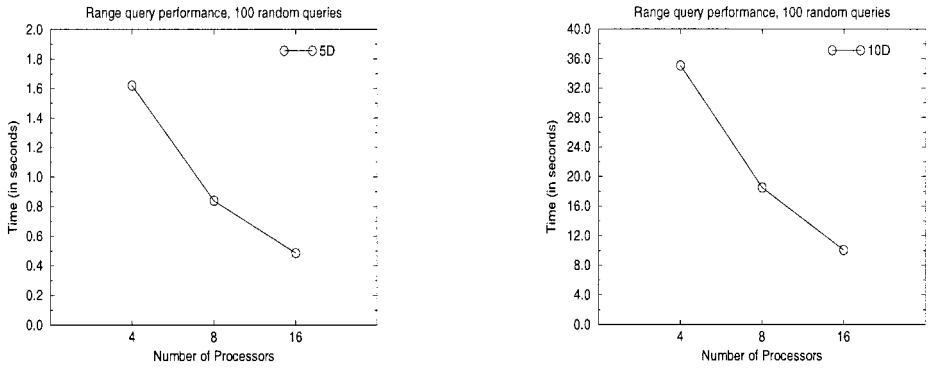


FIG. 25. Range query performance for 100 random queries on the last two levels of the cube lattice for Datasets II (5D) and III(10D).

hierarchies, roll-ups, also scale well with an increasing number of processors as they are also range queries on the lowest level. These queries reflect the queries in the OLAP council benchmarks where variable ranges are specified for various attributes.

8. CONCLUSIONS

In this article we have presented the design and implementation of a scalable parallel system for multidimensional analysis, OLAP, and data mining. Using the multidimensional data model, data are stored in *chunks* which can be either sparse or dense. Sparse chunks are represented by a bit encoding which can be used for efficient aggregation operations on compressed data. For maximum efficiency of operations dense regions can also be stored as multidimensional arrays if the cardinalities of the dimensions involved are not large and the cube size is below a specific threshold. Operations between chunked and multidimensional array cubes are supported.

The data structures to track the different cubes in a data cube, the chunk structures of each cube, and the chunks themselves (using minichunks) use paging to support a large number of cubes, a large number of chunks per cube, and a large chunk size. A combination of these results in supporting a large number of dimensions and large data sizes. For coarse grained parallelism, the cubes can be partitioned using either one largest dimension or a combination of the two largest dimensions. The choice depends on the cardinality of the two largest dimensions and the number of processors. Furthermore cube can be allocated on a single processor if its size is small and the overhead of parallelization is high. Operations between these distributions are supported in our framework.

This framework has been demonstrated for use in OLAP queries which are ad hoc in nature and require fast computation times by preaggregating calculations. Data mining uses some of the precomputed aggregated calculations to compute the probabilities needed for calculating support and confidence measures for association rules and the split point evaluation while building classification trees. Parallelism has been used to support a large number of dimensions and large data sets for effective data analysis and decision making.

REFERENCES

1. S. Agarwal, R. Agrawal, P. M. Deshpande, A. Gupta, J. F. Naughton, R. Ramakrishnan, and S. Sarawagi, On the computation of multidimensional aggregates, in "Proc. 22th Int'l Conference on Very Large Databases, Mumbai, India September 1996."
2. I. Bhandari, "Attribute Focusing: Data Mining for the Layman," Technical Report RC 20136, IBM T. J. Watson Research Center, 1995.
3. I. Bhandari *et al.*, A case study of software process improvement during development, *IEEE Trans. Software Engineering* **19**, 12 (1993).
4. L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone, "Classification and Regression Trees," Wadsworth, Belmont, 1984.
5. P. K. Chan and S. J. Stolfo, Meta-learning for multistrategy and parallel learning, in "Proc. Intl. Workshop on Multistrategy Learning," 1993.
6. G. Colliat, OLAP, relational, and multi-dimensional database systems, in "SIGMOD Record," Vol. 25, 3 September 1996.
7. P. M. Deshpande, K. Ramaswamy, A. Shukla, and J. F. Naughton, Caching multidimensional queries using chunks, in "Proc. ACM SIGMOD Conference on Management of Data," June 1998.
8. U. M. Fayyad, G. Piatetsky-Shapiro, P. Smyth, and R. Uthurusamy, "Advances in Data Mining and Knowledge Discovery," MIT Press, Cambridge, MA, 1994.
9. S. Goil, "High Performance On-Line Analytical Processing and Data Mining on Parallel Computers," Ph.D. thesis, Northwestern University, 1999.
10. S. Goil and A. Choudhary, "BESS: Sparse Data Storage of Multi-dimensional Data for OLAP and Data Mining," Technical Report CPDC-9801-005 Northwestern University, December 1997.
11. S. Goil and A. Choudhary, High performance OLAP and data mining on parallel computers, *J. Data Mining Knowledge Discovery* **1**, 4 (1997).
12. S. Goil and A. Choudhary, High performance data mining using data cubes on parallel computers, in "Proc. International Parallel Processing Symposium," March 1998.
13. S. Goil and A. Choudhary, High performance multidimensional analysis and data mining, in "Proc. SC98: High Performance Networking and Computing Conference," November 1998.
14. S. Goil and A. Choudhary, Efficient parallel classification using dimensional aggregates, in "Proc. ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD 99), Workshop on Large-Scale Parallel KDD Systems," August 1999.
15. D. E. Goldberg, "Genetic Algorithms in Search, Optimization and Machine Learning," Morgan Kaufmann, San Marco, CA, 1989.
16. J. Gray, A. Bosworth, A. Layman, and H. Pirahesh, Data cube: A relational aggregation operator generalizing Group-By, Cross-Tab, and Sub-Totals, in "Proc. 12th International Conference on Data Engineering," 1996.
17. J. Han, Y. Cai, and N. Cercone, Data-driven discovery of quantitative rules in relational databases, *IEEE Trans. Knowledge Data Engineering* **5**, 1 (February 1993).
18. J. Han and Y. Fu, Discovery of multiple-level association rules from large databases, in "Proc. of the 21st VLDB Conference, Zurich," 1995.
19. V. Harinarayan, A. Rajaraman, and J. D. Ullman, Implementing data cubes efficiently, in "Proc. SIGMOD International Conference on Management of Data," 1996.
20. M. Joshi, G. Karypis, and V. Kumar, A new scalable and efficient parallel classification algorithm for mining large datasets, in "Proc. International Parallel Processing Symposium," March 1998.
21. Kenan Software, "An Introduction to Multidimensional Database Technology," available at <http://www.kenan.com/acumate/mdd.htm> 1997.
22. W. Lehner and W. Sporer, On the design and implementation of the multidimensional cube-store storage manager, in "15th IEEE Symposium on Mass Storage Systems (MSS'98)," March 1998.

23. M. Mehta, R. Agrawal, and J. Rissanen, SLIQ: A fast scalable classifier for data mining, in "Proc. of the Fifth Int'l Conference on Extending Database Technology, Avignon," March 1996.
24. D. Michie, D. J. Spiegelhalter, and C. C. Taylor, "Machine Learning, Neural and Statistical Classification," Ellis Horwood, Chichester, 1994.
25. OLAP Council, "OLAP Council Benchmark," available at <http://www.olapcouncil.com>, 1997.
26. J. Ross Quinlan, "C4.5: Programs for Machine Learning," Morgan Kaufmann, San Marco, CA, 1993.
27. S. Sarawagi, Indexing OLAP data, *Data Engineering Bull.* **20**, 1 March 1997.
28. P. Scheuermann, J. Shim, and R. Vingralek, Watchman: A data warehouse intelligent cache manager, in "Proc. 22th Int'l Conference on Very Large Databases, Mumbai, India," June 1998.
29. J. C. Shafer, R. Agrawal, and M. Mehta, SPRINT: A stable parallel classifier for data mining, in "Proc. 22th Int'l Conference on Very Large Databases, Mumbai, India," September 1996.
30. A. Shoshani, OLAP and statistical databases: similarities and differences, in "Proc. Principles of Database Systems," 1997.
31. A. Shukla and P. M. Deshpande, J. Naughton and K. Ramaswamy, Storage estimation for multi-dimensional aggregates in the hierarchies, in "Proc. of the 22nd International VLDB Conference," May 1996.
32. Y. Zhao, P. Deshpande, and J. Naughton, An array-based algorithm for simultaneous multidimensional aggregates, in "Proc. ACM-SIGMOD International Conference on Management of Data," pp. 159-170, 1997.

SANJAY GOIL received his M.Sc.(Tech.) in computer science from the Birla Institute of Technology and Science, Pilani, India, in 1990; his M.S. in computer science from Syracuse University in 1995; and his Ph.D. in computer engineering from Northwestern University in 1999. He was a research associate at the Network computing research group at AT&T Bell Laboratories at Murray Hill, New Jersey, from 1991 to 1993. He has been a member of the Performance Technologies Group in Sun Microsystems, Inc. at Sunnyvale, California, since 1999. His research interests include parallel computing, architecture of high performance parallel systems and servers, and performance modeling. He is a member of the ACM.

ALOK N. CHOUDHARY received his B.E. (with honors) from the Birla Institute of Technology and Science, Pilani, India, in 1982; his M.S. from the University of Massachusetts, Amherst, in 1986; and his Ph.D. from the University of Illinois at Urbana-Champaign in electrical and computer engineering, in 1989. He has been an associate professor in the Electrical and Computer Engineering Department at Northwestern University since September 1996. From 1993 to 1996, he was an associate professor in the Electrical and Computer Engineering Department at Syracuse University; and from 1989 to 1993, he was an assistant professor in the same department. He has spent time visiting various industries—including IBM and Intel—as a research scientist. His main research interests are in high-performance computing and communication systems and their applications in many domains, including multimedia systems, information processing, and scientific computing. In particular, his interests lie in the design and evaluation of architectures and software systems (from system software such as runtime systems, compilers, and programming languages to applications), high-performance servers, high-performance databases, and input-output. He has published more than 100 papers in various journals and conference proceedings in the above areas. He has also written a book and several book chapters on these topics. He received the National Science Foundation's Young Investigator Award in 1993 (for 1993-1999). He has received an IEEE Engineering Foundation Award, an IBM Faculty Development Award, and an Intel Research Council Award. His past and present research has been sponsored by DARPA, NSF, NASA, AFOSR, ONR, DOE, Intel, IBM, and TI. He served as conference co-chair for the International Conference on Parallel Processing and is currently chair of the International Workshop on I/O Systems in Parallel and Distributed Systems. He is an editor of the *Journal of Parallel and Distributed Computing* and has served as a guest editor for *Computer* and *IEEE Parallel and Distributed Technology*. He serves (or has served) on the program committee of many international conferences in architectures, parallel computing, multimedia systems, performance evaluation, distributed computing, etc. He serves in the High-Performance Fortran Forum, a forum of academia, industry, and government laboratories working on standardizing programming languages for portable programming on parallel computers. He is a member of the IEEE, the IEEE Computer Society, and the ACM.