

## Cyclic Allocation of Two-Dimensional Data

S. Prabhakar<sup>†</sup> K. Abdel-Ghaffar<sup>‡</sup> D. Agrawal<sup>†</sup> A. El Abbadi<sup>†</sup>

<sup>†</sup>Department of Computer Science      <sup>‡</sup>Department of Electrical & Computer Engineering  
University of California                  University of California  
Santa Barbara                                  Davis  
CA 93106, U.S.A.                              CA 95616, U.S.A.

### Abstract

*Various proposals have been made for declustering two-dimensionally tiled data on multiple I/O devices. Recently it has been shown that strictly optimal solutions only exist under very restrictive conditions on the tiling of the two-dimensional space or for very few I/O devices. In this paper we explore allocation methods where no strictly optimal solution exists. We propose a general class of allocation methods, referred to as cyclic allocation methods, and show that many existing methods are instances of this class. As a result, various seemingly ad hoc and unrelated methods are presented in a single framework. Furthermore, the framework is used to develop new allocation methods that give better performance than any previous method and that approach the best feasible performance.*

### 1 Introduction

Efficient retrieval of information defined in a two-dimensional space is gaining in importance. An increasingly popular and user-friendly approach for retrieving geographically referenced information is to request the user to define a rectangular bounding box on a map indicating the area of interest. In the Alexandria Digital Library project [CDN<sup>+</sup>97], we are exploring various methods for supporting such geographically referenced data. For example, if a user is interested in more information about the Sierra Nevada mountains, the user can draw a rectangle around the area on a map of the United States. As a result, more detailed maps can be retrieved, or in general, any information related to the Sierra Nevada mountains, for example, amount of rainfall, elevation information, or even papers studying the mountain range.

In general, the amount of information associated with different parts of a two-dimensional map, or image, may be quite large. Hence, the map is often di-

vided into tiles, where each tile corresponds to the information associated with the area defined by the map. Since user queries will, in general, refer to multiple tiles, the response time to the query is proportional to the time it takes to retrieve all information associated with the tiles overlapping the query region. A common approach for improving query response time is to decluster the information of the tiles on different I/O devices, e.g., disks. Hence, for a given query, if all tiles are allocated to different devices, the response time is minimized since data can be retrieved from multiple devices in parallel. On the other hand, if they are allocated to the same device, response time is maximized, since the tiles must be retrieved sequentially instead of concurrently.

Various proposals have been made for the allocation of data on multiple I/O devices [DS82, KP88, FB93]. Some of these studies were conducted in the context of relational databases, where data is defined over several attributes. We address the important but restricted case of two attributes (the  $x$  and  $y$  coordinates of two-dimensional data). Recently in [AE97], it was shown that strictly optimal solutions (defined shortly) for rectangular queries exist only under very restrictive conditions on the tiling or for very few I/O devices (namely 1,2,3 and 5). In this paper we explore allocation methods where no strictly optimal solution exists. We propose a general class of allocation methods, referred to as *cyclic* allocation methods, and show that many existing methods are instances of this class. This allows us to present various seemingly ad hoc and unrelated methods in a unified framework. It also allows us to develop new allocation methods that give better performance than any previous method. We perform extensive analysis of various allocation schemes to understand their relative performance.

The rest of the paper is organized as follows. Section 2 presents the existing approaches and their performance. In Section 3 we develop new allocation schemes and show their superior performance. The

<sup>†</sup>Work supported by a research grant from NSF /ARPA /NASA IRI9411330 and NSF grant CDA-9421978.

<sup>‡</sup>Work supported by NSF grant NCR 96-12354.

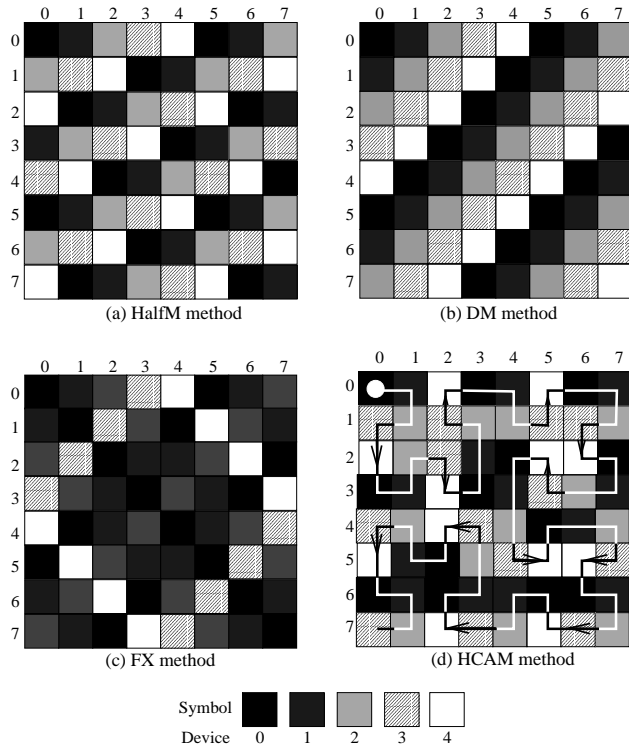


Figure 1: Allocation generated by existing methods

robustness of these new schemes is evaluated in Section 4 and Section 5 concludes the paper.

## 2 Declustering for Parallel Retrieval

In this section we present the more prominent existing allocation approaches and evaluate their performance. First we describe the existing approaches and then show their performance relative to each other and with respect to a Random allocation scheme.

### 2.1 Allocating Two-Dimensional Data

Consider a two-dimensional data set that has been divided into tiles along each dimension. These tiles are to be allocated to several I/O devices so as to minimize the maximum number of tiles accessed from a single I/O device when a group of tiles contained in a rectangular region of the data set is retrieved.

Given a two-dimensional array of  $N_1 \times N_2$  tiles, and  $M$  I/O devices, the allocation function maps each tile  $(i, j)$ ,  $0 \leq i < N_1, 0 \leq j < N_2$  to one of the I/O devices  $0, \dots, M-1$ . To process a rectilinear rectangular query, all tiles that intersect with the query need to be accessed. The cost of executing the query is proportional to the maximum number of tiles accessed from a single I/O device. The minimum possible cost when retrieving  $A$  tiles distributed over  $M$  devices is  $\lceil \frac{A}{M} \rceil$ . An allocation policy is said to be *strictly optimal* if no

rectangular area  $A$  has more than  $\lceil \frac{A}{M} \rceil$  tiles allocated to the same device. In [AE97], the necessary and sufficient conditions for the existence of strictly optimal allocation schemes are derived. In particular it is shown that strictly optimal allocations exist for the following cases only: i) any number of I/O devices if the number of tiles in one of the dimensions is 1 or 2, or ii) any number of tiles in each dimension if the number of devices is 1, 2, 3 or 5, or iii) for some special boundary cases (the number of devices is almost as large as the total number of tiles, i.e.  $M \geq N_1 N_2 - 2$ , or  $M=8$ , and  $N_1=N_2=4$ , or  $M=N_1 N_2 - 4$  and  $\min\{N_1, N_2\} = 3$ ). Strictly optimal allocations for each of these cases are also defined in [AE97]. Of these cases, only the second case is of general interest. To demonstrate sufficiency of these conditions, a strictly optimal allocation method is given, which allocates tile  $(i, j)$  to device  $(\lfloor \frac{M}{2} \rfloor i + j) \bmod M$ . We will refer to this as the *HalfM* method. Figure 1(a) shows the allocation generated by the HalfM method for  $N_1=N_2=8$  and  $M=5$ . Since it has been established that except for these restricted cases, no strictly optimal allocation exists, in this paper we look for methods that give good performance for the cases for which no strictly optimal solution exists. Several allocation methods have been proposed in the literature.

The Disk Modulo (DM) [DS82] approach proposed by Du and Sobolewski and later extended for range queries and dynamic files in [LSR92] allocates tile  $(i, j)$  to device  $(i + j) \bmod M$ . The FX method proposed by Kim and Pramanik [KP88], requires that  $N_1, N_2$  and  $M$  are all powers of 2. The tile  $(i, j)$  is then allocated to the device given by the lowest  $\log_2 M$  bits of the bit-wise exclusive-OR of the binary representations of  $i$  and  $j$ . The method can easily be extended to be applicable to values of  $N_1, N_2$  and  $M$  that are not powers of 2. Thus the tile  $(i, j)$  is allocated to device  $k$  given by:  $k = (i_B \oplus j_B) \bmod M$  where  $i_B$  and  $j_B$  are the binary representations of  $i$  and  $j$  respectively and  $\oplus$  is the bit-wise exclusive-OR operation. The HCAM method [FB93] is based upon the Hilbert space filling curves. Hilbert curves can be used to convert a discrete two-dimensional array into a linear sequence such that spatial proximity in the original array is preserved. Given an  $N_1 \times N_2$  array of tiles, the HCAM method first converts the array into a linear sequence, and then allocates the tiles to devices in a round-robin fashion following the linear sequence. Since spatial proximity is captured by the linear sequence, it is expected that this method will give good declustering of the tiles. Figures 1(b) through (d) show the allocations generated by the DM, FX and HCAM meth-

ods respectively. The linear path corresponding to the Hilbert curve is also shown in Figure 1(d).

## 2.2 Performance Comparison

We now present a comparison of the performance of these methods. The comparison is based upon the exhaustive evaluation of all possible queries. Given an  $N_1 \times N_2$  array of tiles and  $M$  I/O devices, each of the tiles is allocated to an I/O device using the various methods. In all evaluations,  $N_1$  was taken to be equal to  $N_2$  although this restriction is not critical to the evaluation or the performance of the various methods, as shown in Section 4. The effectiveness of the allocation in providing parallel I/O for rectangular queries is then evaluated. Consider a given allocation method. The cost of a query of area  $A$  is taken to be the maximum number of tiles retrieved from the same I/O device divided by the strictly optimal cost,  $\lceil \frac{A}{M} \rceil$ . The average cost of queries of area  $A$  is obtained by averaging the costs of all possible queries of area  $A$ . We then average the cost for each area for all areas from 2 through  $N_1 N_2$ . This gives a measure of performance for the given allocation method. Clearly, this performance measure is at least 1, and the smaller the value, the better the allocation method for providing parallel I/O. For strictly optimal allocation, the performance measure equals 1.

In Figure 2 the performance of the four methods, HalfM, DM, FX and HCAM, is shown. For the purpose of comparison, the performance of a Random approach, where each tile is allocated to a randomly chosen I/O device, is also shown. The value of  $N_1$  and  $N_2$  was fixed at 32, while  $M$  was varied from 2 to 32. The  $x$ -axis of the graph shows the value of  $M$ , the number of I/O devices available, and the  $y$ -axis gives the ratios.

From Figure 2, the following observations can be made. The DM approach gives good performance consistently. It is the best approach for almost half the values of  $M$ . The FX approach does not perform as well as the DM approach for all values of  $M$ , except for the cases when  $M$  is a power of 2. When the value of  $M$  is a power of 2, i.e.  $M = 2, 4, 8, 16$  and 32, FX is the best approach. It should be pointed out that these are the values for which the FX approach was originally designed. We have made a natural extension to the approach for other values of  $M$ , however, as seen from the results, the approach is not as effective for other values. The HCAM approach performs worse than DM for most values of  $M$  and FX for powers of 2. The DM and HCAM methods have better performance than the Random policy for all values of  $M$ . The HalfM and FX methods perform better than the

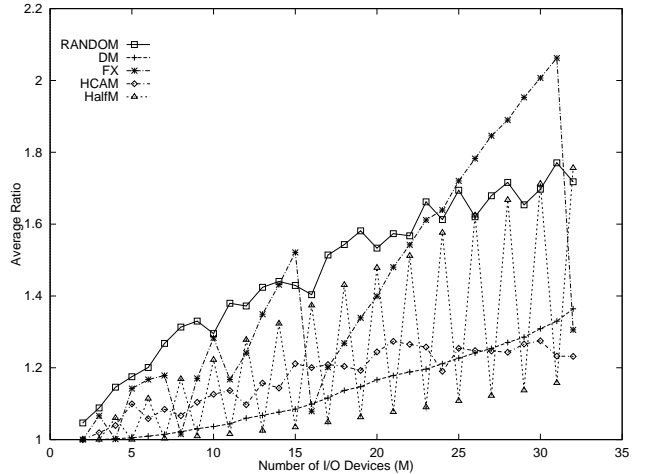


Figure 2: Performance of existing schemes

Random method for most values of  $M$ . As expected, HalfM has the strictly optimal behavior for small values of  $M = 2, 3$  and 5, for which it was designed. On the other hand, for larger values of  $M$ , HalfM has very interesting behavior. For odd values of  $M$  it gives the best performance among all the approaches, but for even values of  $M$  it has worse performance than both DM and HCAM, and in some cases worse than FX. Note that HalfM was developed primarily to establish the sufficiency of the optimal allocation scheme.

Overall, we observe that there is no single existing approach that gives the best performance for all or even most values of  $M$ . It should be noted that as expected from the results obtained in [AE97], none of the policies produces strictly optimal results, except for the special cases. DM is strictly optimal for  $M = 2$  and 3, FX for  $M = 2$  and HalfM for  $M = 2, 3$  and 5.

Experiments were also conducted with  $N_1 = N_2 = 64$  and  $M$  varying from 2 through 64. We found that there were no significant differences in the relative performance seen for  $N_1 = N_2 = 32$  and  $N_1 = N_2 = 64$ . This insensitivity to the number of tiles in each dimension was observed in all our experiments. For this reason, in the paper we present only the results for  $N_1 = N_2 = 32$  since these graphs are less dense. Results of other experiments not reported here are available in [PAAE97]. Also for clarity, in the remainder of the paper, the Random method is not shown because of its poor performance. The performance of FX for values of  $M$  that are not powers of 2 is observed to be quite poor. Since the scheme was originally designed only for values of  $M$  that are powers of 2, in the remainder of the paper, the performance of FX is shown only for these values.

### 3 General Declustering Paradigm

We first develop a general class of allocation methods called *cyclic* allocation methods. A criterion for identifying those methods that will give better performance is established. We then investigate in detail, three different ways of finding such methods.

#### 3.1 Cyclic Declustering Methods

We now introduce a set of allocation methods based upon the idea of *cycles*. Consider the set of  $M$  devices as being mapped onto a circle with  $M$  tick marks. The allocation of tiles to devices can be viewed as hops on this circle as one moves from one tile to another. The goal of the declustering of tiles to devices is to avoid allocating tiles that are close to each other to the same device. Consider the sequence of hops as we move along a row of the data set, corresponding to certain rectangular queries. Clearly it is desirable that the tick marks visited should not be repeated as much as possible. Similarly, when moving down a column, it is desirable that tick marks are not repeated as much as possible. Extending the argument, starting from a given tile, the hops along its row and along its column should be as distinct as possible. Graphically, one can view this as two sequences of hops along the cycle of devices which minimize the number of times any given tick mark is visited as shown in Figure 3.

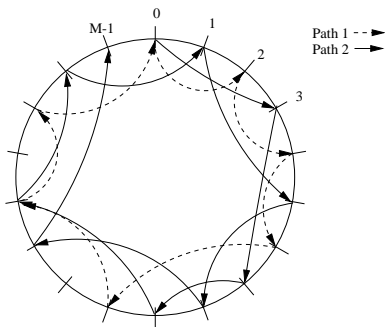


Figure 3: Visualization of paths on device cycle

Based upon the idea of paths on a cycle, a good allocation method should try to visit each tick mark before repeating any of the tick marks and also, the two paths should not visit the same tick marks as far as possible. One way of achieving this is for the paths to consist of equi-distance hops. If the number of devices,  $M$  and the size of the hop are relatively prime, i.e. their greatest common divisor (GCD) is 1, then each path will visit every tick mark before repeating any tick mark. Also, in order to minimize the overlap between the two paths for the row and the column, the hop values for these two should also be relatively

prime with respect to each other. Therefore, given a value of  $M$ , we need to find two other numbers which are relatively prime with respect to each other and to  $M$ . The problem of finding these two numbers can be simplified by picking one of them to be 1. This leads to a set or class of allocation methods for any given value of  $M$ , depending upon the value of the other hop chosen. We call this a class of *cyclic* allocation methods. This class is a subset of the class of allocations defined by the GDM method introduced in [DS82] and studied in [CC92]. As in the cyclic methods, the GDM approach, also requires the two hop values,  $a$  and  $b$ , to be relatively prime with respect to  $M$ . However, it does not require that  $a$  and  $b$  be relatively prime with respect to each other. Moreover, effective means of determining values of  $a$  and  $b$  that give good performance are not described.

Given a value of  $M$ , each cyclic allocation method is distinguished by the value of the hop,  $H$  that it uses. Each cyclic method allocates tiles to devices in the following manner. Tile  $(i, j)$  is allocated to device  $(Hi + j) \bmod M$ . Another way of looking at the allocation is that along any row tiles are allocated to devices sequentially, i.e., in steps of 1, and along a column tiles are allocated to devices in steps of  $H$ .

Thus for each value of  $M$ , we get a class of  $M$  allocation methods with hop values,  $H = 0, \dots, M-1$ . Interestingly, some of the methods that have been proposed earlier are included in this class of methods. In particular, the DM method is a cyclic method with  $H = 1$ , and HalfM is a cyclic method with  $H = \lfloor \frac{M}{2} \rfloor$ .

#### 3.2 HalfM Cycles for Declustering

As mentioned earlier, a cyclic method that chooses a value of the hop,  $H$  ( $H \neq 1$ ), that is relatively prime with respect to  $M$  will perform well. Consequently, the HalfM scheme gives very good performance for odd values of  $M$  and very poor performance for even values of  $M$  (Figure 2). We therefore choose  $H$  to be a relatively prime value close to  $\frac{M}{2}$ . The scheme resulting from this choice is called *Relatively Prime Half M* or RPHM for short. It chooses  $H$  to be

$$H = \begin{cases} \frac{M}{2} + 2 & \text{if } M = 2 + 4k \quad k \in \mathcal{N} \\ \frac{M}{2} + 1 & \text{if } M = 4k \quad k \in \mathcal{N} \\ \lfloor \frac{M}{2} \rfloor & \text{otherwise} \end{cases}$$

The performance of RPHM is shown in Figure 4. It outperforms DM, FX and HCAM for all values of  $M$ .

#### 3.3 Fibonacci Cycles for Declustering

In the previous subsection we were able to modify the HalfM method to obtain a more efficient cyclic allocation method, RPHM. In this subsection we consider another method for identifying a good cyclic

method that uses a hop value that is relatively prime with respect to  $M$ . The method is motivated by an allocation based upon certain properties of Fibonacci numbers proposed by Chor, Leiserson, Rivest and Shearer [CLRS86], for the allocation of two-dimensional screen pixels to memory chips. We will call this method FIB. The FIB method is applicable only to very few values of  $M$ , in particular, only for values of  $M$  that are odd order Fibonacci numbers, i.e.  $M = F_{2r+1}$  where  $r \geq 2$  and  $F_x$  is the  $x$ th Fibonacci number. The sequence of Fibonacci numbers is defined as  $F_0 = 0, F_1 = 1$  and  $F_i = F_{i-1} + F_{i-2}$   $i \geq 2$ . Thus the FIB method is applicable only if the number of I/O devices is 5,13,34,89, etc. The FIB method allocates all tiles  $(i', j')$  which satisfy

$$\begin{pmatrix} i' \\ j' \end{pmatrix} = \begin{pmatrix} i \\ j \end{pmatrix} + a \begin{pmatrix} F_r \\ F_{r+1} \end{pmatrix} + b \begin{pmatrix} -F_{r+1} \\ F_r \end{pmatrix}$$

for integral values of  $a$  and  $b$ , to the same device as the tile  $(i, j)$ , for  $M = F_{2r+1}$ . Along the first row, the tiles are allocated to devices in a round robin fashion, similar to the cyclic methods. Therefore, tile  $(0, j)$  is allocated to device  $j \bmod M$ .

The requirement of the FIB method that the number of I/O devices must be an odd order Fibonacci number is quite restrictive. More importantly, unlike the FX method, there is no obvious technique to extend the FIB method to arbitrary values of  $M$ . We begin by first showing that FIB, like DM and HalfM, is also a cyclic allocation method and then based upon this observation we propose a technique to extend its applicability to all values of  $M$  by developing a cyclic method, called Generalized FIB or GFIB.

We will show that FIB is a cyclic allocation method with a hop  $H = (-1)^r F_{2r}$  for  $M = F_{2r+1}$ . To prove this, we use the following two well-known Fibonacci identities, which can be easily proved by induction:

$$F_r^2 + F_{r+1}^2 = F_{2r+1}, \quad r \geq 0 \quad (1)$$

$$F_{n+1}F_t - F_nF_{t+1} = (-1)^{t+1}F_{n-t}, \quad n \geq t \geq 0 \quad (2)$$

Substituting  $n = 2r$  and  $t = r$  in (2) yields

$$F_r - (-1)^r F_{2r}F_{r+1} \equiv 0 \pmod{F_{2r+1}}. \quad (3)$$

Similarly, substituting  $n = 2r$  and  $t = r - 1$  yields

$$F_{r+1} + (-1)^r F_{2r}F_r \equiv 0 \pmod{F_{2r+1}}. \quad (4)$$

Now we will show that the cyclic allocation method defined by the hop  $H = (-1)^r F_{2r}$ , where  $M = F_{2r+1}$  is equivalent to FIB allocation. Notice that tiles  $(i, j)$

and  $(i', j')$  are assigned to the same device in the cyclic allocation if and only if

$$(-1)^r F_{2r}(i' - i) + (j' - j) \equiv 0 \pmod{F_{2r+1}}. \quad (5)$$

On the other hand the same tiles are assigned to the same device in FIB allocation if and only if the values of  $a$  and  $b$  satisfying the two equations

$$aF_r - bF_{r+1} = i' - i \quad (6)$$

$$aF_{r+1} + bF_r = j' - j \quad (7)$$

are integers. By making use of (1), we get

$$a = \frac{(i' - i)F_r + (j' - j)F_{r+1}}{F_{2r+1}} \quad (8)$$

$$b = \frac{-(i' - i)F_{r+1} + (j' - j)F_r}{F_{2r+1}}. \quad (9)$$

It is easy to check that if (5) holds, then (3) and (4) imply that  $a$  and  $b$  in (8) and (9) are integers. Hence, if  $(i, j)$  and  $(i', j')$  are assigned to the same device in the cyclic allocation, then they are also assigned to the same device in FIB allocation. Conversely, if  $a$  and  $b$  are integers, then adding (6) multiplied by  $(-1)^r F_{2r}$  to (7), and using (3) and (4), gives (5). Hence, if  $(i, j)$  and  $(i', j')$  are assigned to the same device in FIB allocation, then they are also assigned to the same device in the cyclic allocation. This proves the equivalence of the two allocation methods. Note that a negative hop value of  $H$  is equivalent to a positive hop value of  $M + H$  because the allocation is done in modulo  $M$ . Hence we see that the FIB method is equivalent to a cyclic method with hop values chosen to be  $H = F_{2r}$  if  $r$  is even and  $H = M - F_{2r} = F_{2r-1}$  if  $r$  is odd.

We now generalize the FIB method to any value of  $M$ , based upon the above observation. In order to be applicable to more general values of  $M$ , we need to move from the integer domain into the real domain for the Fibonacci sequence and perform a similar function to that done in the discrete domain, i.e. given  $M$  corresponding to a Fibonacci number, we need to find the preceding Fibonacci number. The closed-form equation for the Fibonacci numbers provides a convenient means to do this. In particular, given a value of  $M$ , we find, using the closed-form equation, the corresponding Fibonacci index,  $i$ . We then determine the Fibonacci number corresponding to index  $i - 1$ . Unfortunately, since we are in the continuous domain, the resulting value may not be relatively prime with respect to  $M$ . We therefore use this value as an approximation and then search for the closest relatively prime value. The new method, called GFIB, is a cyclic

method which when given a value of  $M$ , begins by choosing the hop,  $H$  to be  $H = F(F^{-1}(M) - 1)$ , where  $F(x)$  represents the closed-form equation for the  $x$ th Fibonacci number:  $F(x) = \frac{\phi^x - \hat{\phi}^x}{\sqrt{5}}$ , where  $\phi$  is the golden ratio  $\phi = \frac{1+\sqrt{5}}{2}$ ,  $\hat{\phi}$  is its complement,  $\hat{\phi} = \frac{1-\sqrt{5}}{2}$ . Although Fibonacci numbers are defined only for integral values of  $x$ , we allow  $x$  to take on non-integral values. The function  $y = F^{-1}(x)$  represents the (possibly non-integral) value which satisfies  $x = \frac{\phi^y - \hat{\phi}^y}{\sqrt{5}}$ . If the value of  $M$  happens to be a Fibonacci number ( $M = F_k$ ) then the value of hop,  $H$  will be the previous Fibonacci number ( $H = F_{k-1}$ ) and thus  $H$  and  $M$  will be relatively prime. For other values of  $M$ , however, the value of  $H$  may not be relatively prime with respect to  $M$ . If the value of  $H$  given above and  $M$  are relatively prime, then the hop value is taken to be  $H$ , otherwise, we choose a value of  $H$  that is as close to the value given by the above formula as possible and is relatively prime with respect to  $M$ . We do this by considering the neighboring values in the following order:  $H-1, H+1, H-2, H+2, H-3 \dots$ . The performance of GFIB relative to the other methods is shown in Figure 4. We observe that GFIB performs extremely well. It outperforms all other methods except RPHM for some values of  $M$  (and the EXH scheme which will be introduced soon).

### 3.4 The Exhaustive cyclic method

In order to obtain a qualitative measure for the “goodness” of our methods and how close we are to the best achievable performance, we also considered a policy which chooses the best possible value of the hop from among all possible values ( $0, \dots, M-1$ ) by evaluating all the cyclic methods for each value of  $M$  and picking one which gives the best performance. The performance of this, “Exhaustive”, or EXH strategy is shown in Figure 4.

The EXH method represents the current limit on the performance of allocation methods. It indicates how close we can get to the optimal even though no strictly optimal allocation exists. It is interesting to note that the GFIB method has performance that is very close to that of EXH for most values of  $M$ . Although evaluating the best hop value is useful, it is an expensive operation. Due to the extensive computation required to exhaustively search each hop value, only a subset of queries was considered in evaluating the relative performance of the hop values for each value of  $M$ . Only queries with area less than or equal to  $N_1$  were considered, instead of considering all possible queries. It is interesting that using the hop values determined by considering only a subset of the queries,

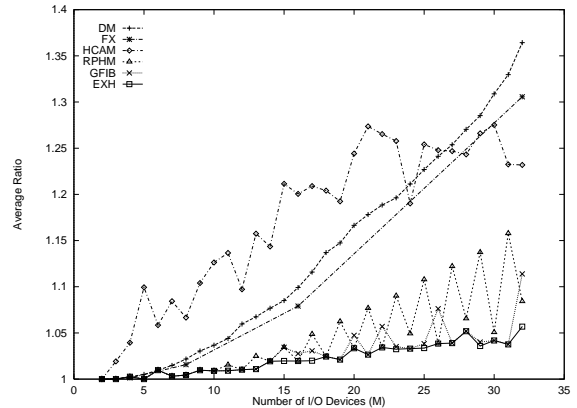


Figure 4: Performance of the new Cyclic schemes

the EXH method gives the best performance when all queries are considered, as seen in Figure 4. Thus we are able to reduce the computation required to determine the best hop value. The search space can be further reduced by considering only values of  $H$  from 0 to  $\lfloor \frac{M}{2} \rfloor$ .

## 4 Robustness of the new approaches

In Section 3, the superiority of the new *cyclic* methods RPHM, GFIB and EXH was demonstrated by their performance for all possible queries. The types of queries that are considered is an important factor in the performance of the allocation scheme. In this section we consider the performance of the new approaches using different types of queries. We also demonstrate the robustness of the new approaches with respect to variations in the numbers of tiles in each dimension, and the number of I/O devices. Results of further experiments demonstrating the robustness of our approach to varying parameters are available in [PAAE97].

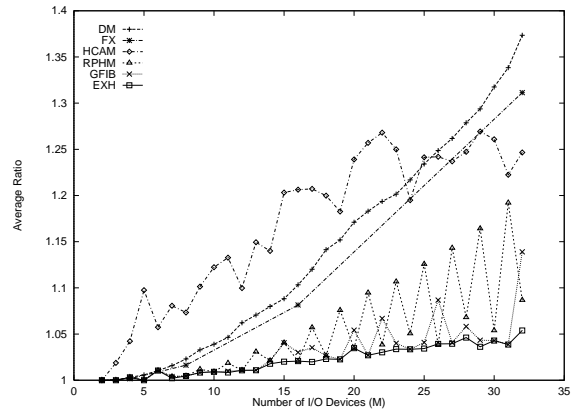


Figure 5: Performance for *Tall* queries only

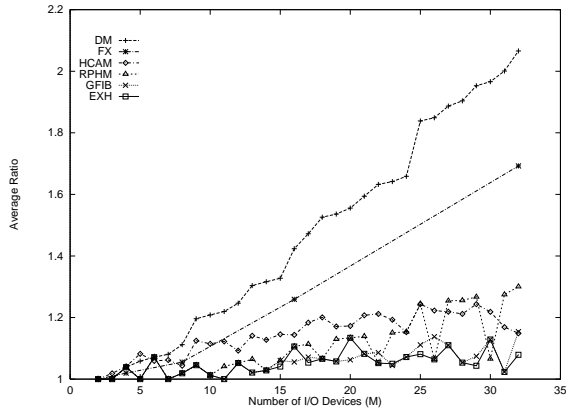


Figure 6: Performance for *Square* queries only

#### 4.1 Effect of Varying the Queries

Let us first consider only those queries for which the number of columns is less than or equal to the number of rows. Since for these the vertical side is longer than the horizontal side, we call them *tall* queries. For the Random, FX and DM approaches, the distinction is unimportant. For the cyclic approaches, the notion is used only to decide which dimension uses a hop of 1, and which uses a different value of the hop depending on the particular cyclic approach. Figure 5 shows the performance of the various approaches for *tall* queries. We note that the relative performance of the strategies does not change. The three new approaches perform as well as they did for all queries. In particular, all the three new methods, RPHM, GFIB and EXH, give better performance than the existing methods FX, DM and HCAM. Similar results were obtained when only long queries (number of rows is greater than or equal to number of columns) and only small queries (with areas less than or equal to  $N_1$ ) were considered.

We now consider the special case of queries with the same number of rows and columns, called *square* queries. This is an interesting case because the HCAM approach [FB93] has been shown to have very good performance for square queries. Figure 6 gives the performance of the various approaches for square queries. We observe that the HCAM approach performs better than DM and FX for almost all values of  $M$ . However, HCAM does not perform better than RPHM, EXH or GFIB. These observations are in contrast to those seen in Figure 2 in Section 2, which showed that when we consider all types of queries, not just square queries, HCAM is poorer than DM and FX.

It is interesting that the relative performance of all three new approaches does not change significantly for the different query types. EXH always gives the best

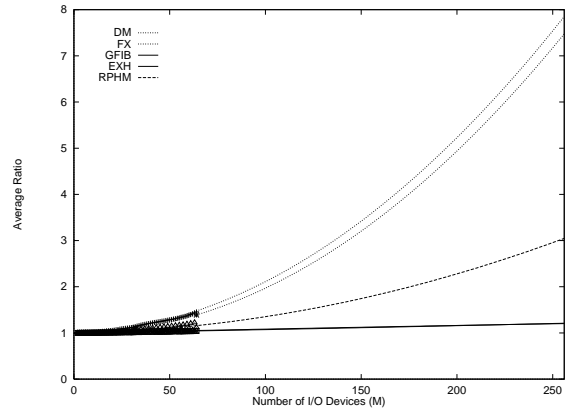


Figure 7: Extrapolation to  $M = 256$

performance and GFIB has a performance very close to that of EXH. The RPHM method performs well for even values of  $M$  and reasonably well for odd values. Though it does not perform as well as EXH or GFIB, it is always much better than all the existing methods, DM, FX and HCAM.

#### 4.2 Effect of Varying $N_1, N_2$ and $M$

In the preceding sections, it was assumed that  $N_1=N_2=32$  and  $M=32$ . In this subsection, we present some of the results with different values for these parameters. Experiments were repeated for  $N_1=N_2=64$  and  $M=2, \dots, 64$ . We found that there is no qualitative difference in the relative performance of the various methods. Most importantly, GFIB and EXH give the best performance and RPHM also performs quite well in comparison to the existing approaches. Thus we see that the new cyclic methods are insensitive to variations in the numbers of tiles in each dimension and also the number of I/O devices in addition to variations in the query types.

In order to experimentally predict the behavior of the various allocation methods for large values of  $M$ , we used the least square fitting method with quadratic polynomials to the various plots. Figure 7 depicts the extrapolation of the fitted curves to larger values of  $M$ . Based upon the results, it can be stated that RPHM, GFIB and EXH diverge from the strictly optimal values at much slower rates than those that have been proposed elsewhere (DM and FX). We did not consider HCAM in this comparison due to its performance sensitivity to different query types.

## 5 Discussion

It has been shown that strictly optimal allocations for two-dimensional data on different I/O devices exist only in very limited cases [AE97]. We proposed a

Number of I/O devices $M$	% Improvement over		
	DM	HCAM	FX
5	100	100	100
10	75.749	92.953	96.86
15	59.459	83.680	93.380
20	71.531	80.589	88.135
25	83.076	84.887	94.669
30	86.468	84.811	95.848
16	72.177	86.234	65.107
32	68.735	50.884	62.741

Table 1: Percentage reduction in gap.

general class of cyclic declustering methods applicable in these cases. We have shown that the Disk Modulo [DS82], HalfM [AE97] and the Fibonacci [CLRS86] allocation methods are special cases of this class. Also, showing that the Fibonacci method is a cyclic allocation method allowed us to generalize it for all numbers of I/O devices (in the original method, this was restricted to odd ordered Fibonacci numbers) resulting in the generalized Fibonacci, or GFIB method. We have also developed two other techniques (RPHM and EXH) for identifying cyclic methods that produce efficient allocations. The GFIB, RPHM as well as the EXH methods are all instances of the cyclic allocation that give superior performance to the best prior proposed methods. Consider for example the improvement over the existing policies that is achieved by the GFIB method. In Table 1, we list the percentage reduction in the difference between each of the existing methods - DM, HCAM and FX - and the strictly optimal (corresponding to a ratio of 1), achieved by GFIB for  $N_1=N_2=32$  and various values of  $M$ . Note that for the last two rows,  $M$  is a power of 2, for which FX was designed. We note that the reductions are significant, especially since no strictly optimal allocation exists for most cases. For those cases where a strictly optimal allocation does exist, (e.g.  $M=5$ ), the gap is completely eliminated. In our experiments, we found that the GFIB method is within 5% of the strictly optimal 90% of the time and the EXH method is within 5% of the strictly optimal 93% of the time. Moreover, the new methods give the best performance irrespective of the types of queries, the number of tiles or the number of I/O devices.

## References

- [AE97] K. A. S. Abdel-Ghaffar and A. El Abbadi. Optimal allocation of two-dimensional data. In *International Conference on Database Theory*, pages 409–418, Delphi, Greece, January 1997.
- [CC92] C. Y. Chen and C. C. Chang. On the gdm allocation method for partial range queries. *Information Systems*, 17(5):381–394, 1992.
- [CDN<sup>+</sup>97] X. Cheng, R. Dolin, M. Neary, S. Prabhakar, K. Ravikanth, D. Wu, D. Agrawal, A. El Abbadi, M. Freeston, A. Singh, T. Smith, and J. Su. Scalable access within the context of digital libraries. In *IEEE International Conference on Advances in Digital Libraries, ADL*, pages 70–81, Washington, D.C., 1997.
- [CLRS86] B. Chor, C. E. Leiserson, R. L. Rivest, and J. B. Shearer. An application of number theory to the organization of rastergraphics memory. *Journal of the Association for Computing Machinery*, 33(1):86–104, January 1986.
- [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 *Proceedings of the 2nd International Conference 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.
- [PAAE97] S. Prabhakar, K. Abdel-Ghaffar, D. Agrawal, and A. El Abbadi. Cyclic allocation of two-dimensional data. Technical Report TRCS97-08, Dept. of Computer Science, Univ. of California, Santa Barbara, June 1997.