

1999

## **(Almost) Optimal Parallel Block Access for Range Queries**

Mikhail J. Atallah  
*Purdue University, mja@cs.purdue.edu*

Sunil Prabhakar  
*Purdue University, sunil@cs.purdue.edu*

**Report Number:**  
99-020

---

Atallah, Mikhail J. and Prabhakar, Sunil, "(Almost) Optimal Parallel Block Access for Range Queries" (1999). *Department of Computer Science Technical Reports*. Paper 1451.  
<https://docs.lib.purdue.edu/cstech/1451>

This document has been made available through Purdue e-Pubs, a service of the Purdue University Libraries. Please contact [epubs@purdue.edu](mailto:epubs@purdue.edu) for additional information.

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

**Mikhail J. Atallah  
Sunil Prabhakar**

**Department of Computer Sciences]  
Purdue University  
West Lafayette, IN 47907**

**CSD TR #99-020  
June 1999**

# (Almost) Optimal Parallel Block Access for Range Queries

Mikhail J. Atallah\*

CERIAS and

Department of Computer Sciences

Purdue University

West Lafayette, IN 47907

U.S.A.

mja@cs.purdue.edu

Sunil Prabhakar

Department of Computer Sciences

Purdue University

West Lafayette, IN 47907

U.S.A.

sunil@cs.purdue.edu

## Abstract

Range queries are an important class of queries for several applications including relational databases, spatial databases, and GIS applications. For large datasets, the performance of range queries is limited by disk I/O. Performance improvements are achieved by tiling the multi-dimensional data and distributing it among multiple disks or nodes. Consequently, in order to process a range query, it is necessary to access only those tiles or blocks that intersect with the query. Given  $k$  disks, a query that accesses  $m$  blocks, the optimal number of parallel block accesses that is theoretically possible is  $OPT = \lceil m/k \rceil$ . Though several schemes for the allocation of tiles to disks have been developed, no scheme with guaranteed worst-case performance 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  $OPT + O(1)$  parallel block accesses. The result generalizes to higher dimensions:  $OPT + f(d)$  parallel block accesses for  $d$ -dimensional queries. 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  $OPT$  performance (on average less than 0.5 away from  $OPT$ , with a worst-case of 3). Although several declustering schemes for multidimensional range queries have been developed, this is the first scheme with a guaranteed non-trivial performance bound.

## 1 Introduction

Range queries are an important class of queries for several application domains including relational databases, spatial databases, 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,

---

\*Portions of this work were supported by sponsors of the CERIAS Laboratory.

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. We will refer to each such disk or node as a disk.

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 [DS82, KP88, FB93, AE97, PAAE98]. These schemes are developed under the following framework. Due to the relatively high cost of disk accesses, the CPU processing time is ignored. Furthermore, since the disk accesses are random, the cost of a single disk access is assumed to be constant. 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 *optimal* or lowest access cost using  $k$  disks is  $\lceil m/k \rceil$ . An allocation of blocks to disks is said to be *strictly optimal* if it is optimal for every possible range query. It has been established that for the 2-dimensional case, strictly optimal allocations exist in only a small number of cases [AE97]. In particular, there are strictly optimal allocations if and only if: 1) the number of disks is 1,2,3 or 5; or 2) there are no more than two blocks in at least one of the dimensions; or 3) the number of disks is almost as large as the total number of blocks; and 4) a special case for a  $4 \times 4$  tiling with 8 disks. The existence of strictly optimal allocations for higher dimensions is expected to be at least as restrictive. The paper [AE97], also describes a scheme that produce strictly optimal allocation for two dimensional data whenever one exists. Next, we briefly describe the most important allocation schemes that have previously been developed.

For ease of exposition, the following notation is used. For a  $d$ -dimensional dataset, each block is described by a set of coordinates  $(x_0, x_1, \dots, x_{d-1})$ . Each coordinate,  $x_j$ , is in the range  $[0, N_j - 1]$  and represents the order of the block in dimension  $j$ , where dimension  $j$  is divided into  $N_j$  blocks. The number of disks is  $k$ , and each disk is identified by a number ranging from 0 to  $k - 1$ . The Disk Modulo (DM) scheme [DS82], developed by Du and Sobolewski and later extended for range queries and dynamic files in [LSR92] allocates block  $(x_0, \dots, x_{d-1})$  to disk  $(x_0 + x_1 + \dots + x_{d-1}) \bmod k$ . The Fieldwise eXclusive (FX) method proposed by Kim and Pramanik [KP88], allocates a block to the disk given by the lowest  $\log_2 k$  bits of the bit-wise exclusive-OR of the binary representations of all the coordinates of the block. The Hilbert Curve Allocation Method (HCAM) [FB93], proposed by Faloutsos and Bhagwat, is based upon the Hilbert space-filling curve. Hilbert curves can be used to convert a discrete multidimensional space into a linear sequence such

that spatial proximity is preserved as much as possible. After mapping the blocks into this linear sequence, the blocks are assigned to disks in a round-robin fashion. The allocation method that achieves strictly optimal allocation for 2-dimensional data, described in [AE97], allocates block  $(x_0, x_1)$  to disk  $(x_0 + \lfloor \frac{M}{2} \rfloor x_1) \bmod k$ . A class of declustering schemes called Cyclic allocation schemes was developed in [PAAE98] and [PAE98b]. These schemes allocate block  $(x_0, \dots, x_{d-1})$  to disk  $(x_0 H_0 + x_1 H_1 + \dots + x_{d-1} H_{d-1}) \bmod k$ , where the values  $H_0, H_1, \dots, H_{d-1}$ , are called *skip* values. Each scheme in the class is defined by a different choice of these values. It is shown that the choice of skip values is critical in determining the quality of the allocation, and three different approaches for determining values that give good performance are also developed. There has also been some recent work on declustering for similarity queries, also known as nearest-neighbor queries [BBB<sup>+</sup>97, PAE98a].

The relative performance of these schemes has been studied experimentally in earlier work [PAAE98, PAE98b]. It is seen that, on the average, the Cyclic schemes outperform the other schemes. However there is no guarantee on the performance of a given range query for any of the existing schemes. The worst-case bound for allocation schemes remains an open question. In this paper, we develop an allocation scheme which has guaranteed worst-case performance. We begin with the 2-dimensional case and later generalize to higher dimensions.

Our scheme requires that the number of disks available is  $k = 2^t$ . They are 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 4 proves that, for any 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  $\gamma \leq 7$ . In practice the 7 is a considerable overestimate for  $\gamma$ : The experimental data of Section 6 reveals that the  $\gamma$  is typically no larger than 3. Section 5 describes a generalization of our scheme to higher dimensions. Finally, Section 7 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, so the name acts as a mnemonic.)

For example, a 1-swap of a column’s coloring consists of interchanging the colors of cells  $2\ell - 1$  and  $2\ell$ , 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  $t = 2$  (i.e., 16 colors). All columns are generated from column 1. For example, column 16 is a 1-swap of column 8, which is a 2-swap of column 4, which is a 4-swap of column 2, which is a 8-swap of column 1. Similarly, 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.

### 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**, i.e., the properties that hold within each group:

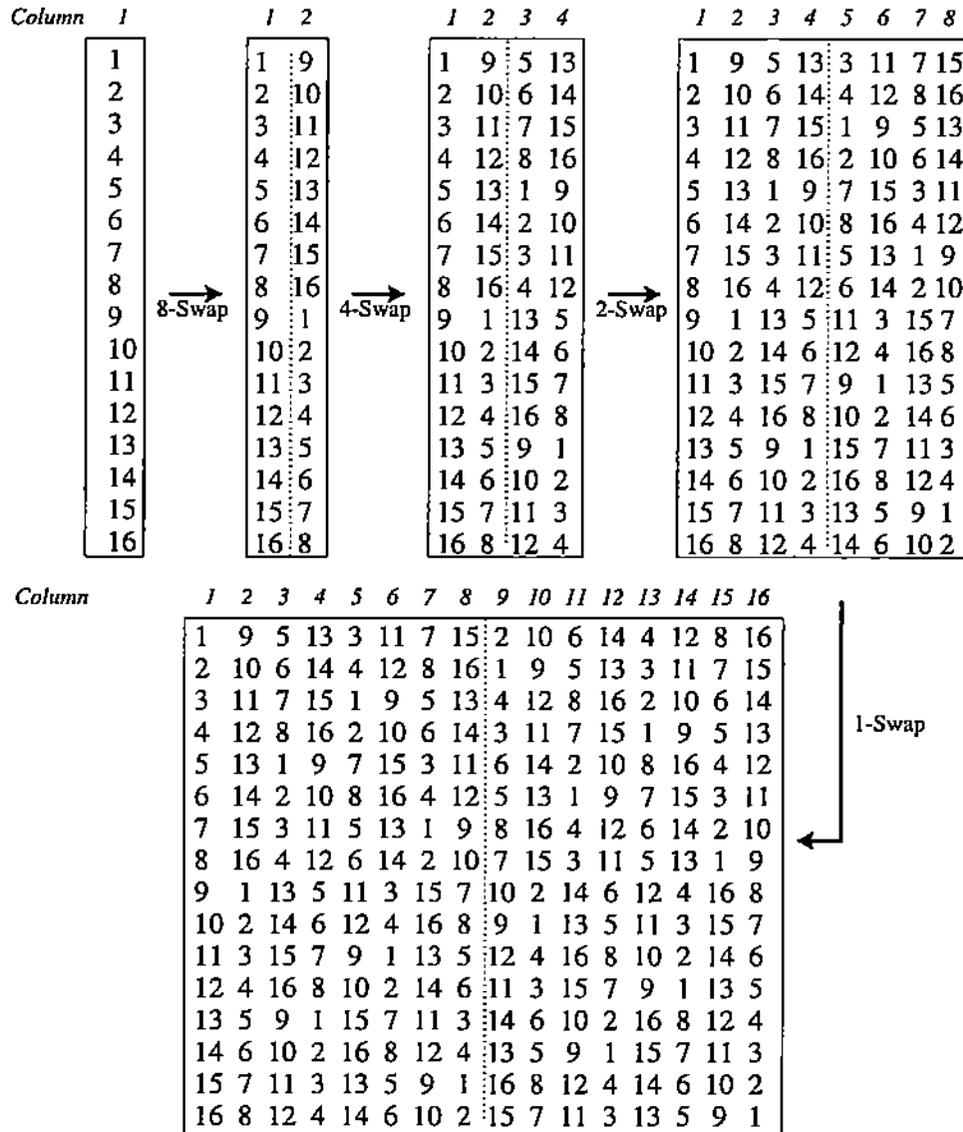


Figure 1: An example of the allocation scheme for 16 disks

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 optimal for that query (i.e., results in  $\lceil m/k \rceil$  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 optimal for that query (i.e., results in  $\lceil m/k \rceil$  parallel block accesses where  $m$  is the number of blocks touched by the query).

The above superblock properties are proved below (together with group property 1).

### 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).

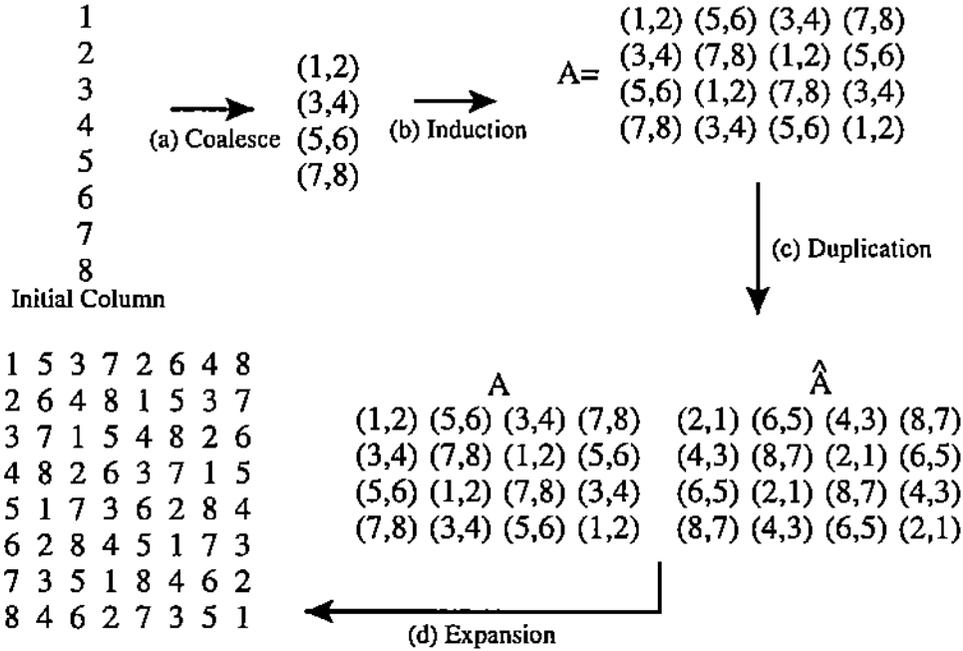


Figure 2: An example of the inductive step for  $t = 2$

1. (*Coalesce step*) For the initial column (of size  $2^{t+1}$ ), *coalesce* entry  $2\ell - 1$  and entry  $2\ell$ ,  $1 \leq \ell \leq 2^t$ . 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, of course).
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).

Figure 2 gives an example of the above four steps for  $t = 2$ .

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

*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 right-most 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 version of  $\hat{A}$  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 2:* The superblocks in  $S$  are either all in the expanded version of  $A$ , or all in the expanded version of  $\hat{A}$ . In either case, the local optimality follows from the local optimality of  $A$  or (respectively)  $\hat{A}$ : If  $m'$  is the relevant number of cells of  $A$  (respectively,  $\hat{A}$ ) then in the expanded version the relevant number of cells is double ( $= 2m'$ ) and so is the number of colors (hence the ratio is the same as before expansion, i.e., optimal).

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

## Proof of superblock property 1

For any  $\alpha < t$ , consider a partition of the leftmost column of a group into  $2^\alpha$  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  $\alpha$ . 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.

## 4 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 four subqueries: Two 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 two subqueries that are completely contained in subblocks can each be done in one parallel block access. The other two subqueries are each done with a number of block accesses that is optimal *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.

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 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 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 and 4 subqueries that span sets of superblocks that are contiguous along a dimension. The total possible deviation from  $OPT$  is then  $4 + 4 - 1 = 7$  (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 7, as becomes apparent in the Section 6.

## 5 Higher Dimensions

For  $d$ -dimensional range queries, we assume a  $2^q \times 2^q \times \dots \times 2^q$  grid of  $2^{dq}$  blocks. We assume  $k = 2^{(d-1)t}$  for some even integer  $t$ .

We partition the grid into groups, where each group is a  $d$ -dimensional square of volume  $k^{d/(d-1)}$ , i.e., it is a  $k^{1/(d-1)} \times \dots \times k^{1/(d-1)}$  grid. Note that a group consists of  $k^{1/(d-1)}$  copies of a  $(d-1)$ -dimensional square of  $k$  blocks. We next describe the coloring scheme for the blocks in a group (the same coloring scheme is used for all the groups). We start with some definitions.

Let  $A$  be a  $(d-1)$ -dimensional square of  $k$  blocks whose blocks have been colored, and let  $A'$  be another similarly shaped square whose coloring is to be derived from that of  $A$ . We say that the coloring of  $A'$  is a  $k^{1/(d-1)}/2$ -swap of the coloring of  $A$  if we first copy the coloring of  $A$  into  $A'$  and then we do the following: We partition  $A'$  into  $2^{d-1}$   $(d-1)$ -dimensional sub-squares of size  $k/2^{d-1}$  each, and for each of pair of diagonal sub-squares we “swap” the colorings of the pair. For example, if  $d = 3$  then  $A$  and  $A'$  are two-dimensional squares, and the coloring of  $A'$  is obtained by first copying the coloring of  $A$  into it and then interchanging the colorings of the top-right and bottom-left quadrants, and of its top-left and bottom-right quadrants. For example, if  $A$  is  $2 \times 2$  (hence  $k = 4$ ) and the coloring of  $A$  is 1, 2 for row 1 and 3, 4 for row 2, then the coloring of  $A'$  is 4, 3 for row 1 and 2, 1 for row 2.

More generally, a  $k^{d/(d-1)}/2^i$ -swap of the coloring of  $A$ , for an integer  $i \leq t$ , is defined by partitioning  $A$  into  $2^{d-i-1}$   $(d-1)$ -dimensional squares of size  $k/2^{d-i-1}$  each and then, *within each square*, doing what we did above (that is, partitioning each square into  $2^{d-1}$  sub-squares and swapping the colorings of the pairs of diagonal sub-squares within each square).

We are now ready to describe the coloring of a group of  $k^{1/(d-1)} \times \dots \times k^{1/(d-1)}$  blocks. We view the group as consisting of the “stacking” on top of each other, along the “height” dimension, of  $k^{1/(d-1)}$   $(d-1)$ -dimensional squares of  $k$  blocks each (each of which is therefore perpendicular to the height dimension).

1. (Initial coloring)

Assign the colors  $1, \dots, k$  (in any order) to the  $k$  cells of the “bottom”  $(d-1)$ -dimensional square of  $k$  blocks.

2. (“Spreading” the coloring along the height dimension)

For  $\mu = 1, \dots, t$  in turn, do the following: For  $j = 1, \dots, 2^{\mu-1}$  in turn, assign to the  $(d-1)$ -dimensional square of height  $2^{\mu-1} + j$  a coloring that is a  $k^{1/(d-1)}/2^\mu$ -swap of the coloring of the  $(d-1)$ -dimensional square of height  $j$ .

Proofs similar to the ones given in the previous section for the case  $d = 2$ , give the following.

1. The coloring obtained above is such that, for *any*  $(d-1)$ -dimensional square of  $k$  blocks (i.e., even one that is not perpendicular to the height dimension), a group contains the  $k$  distinct colors.
2. Suppose we are given any coloring obtained in the way we described. Then by starting with *any* already colored  $(d-1)$ -dimensional square of  $k$  blocks as the initial one (i.e., not necessarily a square perpendicular to the height dimension), and “spreading” the coloring along the remaining dimension, we still obtain the same coloring as before.

The proofs are very similar to the ones given for the case  $d = 2$ , and are therefore omitted. The notion of a superblock also extends naturally: It is now a  $[k^{1/d}] \times \dots \times [k^{1/d}]$  of dimensionality  $d$  that contains  $k' \leq k$  distinct colors. To prove the performance bound, we use a query-decomposition technique similar to the one we used for the 2-dimensional case. The number of queries that can introduce deviations from  $OPT$  are, as before, the “corner” subqueries and the subqueries that span sets of superblocks that are contiguous along a dimension (each such set is locally optimally colored for that subquery). Each such subquery can introduce a deviation of (at most) one, and the number of such subqueries depends only on  $d$  (because the number of corners and faces of a  $d$ -dimensional hyperrectangle depends only on  $d$ , not on  $m$  or  $k$ ). Therefore the number of parallel block accesses is  $OPT + f(d)$  for some function  $f$ . As mentioned earlier, we have experimentally found  $f(d)$  to be quite small, so that the naive  $2^d$  upper bound on  $f(d)$  seems to be a large overestimate. A mathematical worst-case and average-case analysis of  $f(d)$  is an interesting area of future investigation.

## 6 Experimental results

In this section, we present the performance of our new allocation scheme on sample datasets. The objective of the experiments is to observe by how much the new allocation scheme deviates from the optimal, OPT. We also tested the major existing allocation schemes described in Section 1, viz. Disk Modulo (DM) [DS82], Fieldwise eXclusive (FX) [KP88], Hilbert Curve Allocation Method (HCAM) [FB93], and the Generalized Fibonacci (GFIB) scheme from the Cyclic allocation schemes [PAAE98, PAE98b]. We conducted the experiments for 2 and 3 dimensions. Experiments for higher dimensions were not conducted due to the large amounts of computation required. For 2 dimensions, tests were conducted with 4, 16, and 64 disks. In the case of 3 dimensions, tests were conducted with 8 and 64 disks. For each combination of number of dimensions and disks, several tilings were tested. In each test, we measured the maximum and average deviation of queries from OPT. In all experiments (with one exception), all possible range were considered. In the case of 64 disks with 3 dimensions, evaluating all queries was too time-consuming, therefore we considered a random set of 100,000 queries. This was repeated with three different random sets to gain confidence in the results.

The results for 2 dimensions are shown in Tables 1 through 3. Table 1 gives the results for 4 disks, Table 2 gives the results for 16 disks, and Table 3 shows the results for 64 disks. The values for the newly developed scheme are shown under the “NEW” columns. As can be seen, the maximum value of the deviation for the new scheme is no more than, and the average deviation is less than 0.5 from the optimal OPT. It should be noted that even tilings that are not multiples of  $2^l$  have similar performance for most schemes (only the HCAM scheme is affected significantly).

<i>Tiling</i>	<i>Maximum</i>					<i>Average</i>				
	HCAM	DM	FX	GFIB	NEW	HCAM	DM	FX	GFIB	NEW
$4 \times 4$	1	1	1	1	1	0.22	0.09	0.05	0.09	0.01
$16 \times 16$	6	1	1	1	1	0.637	0.070	0.035	0.070	0.014
$32 \times 32$	12	1	1	1	1	0.998	0.066	0.033	0.066	0.015
$33 \times 29$	14	1	1	1	1	1.049	0.066	0.033	0.066	0.017

Table 1: Results for various tilings of 2-dimensional data with 4 disks

For higher dimensions, we conducted several experiments with 3-dimensional datasets. In each experiment, the tiling was chosen to be  $\sqrt{k} \times \sqrt{k} \times \sqrt{k}$  with  $k$  disks. Table 4 shows the maximum and average deviations for each of the schemes. As can be seen, the maximum observed deviation for the new scheme is 4, and the average is no more than 0.531. In fact, the new scheme consistently outperformed all other schemes in our experiments.

Tiling	Maximum					Average				
	HCAM	DM	FX	GFIB	NEW	HCAM	DM	FX	GFIB	NEW
$4 \times 4$	0	3	3	1	0	0.0	0.46	0.42	0.06	0.0
$16 \times 16$	5	4	4	2	2	0.697	1.091	0.876	0.300	0.181
$32 \times 32$	10	4	4	2	2	1.430	0.994	0.795	0.276	0.179
$64 \times 64$	23	4	4	2	2	2.658	0.954	0.763	0.267	0.178
$61 \times 28$	15	4	4	2	2	1.797	0.971	0.774	0.271	0.180

Table 2: Results for various tilings of 2-dimensional data with 16 disks

Tiling	Maximum					Average				
	HCAM	DM	FX	GFIB	NEW	HCAM	DM	FX	GFIB	NEW
$16 \times 16$	2	12	12	1	1	0.350	2.608	2.392	0.230	0.127
$32 \times 32$	6	16	16	2	2	0.881	4.464	4.040	0.360	0.336
$64 \times 64$	12	16	16	2	3	1.850	5.347	4.515	0.447	0.468

Table 3: Results for various tilings of 2-dimensional data with 64 disks

## 7 Conclusion

Range queries are an important class of queries for several applications including relational databases, spatial databases, and GIS applications. For large datasets, the performance of range queries is limited by disk I/O. 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, no scheme with guaranteed worst-case performance is known. In this paper we have developed an allocation with guaranteed worst-case performance. 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  $OPT + O(1)$  parallel block accesses. The result is generalized to higher dimensions:  $OPT + f(d)$  parallel block accesses for  $d$ -dimensional queries. Experimental data show that the algorithm achieves very close to  $OPT$  performance (on average less than 0.5 away from  $OPT$ , with a worst-case of 3). Although several declustering schemes for multidimensional range queries have been developed, this is the first scheme with a guaranteed non-trivial performance bound.

## References

- [AE97] 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.

Tiling	Disks	Maximum					Average				
		HCAM	DM	FX	GFIB	NEW	HCAM	DM	FX	GFIB	NEW
$2 \times 2 \times 2$	4	1	1	2	1	0	0.074	0.259	0.296	0.074	0.0
$4 \times 4 \times 4$	16	3	8	12	1	1	0.205	1.296	1.252	0.175	0.030
$8 \times 8 \times 8$	64	5	40	56	3	2	0.740	4.832	4.454	0.442	0.215

Table 4: Results for 3-dimensional data

- [BBB<sup>+</sup>97] S. Berchtold, C. Bohm, B. Braunnmuller, D. A. Keim, and H-P. Kriegel. Fast parallel similarity search in multimedia databases. In *Proc. ACM SIGMOD Int. Conf. on Management of Data*, pages 1–12, Arizona, U.S.A., 1997.
- [DS82] 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.
- [FB93] 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.
- [KP88] 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.
- [LSR92] J. Li, J. Srivastava, and D. Rotem. CMD: a multidimensional declustering method for parallel database systems. In *Proceedings of the Int. Conf. on Very Large Data Bases*, pages 3–14, Vancouver, Canada, August 1992.
- [PAAE98] 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.
- [PAE98a] S. Prabhakar, D. Agrawal, and A. El Abbadi. Efficient disk allocation for fast similarity searching. In *Proc. of the 10th Int. Sym. on Parallel Algorithms and Architectures (SPAA'98)*, pages 78–87, Puerto Vallarta, Mexico, June 1998.
- [PAE98b] S. Prabhakar, D. Agrawal, and A. El Abbadi. Efficient retrieval of multidimensional datasets through parallel I/O. In *Proc. of the 5th International Conference on High Performance Computing, (HiPC'98)*, Chennai, India, December 1998.