

CERIAS Tech Report 2003-08

**(ALMOST) OPTIMAL PARALLEL BLOCK
ACCESS FOR RANGE QUERIES**

by Mikhail J. Atallah and Sunil Prabhakar

Center for Education and Research in
Information Assurance and Security,
Purdue University, West Lafayette, IN 47907

(Almost) Optimal Parallel Block Access for Range Queries [★]

Mikhail J. Atallah¹

CERIAS and Department of Computer Sciences, Purdue University, West Lafayette, IN 47907, U.S.A.

Sunil Prabhakar^{2,*}

Department of Computer Sciences, Purdue University, West Lafayette, IN 47907, U.S.A.

Abstract

Range queries are an important class of queries for several applications. can be achieved by tiling the multi-dimensional data and distributing it among multiple disks or nodes. It has been established that schemes that achieve optimal parallel block access exist only for a few special cases. Though several schemes for the allocation of tiles to disks have been developed, no scheme with non-trivial worst-case bound is known. We establish that any range query on a $2^q \times 2^q$ -block grid of blocks can be performed using $k = 2^t$ disks ($t \leq q$), in at most $\lceil m/k \rceil + O(\log k)$ parallel block accesses. We achieve this result by judiciously distributing the blocks among the k nodes or disks. Experimental data show that the algorithm achieves very close to $\lceil m/k \rceil$ performance (on average less than 0.5 away from $\lceil m/k \rceil$, with a worst-case of 3). Although several declustering schemes for range queries have been developed, prior to our work no additive non-trivial performance bounds were known. Our scheme guarantees performance within a (small) *additive* deviation from $\lceil m/k \rceil$. Subsequent to this work, Bhatia et al. [4] have proved that such a performance bound is essentially optimal for this kind of scheme, and have also extended our results to the case where the number of disks is a product of the form $k_1 * k_2 * \dots * k_t$ where the k_i s need not all be 2.

Key words: Multi-Dimensional Data, Parallel I/O, Declustering, Range Queries

1 Introduction

Range queries are an important class of queries for several application domains including relational databases, spatial databases, visualization, and GIS applications. Given a multidimensional dataset, a range query specifies a range of values for each dimension. The result of the range query is the set of all items in the dataset that have values within the specified range in each dimension. As the size of the dataset grows, the amount of data that needs to be accessed to answer range queries also increases, resulting in poor performance. In order to improve performance, the dataset is typically tiled along each dimension. Therefore in order to process a range query, it is necessary to access only those tiles or blocks that intersect the query, resulting in reduced I/O and improved performance. Even with such tiling, the performance is limited by the disk I/O. To further improve performance, multiple disks or processing nodes can be used to access the blocks in parallel. The blocks of the dataset are distributed among the disks. The key to achieving gains from parallelism is in the allocation of the blocks to the disks. Note that the data could be placed on multiple disks connected to a single processor or stored on parallel nodes, each with local disk space. Similarly, at the primary storage level, access to rectangular ranges from main memory are needed to refresh rectangular windows in GUI displays, and for visualizing three dimensional objects. Main memory typically consists of several memory banks, each of which can retrieve a single data item in one memory cycle. These range accesses can be optimized by accessing the required data in parallel from the data banks. Once again, the key is to place the data on the multiple data banks so as to reduce the overall number of memory cycles needed to access the data. The allocations described in this paper are applicable to each of these problems. We will use the term disk to refer to such data banks, parallel disks, or nodes.

The goal of the allocation is to achieve optimal parallel access for each range query. The design of these allocation schemes has been an active research area, resulting in the development of several allocation schemes [2,6,7,11,15,16,10]. These schemes are developed under the following framework. Due to the relatively high cost of disk accesses, the CPU processing time is ignored. Fur-

* An earlier version of this paper appeared in the Proceedings of the 19th ACM Symposium on Principles of Database Systems (PODS), Dallas, Texas, May 2000

* Corresponding Author: Fax: +1-765 4940739

Email addresses: mja@cerias.purdue.edu (Mikhail J. Atallah), sunil@cs.purdue.edu (Sunil Prabhakar).

¹ Portions of this work were supported by Grant EIA-9903545 from the National Science Foundation, Contract N00014-02-1-0364 from the Office of Naval Research, and by sponsors of the Center for Education and Research in Information Assurance and Security.

² Portions of this work were supported by NSF CAREER grant IIS-9985019.

thermore, since the disk accesses are random, the cost of a single disk access is assumed to be constant. Note that for the case of main memory banks, this model is completely accurate. Thus, given a query, the cost of executing the query is taken to be proportional to the number of disk accesses performed. When the data are accessed from multiple disks in parallel, the cost is proportional to the largest number of accesses performed on a single disk. For a query that intersects m blocks, the access cost using k disks is at least $OPT = \lceil m/k \rceil$. The worst-case performance of an allocation is, expressed as a function of m , the maximum number of parallel accesses over all possible queries that intersect m blocks. The non-existence of general solutions that achieve optimal allocation for all queries has been established [2,9,16]. Schemes that achieve OPT for every query exist in only a small number of special cases. The existence of allocations that achieve OPT for higher dimensions is expected to be at least as restrictive. Next, we describe the most prominent schemes that have previously been developed. In view of the provable impossibility of achieving OPT except in a small number of special cases, the notation “ OPT ” should be viewed as merely a shorthand for $\lceil m/k \rceil$ rather than as an achievable performance bound.

A wide range of declustering schemes have been proposed in the literature for range queries over parallel disks. The techniques can be divided into two broad classes. Data independent schemes generate an allocation of tiles to disks based upon a mapping function [1–3,5–7,10,11,13,15,16]. Data-dependent techniques typically transform the declustering into a graph problem (e.g. mincut) based upon the datasets considered and possibly also the queries [8,14,12].

Prior to our work, there is no guarantee on the performance of a given range query for any of the existing schemes, except the GRS scheme. For two dimensions, the GRS scheme has been shown to guarantee performance within a *multiplicative* factor of 3 from OPT (i.e., $3OPT$). In this paper, we develop an allocation scheme which has guaranteed worst-case performance within an *additive* deviation from OPT : Within $OPT + O(\log k)$ for two dimensions. Our scheme requires that the number of disks available is $k = 2^t$, numbered from 1 to k . To generate the allocation for a dataset that has been divided into $N_1 \times N_2$ blocks along the two dimensions, we first extend the number of blocks in each dimension such that we have 2^q tiles in each dimension, where $q \geq t$. After generating the allocation for this larger grid of blocks, we simply ignore the extra blocks that were added, resulting in the allocation for the $N_1 \times N_2$ dataset. The disk allocation problem can be viewed as that of coloring the $N = 2^{2q}$ blocks by using colors numbered 1 to k , with the interpretation that a block of color i is to be stored in disk i . The number of parallel block accesses for processing a range query is the maximum occurrence of any color in the rectangular region of blocks defined by the range query.

The rest of this paper is organized as follows. Section 2 describes the coloring

(allocation) scheme used. Section 3 discusses some properties of that coloring scheme. Section 3 proves that, for any 2-dimensional range query, if m is the number of blocks for that range query, then the coloring scheme we use can result in no more than $\lceil m/k \rceil + \gamma$ parallel block accesses where $\gamma = O(\log k)$. Finally, Section 4 concludes the paper.

2 The coloring scheme

We partition the $2^q \times 2^q$ grid of blocks into a $2^{q-t} \times 2^{q-t}$ grid of *groups* each of which is itself a $k \times k$ grid of blocks (recall that $k = 2^t$). We next describe the coloring scheme for the blocks in a group (the same coloring scheme is used for all the groups). Row and column numbers in that description are *relative to that group* (not relative to the whole grid). We start with some definitions.

Let j be a column whose blocks have been colored, and let j' be another column whose coloring is to be derived from that of column j . We say that the coloring of column j' is a $k/2$ -swap of the coloring of column j if we first copy the coloring of j into j' and then we “swap” the coloring of the upper half of column j' with the coloring of its lower half. For example, if the colors of the cells of column j are (in row order) $1, 2, \dots, k$ then the colors for column j' would be $(k/2) + 1, \dots, k, 1, \dots, (k/2)$. More generally, a $k/2^i$ -swap of a column’s coloring, for an integer $i \leq t$, is defined as follows: *Partition the column into 2^{i-1} contiguous, non overlapping pieces of size $k/2^{i-1}$ each. Then, for each piece, swap the coloring of the piece’s upper half with the coloring of the piece’s lower half. (We call it a “ $k/2^i$ ” swap because that is the size of each portion being swapped, as a mnemonic.)*

For example, a 1-swap of a column’s coloring consists of interchanging the colors of cells $2\ell - 1$ and 2ℓ , for all $1 \leq \ell \leq k/2$. We are now ready to describe the coloring of a group of $k \times k$ blocks.

- (1) Assign the colors $1, \dots, k$ to the cells of column 1.
Comment. Although we assign the colors in sorted order, in fact any permutation would also work (as will soon become apparent).
- (2) For $\mu = 1, \dots, t$ in turn, do the following: For $j = 1, \dots, 2^{\mu-1}$ in turn, assign to column $2^{\mu-1} + j$ a coloring that is a $k/2^\mu$ -swap of the coloring of column j .

Figure 1 gives an example of the above coloring for the case $t = 4$ (i.e., 16 colors). All columns are generated from column 1. For example, column 15 is a 1-swap of column 7, which is a 2-swap of column 3, which is a 4-swap of column 1.

1	1 2	1 2 3 4	1 2 3 4 5 6 7 8	1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16
1	1 9	1 9 5 13	1 9 5 13 3 11 7 15	1 9 5 13 3 11 7 15 2 10 6 14 4 12 8 16
2	2 10	2 10 6 14	2 10 6 14 4 12 8 16	2 10 6 14 4 12 8 16 1 9 5 13 3 11 7 15
3	3 11	3 11 7 15	3 11 7 15 1 9 5 13	3 11 7 15 1 9 5 13 4 12 8 16 2 10 6 14
4	4 12	4 12 8 16	4 12 8 16 2 10 6 14	4 12 8 16 2 10 6 14 3 11 7 15 1 9 5 13
5	5 13	5 13 1 9	5 13 1 9 7 15 3 11	5 13 1 9 7 15 3 11 6 14 2 10 8 16 4 12
6	6 14	6 14 2 10	6 14 2 10 8 16 4 12	6 14 2 10 8 16 4 12 5 13 1 9 7 15 3 11
7	7 15	7 15 3 11	7 15 3 11 5 13 1 9	7 15 3 11 5 13 1 9 8 16 4 12 6 14 2 10
8	8 16	8 16 4 12	8 16 4 12 6 14 2 10	8 16 4 12 6 14 2 10 7 15 3 11 5 13 1 9
9	9 1	9 1 13 5	9 1 13 5 11 3 15 7	9 1 13 5 11 3 15 7 10 2 14 6 12 4 16 8
10	10 2	10 2 14 6	10 2 14 6 12 4 16 8	10 2 14 6 12 4 16 8 9 1 13 5 11 3 15 7
11	11 3	11 3 15 7	11 3 15 7 9 1 13 5	11 3 15 7 9 1 13 5 12 4 16 8 10 2 14 6
12	12 4	12 4 16 8	12 4 16 8 10 2 14 6	12 4 16 8 10 2 14 6 11 3 15 7 9 1 13 5
13	13 5	13 5 9 1	13 5 9 1 15 7 11 3	13 5 9 1 15 7 11 3 14 6 10 2 16 8 12 4
14	14 6	14 6 10 2	14 6 10 2 16 8 12 4	14 6 10 2 16 8 12 4 13 5 9 1 15 7 11 3
15	15 7	15 7 11 3	15 7 11 3 13 5 9 1	15 7 11 3 13 5 9 1 16 8 12 4 14 6 10 2
16	16 8	16 8 12 4	16 8 12 4 14 6 10 2	16 8 12 4 14 6 10 2 15 7 11 3 13 5 9 1

8-Swap
4-Swap
2-Swap
1-Swap

Fig. 1. An example of the allocation scheme for 16 disks

3 Properties of coloring scheme

A coloring of a group is said to be *left-to-right legal* if it is obtained according to the process described in the previous section except that the process can be initiated with the first column holding *any* permutation of the k colors (not necessarily the sorted one we used in the previous section). The coloring is *right-to-left legal* if we do the same thing except that we start with the rightmost column of the group and proceed leftward from there. The notions of top-to-bottom legal and of bottom-to-top legal are defined using a similar coloring process that operates by rows rather than by columns.

We begin with the **group properties** that hold within each group:

- (1) Any of the following four properties of a group coloring implies the other three: { left-to-right legal, right-to-left legal, top-to-bottom legal, bottom-to-top legal }. Therefore the left-to-right legal coloring we produced for a block also has the other three properties. We henceforth use the word *legal coloring* as an abbreviation for these.
- (2) Each column of a group contains a permutation of the k colors.
- (3) Each row of a group contains a permutation of the k colors.

Group property 2 is an immediate consequence of the way a column is colored (because copying then permuting the coloring of a column results in another permutation of the colors). Group property 1 will be proved (together with other properties) at the end of this section. Group property 3 follows from group properties 1 and 2.

We now define a finer partition of the input grid than its partition into groups. This is not needed algorithmically and is done purely for the sake of the analysis. To avoid unnecessarily cluttering the analysis with “[.]” notation,

we assume t is even (it is easy to modify the analysis for odd t).

Partition each $2^t \times 2^t$ group into a $2^{t/2} \times 2^{t/2}$ grid of *superblocks* each of which is itself a $2^{t/2} \times 2^{t/2}$ grid of blocks (observe that $2^{t/2} = \sqrt{k}$). A range query is said to *vertically span* a superblock if it does not completely contain that superblock, and its intersection with that superblock is a contiguous set of columns of that superblock (horizontal span is defined similarly with respect to rows). Let S be a set of superblocks that are vertically contiguous to each other (i.e., each of them is “on top” of another one of them). A range query is said to vertically span S if it is contained in S and it vertically spans all the superblocks of S *at corresponding sets of columns*, i.e., if its intersection with a superblock x of S is the same interval of columns $[j, j']$ for all such x (horizontal span is analogously defined).

The following **superblock properties** hold:

- (1) Each superblock of $\sqrt{k} \times \sqrt{k}$ blocks contains the k colors (i.e., one occurrence of each color).
- (2) Let S be a set of superblocks that are vertically contiguous to each other. For any range query that vertically spans S , the legal coloring described in the previous section is within an additive $2^{-1} \log k$ of OPT for that query (i.e., results in at most $\lceil m/k \rceil + 2^{-1} \log k$ parallel block accesses where m is the number of blocks touched by the query).
- (3) Let S be a set of superblocks that are horizontally contiguous to each other. For any range query that horizontally spans S , the legal coloring described in the previous section is within an additive $2^{-1} \log k$ of OPT for that query (i.e., results in $\lceil m/k \rceil + 2^{-1} \log k$ parallel block accesses).

The above superblock properties and group property 1 are proved below.

Proof of group property 1 and superblock properties 2 and 3

We give the proof for the general case where the coloring process is initiated with an arbitrary permutation of k distinct symbols (rather than the particular sorted permutation of the integers 1 to k we used in Section 2).

The proof is by induction on t . The basis, $t = 1$, is trivial. We assume inductively that the properties hold for t . To show that they hold for $t + 1$, we observe that the coloring of a $2^{t+1} \times 2^{t+1}$ grid can be thought of as consisting of the following four steps (which we describe assuming a coloring that starts with an initial column and operates left-to-right – essentially the same argument can be made for a coloring process that starts with an initial row and operates row-wise).

- (1) (*Coalesce step*) For the initial column (of size 2^{t+1}), *coalesce* entry $2\ell - 1$ and entry 2ℓ , $1 \leq \ell \leq 2^t$. This “shrinks” the column into one of half the size ($= 2^t$), with 2^t distinct *new* colors; each new color c corresponds to an ordered pair (c', c'') of old colors.
- (2) (*Induction step*) Perform the iterative coloring process on a $2^t \times 2^t$ array A with the (shrunk) column of size 2^t as the initial column (and using the new colors). This results in a $2^t \times 2^t$ colored array A that (by the induction hypothesis) has the desired properties (relative to the new colors).
- (3) (*Duplication step*) Duplicate the colored $2^t \times 2^t$ array A and, in the duplicate copy \hat{A} , replace every new color $c = (c', c'')$ by its *complement* $\hat{c} = (c'', c')$ (i.e., complementing c consists of interchanging the ordering of the two old colors c' and c'' that define it). Array \hat{A} has the desired properties (relative to the complemented colors).
- (4) (*Expansion step*) Append \hat{A} to the right of A , resulting in a $2^t \times 2^{t+1}$ array. This array is turned into a $2^{t+1} \times 2^{t+1}$ one by “expanding” each color c to the two old colors corresponding to it (thus doubling the size of each column).

We must show that the last (“expansion”) step results in an array that satisfies the claimed properties. We do so separately for each property.

Group property 1: Because we assumed a coloring that is left-to-right legal, we must show that the final array is also right-to-left legal, top-to-bottom legal, and bottom-to-top legal (the proof would be very similar if we had assumed one of the other three legal colorings rather than left-to-right, so we avoid repeating this argument four times).

That the left-to-right legal array is right-to-left legal can be seen by looking at the resulting rightmost colored column (after the expansion step) and trying to use it as the starting column for a right-to-left legal coloring: In this “right-to-left” process, we would first generate the expanded version of \hat{A} because the unexpanded \hat{A} itself has group property 1 (by the induction hypothesis). Next, the last step of the right-to-left process would duplicate the expanded \hat{A} and append a 1-swapped copy of it to the left of the expanded \hat{A} : But a 1-swapped version of the expanded \hat{A} is the same as the expanded version of A . Thus the right-to-left process would generate exactly the same array.

We now prove that the left-to-right array is bottom-to-top legal. Starting with the bottom row (after the expansion step), it is easy to see that the next-to-bottom row looks just like a copy of the bottom row but with the left and right halves interchanged. From that point on, the bottom-to-top process does not cause any interaction between the left half of the rows and their right half, and can be viewed as two separate bottom-up processes (one for each half). That the left-half process gives rise to the expanded version of A follows from the fact that A itself is bottom-to-top legal (by the induction hypothesis). Similarly, the right-half process gives rise to the expanded version of \hat{A} because \hat{A} is bottom-to-top legal.

To prove that the left-to-right array is top-to-bottom legal, we use the fact (just proved) that it is bottom-to-top legal, followed by an almost identical argument to the one we used for showing that left-to-right legal implies right-to-left legal (except that the roles of rows and columns are interchanged, “bottom” replaces “left” and “top” replaces “right”).

This completes the proof of group property 1.

Superblock property 3: At the bottom of the recursive construction the additive deviation from optimal is zero (verify this – in fact the first deviation from optimality by an additive 1 occurs for $k = 64$). Because the superblocks in S can extend partly over A and partly over \hat{A} , each “expansion” step in the construction introduces one more additive unit deviation from optimality. To see that this is so, before expansion let m_1 (resp., m_2) be the number of cells of S in A (resp., \hat{A}), γ be the total number of parallel block accesses we use for those m_1 cells in A (independently of \hat{A}) and m_2 cells in \hat{A} (independently of A), and let δ be the deviation of γ from the sum of the two individual OPT values of the two pieces of S (that is, the deviation from $\lceil 2m_1/k \rceil + \lceil 2m_2/k \rceil$ where we used the fact that the effective number of colors on each side before expansion is $k/2$). When we “expand” we effectively double m_1 , m_2 , and the number of colors (which becomes k). We can still process S , after expansion, with γ parallel block accesses, but the deviation of that γ from the new (combined) OPT value can increase by 1 unit because we could have $\lceil (2m_1 + 2m_2)/k \rceil = \lceil 2m_1/k \rceil + \lceil 2m_2/k \rceil - 1$. That is, the number of parallel block accesses needed has stayed same but the overall (combined) OPT has gone down by 1 compared to the sum of the two individual OPT values (the OPT for the portion of S in A + the OPT for the portion of S in \hat{A}). Because a superblock has dimensions $\sqrt{k} \times \sqrt{k}$, the total deviation is $\log \sqrt{k} = 2^{-1} \log k$.

Superblock property 2: Immediately follows from group property 1 and superblock property 3 (because it is the “vertical” equivalent of superblock property 3).

Proof of superblock property 1

For any $\alpha < t$, consider a partition of the leftmost column of a group into 2^α pieces of size $2^{t-\alpha}$ each. Call these pieces $1, \dots, 2^\alpha$.

Claim. For any piece i ($1 \leq i \leq 2^\alpha$), the $2^\alpha \times 2^{t-\alpha}$ rectangle R of k blocks whose left side is piece i , contains all k colors (i.e., each color exactly once).

Before proving the above claim, we note that it would automatically imply superblock property 1 for $\sqrt{k} \times \sqrt{k}$ superblocks that are “left-adjusted” in the sense that their left side is on the leftmost column of their group (simply by choosing $\alpha = t$ in the claim). That the same is true for superblocks that are

further to the right within the block, follows from the observation that the left-to-right legal coloring process maintains the same set of colors from one superblock R to the next superblock immediately to the right of R (it merely permutes the colors). Therefore it suffices to prove the above claim.

We prove the claim by induction on α . The basis ($\alpha = 0$) holds because of group property 3. Now, assume inductively that the claim holds for $\alpha - 1$. We partition the piece i into two halves U (“upper”) and L (“lower”): By the induction hypothesis, the $2^{\alpha-1} \times 2^{t-\alpha+1}$ rectangle R_U (respectively, R_L) whose left side is U (respectively, L) satisfies the claim. Now, the colors in the right half of R_U (respectively, R_L) are the same as the colors in the left half of R_L (respectively, R_U) because the last step in the coloring of R consisted of a “swap” that copied the colors of the left half of R_L (respectively, R_U) into the right half of R_U (respectively, R_L). Therefore the set of colors that appear in R is the same as the set of colors that appear in R_U (or R_L), namely the full set of k colors (each color once). This completes the proof of the claim.

Proof of performance bound

If the query is entirely contained in one superblock then our coloring implies a single parallel block access (because no color appears twice in a superblock), which is optimal. We henceforth assume that the query is not entirely contained in a superblock. We distinguish two cases, depending on whether the query completely contains a superblock or not. We begin with the case where it contains one or more superblocks.

If the query does not completely contain any superblock, then it can be decomposed into at most six subqueries: Four that are each completely contained in a superblock, and two each of which spans (either horizontally or vertically) a set S of superblocks that are (vertically or horizontally) contiguous. The four subqueries that are completely contained in superblocks can each be done in one parallel block access. The other two subqueries are each done with a number of block accesses that is within an additive $2^{-1} \log k$ of the *OPT* value for that individual subquery (by superblock properties 2 and 3 of the previous section). Therefore the total number of parallel block accesses for such queries cannot exceed *OPT* by more than $3 + \log k$.

If the query completely contains one or more superblocks, then we can partition it into nine subqueries:

- (1) Four subqueries that are each completely contained in a superblock (these are the four “corners” of the rectangle defining the original query). These subqueries can each be done in one parallel block access.
- (2) One subquery that consists of all the superblocks that are completely

contained in the original query, say, a rectangle of m' superblocks. These can be done in m' parallel block accesses with full disk utilization (i.e., no disk is idle during any of these m parallel steps). This follows from the fact that each color appears exactly once in a superblock.

- (3) Two subqueries each of which vertically spans a set of superblocks that are vertically contiguous (one of them is just below the top-left corner subquery, the other just below the top-right corner subquery). Each such subquery is done with a number of block accesses that is within an additive $2^{-1} \log k$ of optimal for that individual subquery (by superblock property 2 of the previous section).
- (4) Two subqueries each of which horizontally spans a set of superblocks that are horizontally contiguous (one of them is just to the right of the top-left corner subquery, the other just to the right of the bottom-left corner subquery). Each such subquery is done with a number of block accesses that is within an additive $2^{-1} \log k$ of optimal for that individual subquery (by superblock property 3 of the previous section).

The above implies that the number of subqueries that can (each) introduce a deviation of 1 from OPT are (at most) 4 “corner” subqueries, the number of subqueries that can (each) introduce an additive $2^{-1} \log k$ from optimal are 4 subqueries that span sets of superblocks that are contiguous along a dimension. The total possible deviation from OPT is then $4 + 4 * 2^{-1} \log k - 1 = 3 + 2 \log k$ (where we subtracted one because even in an optimal coloring at least one parallel block access is needed for these 8 subqueries).

This completes the proof. □

In practice, the deviation from OPT is much less than $3 + 2 \log k$. In particular, experimentally, we found the deviation to be no more than 3, and the average deviation less than 0.5 from the (unachievable) OPT .

4 Conclusion

Range queries are an important class of queries for several applications. Performance improvements are typically achieved through parallel I/O by tiling the data set and distributing it among multiple disks or processing nodes. Therefore, in order to process a range query, it is necessary to access only those tiles or blocks that intersect with the query. Though several schemes for the allocation of tiles to disks have been developed, prior to our work this worst-case performance was known only for one scheme with two-dimensional data. Moreover, this bound was a *multiplicative* factor of 3 from OPT . In this paper we developed a novel allocation scheme with guaranteed worst-case performance for any number of dimensions. We showed that any range query on a $2^q \times 2^q$ -block grid of blocks can be performed using $k = 2^t$ disks ($t \leq q$),

in at most $OPT + O(\log k)$ parallel block accesses. Experimental data show that the algorithm achieves very close to OPT performance.

References

- [1] K. A. S. Abdel-Ghaffar and A. El Abbadi. Optimal disk allocation for partial match queries. *Transactions of Database Systems*, 18(1):132–156, March 1993.
- [2] K. A. S. Abdel-Ghaffar and A. El Abbadi. Optimal allocation of two-dimensional data. In *Int. Conf. on Database Theory*, pages 409–418, Delphi, Greece, Jan. 1997.
- [3] R. Bhatia, R. K. Sinha, and C.-M. Chen. Declustering using golden ratio sequences. In *Proc. of Int'l. Conference on Data Engineering (ICDE)*, San Diego, California, March 2000.
- [4] R. Bhatia, R. K. Sinha, and C.-M. Chen. Hierarchical declustering schemes for range queries. In *Proc. of Int'l. Conference on Extending Database Technology*, pages 525–537, Konstanz, Germany, March 2000.
- [5] B. Chor, C. E. Leiserson, R. L. Rivest, and J. B. Shearer. An application of number theory to the organization of raster-graphics memory. *Journal of the Association for Computing Machinery*, 33(1):86–104, January 1986.
- [6] H. C. Du and J. S. Sobolewski. Disk allocation for cartesian product files on multiple-disk systems. *ACM Transactions of Database Systems*, 7(1):82–101, March 1982.
- [7] C. Faloutsos and P. Bhagwat. Declustering using fractals. In *Proc. of the 2nd Int. Conf. on Parallel and Distributed Information Systems*, pages 18 – 25, San Diego, CA, Jan 1993.
- [8] S. Ghandeharizadeh and D. J. DeWitt. A multiuser performance analysis of alternative declustering strategies. In *Proceedings of the International Conference on Data Engineering (ICDE)*, pages 466–475, Los Angeles, California., February 1990.
- [9] B. Himatsingka and J. Srivastava. Performance evaluation of grid based multi-attribute record declustering methods. In *Proceedings of the International Conference on Data Engineering (ICDE)*, pages 356–365, Houston, Texas, February 1994.
- [10] K. Kim and V. K. Prasanna. Latin squares for parallel array access. *IEEE Transactions on Parallel and Distributed Systems*, 4:4, 1993.
- [11] M. H. Kim and S. Pramanik. Optimal file distribution for partial match retrieval. In *Proc. ACM SIGMOD Int. Conf. on Management of Data*, pages 173–182, Chicago, 1988.

- [12] D.-R. Liu and S. Shekhar. A similarity graph-base approach to declustering problems and its application toward parallelizing grid files. In *Proceedings of the 11th International Conference on Data Engineering*, pages 373–381, Taipei, Taiwan, 1995.
- [13] Y.-L. Lo, K. A. Hua, and H. C. Young. A general multidimensional data allocation method for multicomputer database systems. In *8th Int. Conf. on Database and Expert Systems Applications*, pages 357–66, Toulouse, France, September 1997.
- [14] B. Moon, A. Acharya, and J. Saltz. Study of scalable declustering algorithms for parallel grid files. In *Tenth International Parallel Processing Symposium*, pages 434–440, 1996.
- [15] S. Prabhakar, K. Abdel-Ghaffar, D. Agrawal, and A. El Abbadi. Cyclic allocation of two-dimensional data. In *Proc. of the International Conference on Data Engineering (ICDE'98)*, pages 94–101, Orlando, Florida, Feb 1998.
- [16] Y. Zhou, S. Shekhar, and M. Coyle. Disk allocation methods for parallelizing grid files. In *Proceedings of the International Conference on Data Engineering (ICDE)*, pages 243–252, February 1994.