

- [15] W. Liu and V.K. Prasanna. Design of Application Software for Embedded Signal Processing. In *IEEE Signal Processing Magazine*, September 1998.
- [16] Message Passing Interface Forum. MPI: A Message-Passing Interface Standard. *International Journal of Supercomputer Applications and High Performance Computing*, Vol.8, No. 3-4, 1994
- [17] C. H. Papadimitriou, K. Steiglitz. Combinatorial Optimization: Algorithms and Complexity. Prentice Hall, 1982, pp. 225-226.
- [18] L. Prylli and B. Tourancheau. Fast Runtime Block Cyclic Data Redistribution on Multiprocessors. *Journal of Parallel and Distributed Computing*, volume 45, August 1997
- [19] S. Ramaswamy and P. Banerjee. Automatic Generation of Efficient Array Redistribution Routines for Distributed Memory Multicomputers. In *Proc. of 5th Symp. Frontiers of Massively Parallel Computation, McLean, VA*, pages 342-349, February 1995.
- [20] J. C. Setubal. Sequential and Parallel Experimental Results with Bipartite Matching Algorithms. Technical Report IC-96-09, Institute of Computing, State University of Campinas (Brazil), 1996.
or <http://www.cs.sunysb.edu/~algorithm/implement/bipm/implement.shtml>
- [21] J. Suh, M. Ung, and V.K. Prasanna. Parallel Implementation of Synthetic Aperture Radar on High Performance Computing Platforms. *International Conference on Algorithms And Architectures for Parallel Processing '97*, December 1997.
- [22] R. Thakur, A. Choudhary, and G. Fox. Runtime Array Redistribution in HPF Programs. In *Proc. of Scalable High Performance Computing Conference*, pages 309-316, May 1994.
- [23] R. Thakur, A. Choudhary, and J. Ramanujam. Efficient Algorithms for Array Redistribution. *IEEE Trans. on Parallel and Distributed Systems*, Vol. 7, No. 6, pages 587-594, June 1996.
- [24] D.W. Walker and S.W. Otto. Redistribution of Block-Cyclic Data Distributions Using MPI. *Concurrency:Practice and Experience*, Vol. 8, No. 9, pages 707-728, November 1996.
- [25] C.-L. Wang, P.B. Bhat, and V.K. Prasanna. High-Performance Computing for Vision. *Proceedings of IEEE*, 84:931-946, 1996.

References

- [1] L. Blackford, J. Choi, A. Cleary, E. D’Azevedo, J. Demmel, I. Dhillon, J. Dongarra, S. Hammarling, G. Henry, A. Petitet, K. Stanley, D. Walker, and R. C. Whaley. *ScaLAPACK Users’ Guide*. SIAM Publications, 1997.
- [2] J. Bruck, C.-H. Ho, S. Kipnis, and Weathersby. Efficient Algorithms for All-to-All Communications in Multi-Port Message-Passing Systems. In *6th Annual ACM Symp. on Para. Alg. and Arch.*, pages 298–309, July 1994.
- [3] J. Choi, J. Dongarra, and D. Walker. Parallel Matrix Transpose Algorithms on Distributed Memory Concurrent Computers. *Proceedings of Fourth Symposium on the Frontiers of Massively Parallel Computation (McLean, Virginia)*, 1993.
- [4] Y.C. Chung, C.H. Hsu, and S.W. Bai. A Basic-Cycle Calculation Technique for Efficient Dynamic Data Redistribution. In *IEEE Trans. on Parallel and Distributed Systems*, Vol. 9, No. 4, April 1998.
- [5] F. Desprez, J. Dongarra, A. Petitet, C. Randriamaro, and Y. Robert. Scheduling Block-Cyclic Array Redistribution. In *IEEE Trans. on Parallel and Distributed Systems*, Vol. 9, No. 2, February 1998.
- [6] E. A. Dinic. Algorithm for Solution of Maximum Flow in a Network with Power Estimation. *Soviet Math Dokl.*, Vol. 11, 1970, pp. 1277-1280.
- [7] S. Hiranandani, K. Kennedy, J. Mellor-Crummey, and A. Sethi. Compilation Techniques for Block-Cyclic Distributions. In *Proc. of Int’l. Conf. on Supercomputing*, pages 392–403, July 1994.
- [8] E.T. Kalns and L.M. Ni. Processor Mapping Techniques Toward Efficient Data Redistribution. *Proc. of International Parallel Processing Symposium*, April 1994.
- [9] S.D. Kaushik, C.-H. Huang, R.W. Johnson, and P. Sadayappan. An Approach to Communication Efficient Data Redistribution. In *Proc. of Int’l. Conf. on Supercomputing*, pages 364–373, July 1994.
- [10] S.D. Kaushik, C.-H. Huang, J. Ramanujam, and P. Sadayappan. Multiphase Array Redistribution: Modeling and Evaluation. In *Proc. of Int’l. Parallel Processing Symposium*, pages 441–445, 1995.
- [11] C. Koelbel, D. Loveman, R. Schreiber, G. Steele Jr., and M. Zosel. *The High Performance Fortran Handbook*. The MIT Press, 1994.
- [12] Y.W. Lim, P.B. Bhat, and V.K. Prasanna. Efficient Algorithms for Block-Cyclic Redistribution of an Array. In *Proc. IEEE Symposium on Parallel and Distributed Processing*, October 1996.
- [13] Y.W. Lim and V.K. Prasanna. Scalable Portable Implementations of Space-Time Adaptive Processing. In *10th Inter. Conf. High Perf. Comp.*, 1996.
- [14] W. Liu, W. Kostis, and V.K. Prasanna. Communication Issues in Heterogeneous Embedded Systems. In *Proc. of Workshop on Para. and Dist. Real Time Sys.*, April 1996.

Table 6: Comparison of schedule computation time ($\mu secs$) with data transfer time

Case	Array Size ($\times 10^6$)	Data Transfer (DT)	Our algorithm		Bipartite Matching	
			Schedule (OST)	Ratio(%) (OST/DT)	Schedule (BST)	Ratio(%) (BST/DT)
(a) $\mathfrak{R}_{16}(18,9,78)$	0.809	21.5	0.165	0.767	7.59	35.30
	8.09	49.5	0.165	0.333	7.59	15.33
	16.2	71.4	0.165	0.231	7.59	10.63
(b) $\mathfrak{R}_{16}(30,15,66)$	1.43	20.2	0.144	0.713	7.74	35.34
	14.3	47.7	0.144	0.302	7.74	14.83
	28.5	69.6	0.144	0.211	7.74	11.12
(c) $\mathfrak{R}_{16}(40,25,56)$	0.896	19.60	0.173	0.883	10.18	51.94
	8.96	35.4	0.173	0.489	10.18	28.76
	17.9	46.80	0.173	0.370	10.18	21.75

Note: Array size is the number of array elements in single precision.

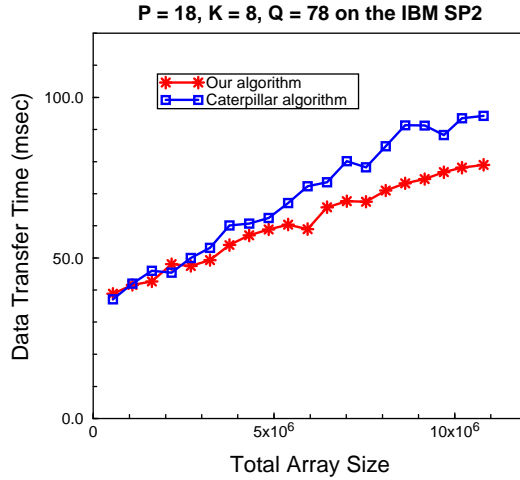
6 Conclusions

In this paper, we derive an efficient algorithm for performing redistribution from $cyclic(x)$ on P processors to $cyclic(Kx)$ on Q processors. The proposed algorithm is based on a generalized circulant matrix formalism. Our algorithm minimizes the number of communication steps and avoids destination node contention in each communication step. The network bandwidth is fully utilized by ensuring that messages of the same size are transferred in each communication step. Therefore, the total data transfer cost is minimized.

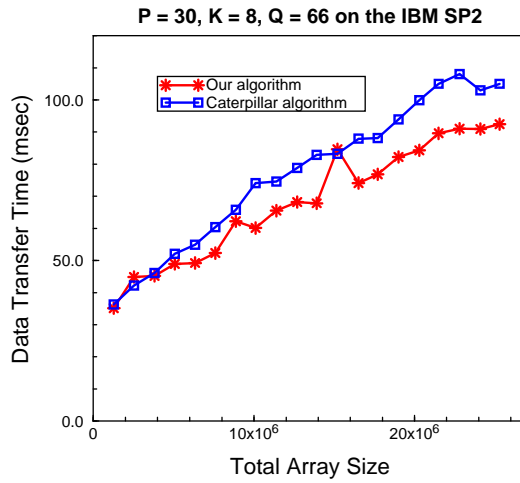
The schedule and index computation costs are also important in performing run-time redistribution. In our algorithm, the schedule and the index sets are computed in $O(\max(P, Q))$ time. The schedule and index computation using our algorithm is significantly faster compared with the bipartite matching scheme in [5]. Our schedule and index computation times are small enough to be negligible compared with the data transfer time, making our algorithms suitable for run-time data redistribution.

Acknowledgment

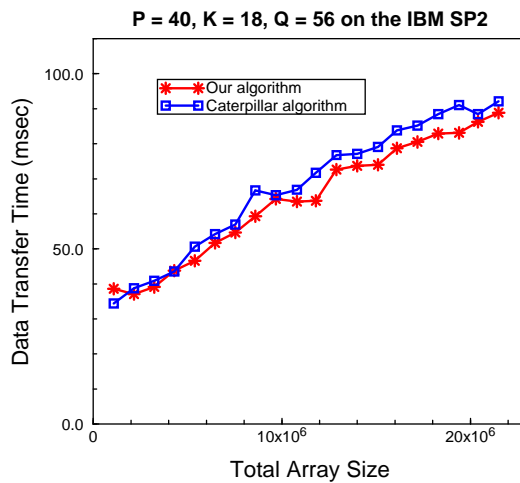
We would like to thank the staff at MHPCC for their assistance in evaluating our algorithms on the IBM SP-2. We also would like to thank Manash Kirtania for his assistance in preparing this manuscript.



(a) $\mathfrak{R}_{16}(18,8,78)$



(b) $\mathfrak{R}_{16}(30,8,66)$



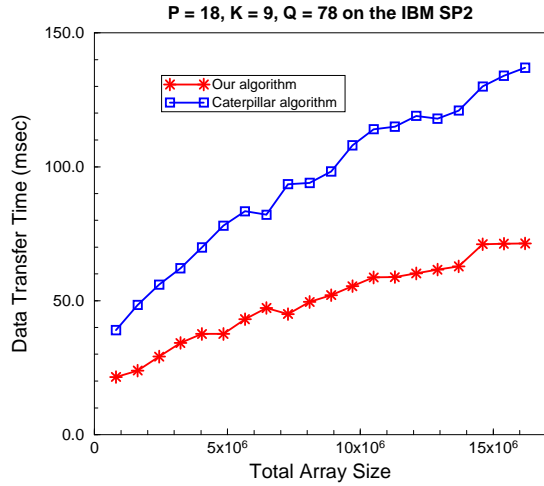
(c) $\mathfrak{R}_8(40,18,56)$

Figure 16: Data transfer time for all-to-all communication examples with different message sizes.

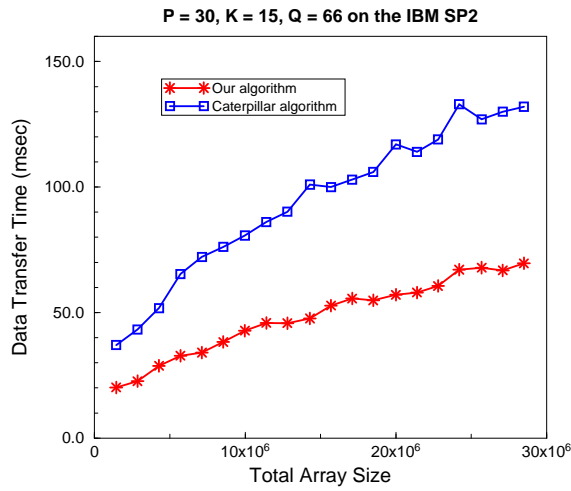
Figure 16 reports the experimental results for examples of all-to-all communication with different message sizes. The data transfer time in the all-to-all communication case is sensitive to network contention since every source processor communicates with every destination processor. For $\mathfrak{R}_{16}(18, 8, 78)$, both algorithms have the same number of steps (78). Within a superblock, a third of the messages are two blocks in size and two thirds are one block in size. The Caterpillar algorithm does not attempt to send equal-sized messages in each communication step. Therefore, the data transfer time for a step is determined by the time to transfer the largest message in that step. Theoretically, the data transfer time of our algorithm is reduced by 33% when compared with that of the Caterpillar algorithm. In our experiments with large message sizes, we observed up to 19.82% reduction. With small messages, both algorithms have approximately the same performance since the start-up time dominates the data transfer time. Similar experimental results are shown in Figure 16(b) and (c) for two other scenarios.

5.3 Schedule computation time

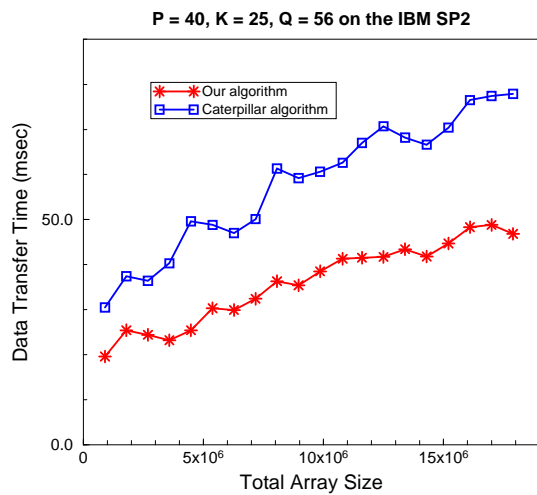
The time for computing the schedule in the Caterpillar algorithm as well as in our algorithm is negligible compared with the total redistribution time. Even though the schedule in the Caterpillar algorithm is simpler than ours, the Caterpillar algorithm needs time for index computation to identify the blocks to be packed in a communication step. This time is approximately the same as our schedule computation time. The schedule computation time of the bipartite matching scheme [5] is much higher than that of the Caterpillar algorithm and that of our algorithm. It can be a significant fraction of the data transfer time. To compute the schedule of the bipartite matching scheme in our experiments, we used a unit weight bipartite matching code down-loaded from [20]. In Table 6, the data transfer time in 3 non all-to-all communication examples is compared with the schedule computation time of our algorithm and that of the bipartite matching scheme. The schedule computation time of bipartite matching scheme can be observed to be as much as 50% of the data transfer time. On the other hand, the schedule computation time of our algorithm is less than 1% of the data transfer time. This makes our algorithm attractive for run-time data redistribution.



(a) $\mathfrak{R}_{16}(18,9,78)$



(b) $\mathfrak{R}_{16}(30,15,66)$



(c) $\mathfrak{R}_{16}(40,25,56)$

Figure 15: Data transfer time for non all-to-all communication examples.

```

for (j=0; j<n1; j++) {
  /* redistribution routine */
  compute schedule and index set
  node_tr[j] = 0;
  for (i=0; i<n2; i++) {
    if (source processor) { /* source processor */
      pack message
      ts = MPI_Wtime()
      send message to a destination processor
      node_tr[j] = tr[j] + MPI_Wtime() - ts
    } else { /* destination processor */
      ts = MPI_Wtime()
      receive message from a source processor
      node_tr[j] = tr[j] + MPI_Wtime() - ts
      unpack message
    }
  }
  compute tavg from node_tr of each node
  Tr[j] = tavg
}

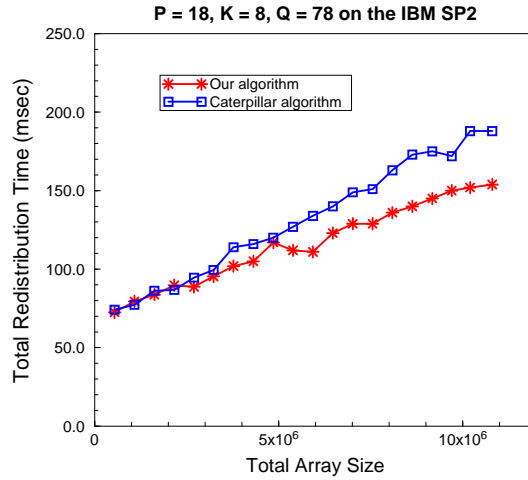
compute Tmax = max{Tr[j]}, Tmin = min{Tr[j]},
             Tmed = median{Tr[j]}, Tavg = average{Tr[j]}

```

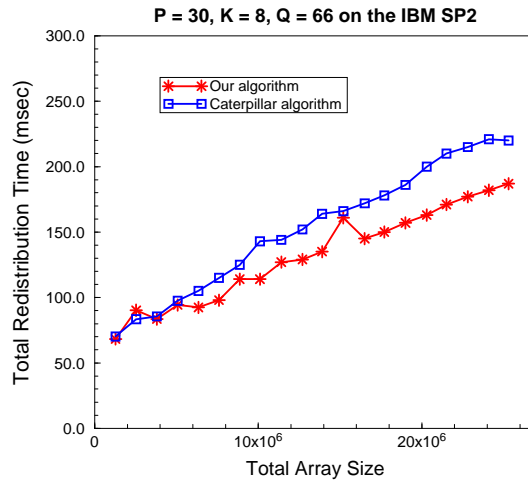
Figure 14: Steps in measuring the data transfer time.

In this subsection, we report the experimental results of the data transfer time of our algorithm and the Caterpillar algorithm. In this section, the bipartite matching scheme is not considered, since it has the same complexity in the non all-to-all communication case compared with our algorithm and does not consider the all-to-all communication case with different message sizes. The experiments were performed in the same manner as discussed in Subsection 5.1. The data sets used in these experiments are the same as those used in the previous subsection. The data transfer time of each communication step is measured first. Then the total data transfer time is computed by summing up the measured time for all the communication steps. The methodology for measuring the time is shown in Figure 14.

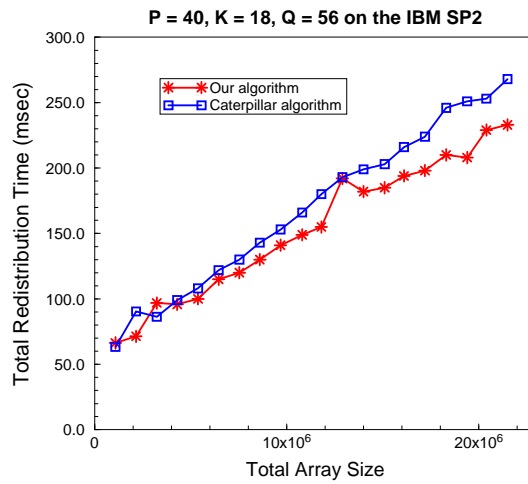
The redistribution $\mathfrak{R}_{18}(18, 9, 78)$ is an example of the non all-to-all communication case. The messages in each communication step are of the same size. The total number of communication steps is 36 using our algorithm, where as the total number of steps is 78 using the Caterpillar algorithm. Therefore, the data transfer time of our algorithm is theoretically 50% of that of the Caterpillar algorithm. In the experimental results (see Figure 15(a)), the data transfer time of our algorithm is between 48.13% and 57.61% of that of the Caterpillar algorithm. Figure 15 (b) and (c) show the experimental results for two other non all-to-all communication examples. Similar reductions in time were observed in these experiments.



(a) $\mathfrak{R}_{16}(18,8,78)$



(b) $\mathfrak{R}_{16}(30,8,66)$



(c) $\mathfrak{R}_8(40,18,56)$

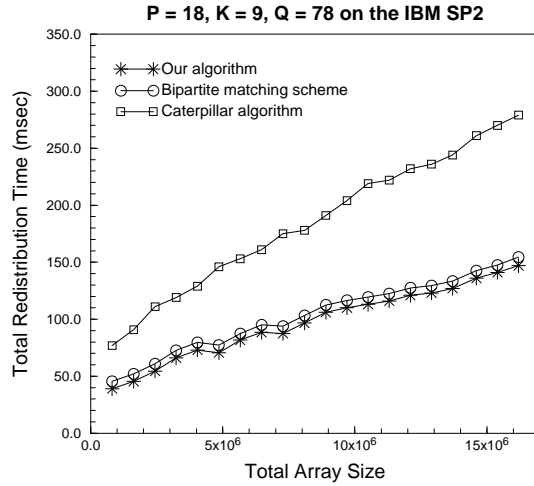
Figure 13: Total redistribution time for all-to-all communication examples with different message sizes.

the Caterpillar algorithm. Therefore, the redistribution time of our algorithm is theoretically 50% of that of the Caterpillar algorithm. In the experimental results shown in Figure 12(a), the redistribution time of our algorithm is between 48.42% and 56.67% of that of the Caterpillar algorithm. In comparison with the bipartite matching scheme, our algorithm has faster schedule computation. In experiments with small arrays, the schedule computation time of the bipartite matching scheme can be observed up to 17% of the total redistribution time.

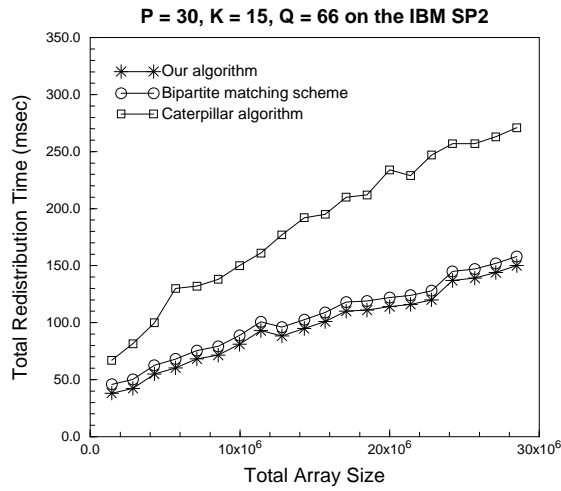
Figure 12 (b) shows experimental results for $P = 30$, $Q = 66$ and $K = 15$ and (c) shows experimental results for $P = 40$, $Q = 56$ and $K = 25$, respectively. For these two instances, the number of communication steps using our algorithm is 33 and 35, respectively. The number of communication steps using the Caterpillar algorithm is 66 and 56, respectively. Therefore, the redistribution time of our algorithm can be expected to be reduced by 50% and 37% when compared with that of the Caterpillar algorithm. Our experimental results confirm these. The schedule computation time of the bipartite matching scheme is around 10 *msecs* in these instances.

Figure 13 compares the overall redistribution time for the all-to-all communication case with different message sizes. Figure 13(a) reports the experimental results for $\mathfrak{R}_{16}(18, 8, 78)$. The array size was varied from 539,136 points (2.16 Mbytes) to 10,782,720 points (43.13 Mbytes). For this case, both the algorithms have the same number of steps (78). Within a superblock, a third of the messages are two blocks in size while the others are one block in size. In our algorithm, equal-sized messages are transferred from any node in each communication step. Therefore, during a third of the communication steps, messages of size two blocks are sent while messages of size one block are sent during other two thirds of the communication steps. The Caterpillar algorithm does not attempt to schedule the communication operations so that equal-sized messages are sent in each communication step. Therefore, the redistribution time in a step is determined by the time to transfer the largest message in that step. Theoretically, the total redistribution time of our algorithm is reduced by 33.3% compared with that of the Caterpillar algorithm. In our experiments, we observed up to 19.15% reduction in redistribution time. When the array size is small, both algorithms have approximately the same performance since the start-up cost dominates the overall data transfer time. As the array size increases, the reduction in the time to perform the distribution using our algorithm improves. For two other scenarios, we obtained similar results (See Figure 13(b) and (c)).

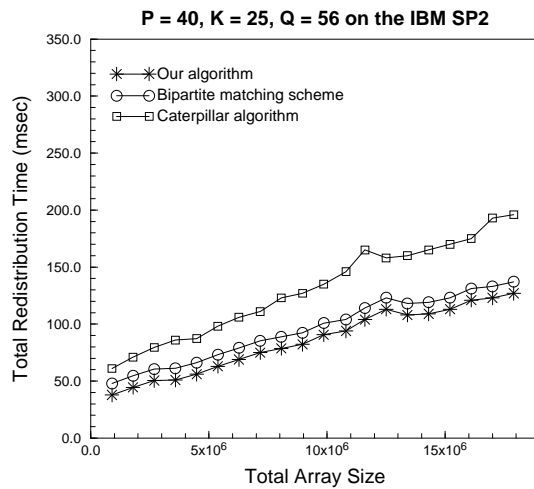
5.2 Data transfer time



(a) $\mathfrak{R}_{16}(18,9,78)$



(b) $\mathfrak{R}_{16}(30,15,66)$



(c) $\mathfrak{R}_{16}(40,25,56)$

Figure 12: Total redistribution time for non all-to-all communication examples.

the bipartite matching scheme. The bipartite matching scheme and our scheme optimize the number of communication steps and minimize the data transfer cost in non all-to-all communication. They will differ only in the schedule computation time. Thus, we can obtain the total redistribution time of the bipartite matching scheme by simply adding the schedule computation time of the bipartite matching scheme to the total redistribution time instead of our schedule computation time. In our experiments, the source and the destination processor sets were disjoint. In each communication step, each sender packs a message before sending it and each receiver unpacks the message after receiving it. Pack operations in the source processors and unpack operations in the destination processors were overlapped, *i.e.*, after sending their message in communication step i , senders start to pack a message for communication in step $(i + 1)$ and receivers start to unpack the message received in step i .

Our methodology for measuring the total redistribution time is shown in Figure 11. The time was measured using the `MPI_Wtime()` call. `n1` is the number of runs. A run is an execution of redistribution. `n2` is the number of communication steps. Each processor measures `node_time[j]` in the j^{th} run. Generally, source and destination processors which do not perform an interprocessor communication in the last step, complete the redistribution earlier than the processors which receive a message and unpack it. A barrier synchronization, `MPI_Barrier()`, was performed at the end of the redistribution. After measuring `node_time`, the average `node_time` over $(P + Q)$ processors is computed and saved as `tavg`. The measured value is stored in an array `T`, as shown in Figure 11. After the redistribution is performed `n1` times, the maximum, minimum, median, and average total redistribution time are computed over `n1` runs. In our experiments, `n1` was set to 20. Over `n1` runs, a variation in the measured redistribution times was observed. Since the minimum time has the smallest component due to OS interference and other effects related to the environment, it provides a more accurate observation of the redistribution time than others. So, we use `Tmin` only in our experimental results.

Figure 12 shows experimental results for the non all-to-all communication case. In these experiments, 96 nodes were used. Figure 12 shows results of the redistribution $\mathfrak{R}_{16}(18, 9, 78)$, where 18 source processors and 78 destination processors were used and K was set to 9. The total number of array elements (in single precision) was varied from 808,704 (3.23 Mbytes) to 14,174,080 (64.7 Mbytes). In the non all-to-all communication case, the messages in each communication step are of the same size. The total number of communication steps is 39 in our algorithm, while it is 78 in

```

for (j=0; j<n1; j++) {
  ts = MPI_Wtime()
  /* redistribution routine */
  compute schedule and index set
  for (i=0; i<n2; i++) {
    if (source processor) { /* source processor */
      pack message
      send message to a destination processor
    } else { /* destination processor */
      receive message from a source processor
      unpack message
    }
  }
  node_time[j] = MPI_Wtime() - ts
  compute tavg from node_time of each node
  T[j] = tavg
}

compute Tmax = max{T[j]}, Tmin = min{T[j]},
              Tmed = median{T[j]}, Tavg = average{T[j]}

```

Figure 11: Steps for measuring the redistribution time.

cost in the bipartite matching scheme is much greater even for the problem considered here.

To evaluate the total redistribution cost and the data transfer cost, we consider 3 different scenarios corresponding to the relative size of P and Q : (Scenario 1); $P \ll Q$, (Scenario 2); $Q \geq 2P$, and (Scenario 3); $P < Q < 2P$. In our experiments, we choose $P = 18$ and $Q = 78$ for Scenario 1, $P = 30$ and $Q = 66$ for Scenario 2, and $P = 46$ and $Q = 50$ for Scenario 3. The array consists of single precision integers. The size of each array element is 4 bytes. The array size was chosen to be a multiple of the size of a superblock to avoid padding using dummy data.

The rest of this section is organized as follows. Subsection 5.1 reports experimental results of the overall redistribution time of our algorithm, the Caterpillar algorithm, and bipartite matching scheme. Subsection 5.2 shows experimental results for the data transfer time of our algorithm and the Caterpillar algorithm. Subsection 5.3 compares our algorithm and the bipartite matching scheme with respect to the schedule computation time.

5.1 Total redistribution time

In this subsection, we report experimental results for the total redistribution time of our algorithm, the Caterpillar algorithm, and the bipartite matching scheme. The total redistribution time consists of the schedule computation time, index computation time, packing/unpacking time, and data transfer time. The experimental results of bipartite matching scheme are reported only in the case of non all-to-all communication, since all-to-all communication instances are not considered in

Table 5: Comparison of data transfer cost and schedule and index computation costs of the Caterpillar algorithm, bipartite matching scheme and our algorithm.

	Non all-to-all communication		All-to-all communication with different message sizes	
	Data transfer cost	Schedule and index computation cost	Data transfer cost	Schedule and index computation cost
Caterpillar algorithm [19]	$Q\left(T_s + \tau_d \frac{M}{L_s}\right)$	$O(Q)$	$QT_s + \tau_d \sum_{i=0}^{Q-1} m_i$	$O(Q)$
Bipartite matching scheme [4]	$L_s\left(T_s + \tau_d \frac{M}{L_s}\right)$	$O((P+Q)^{3.5})$	N/A	N/A
Our algorithm	$L_s\left(T_s + \tau_d \frac{M}{L_s}\right)$	$O(Q)$	$QT_s + \tau_d M$	$O(Q)$

Note: $L_s < Q$ for the non all-to-all communication case, $M = N/P$ and m_i is the maximum transferred data size in communication step i .

putation costs. For the all-to-all communication case with equal-sized messages, the data transfer cost is the same in each communication step for all three algorithms. Also, the schedule computation can be performed in a simple way. Hence, the all-to-all communication case with equal-sized messages is not considered in Table 5. In Table 5, M is the size of the array assigned to each source processor ($M = \frac{N}{P}$). For the non all-to-all communication case, $L_s < Q$, where $L_s = \frac{lcm(P,KQ)}{P}$. Our algorithm as well as the bipartite matching scheme perform fewer communication steps than the Caterpillar algorithm. For the all-to-all communication case with different message sizes, the messages transmitted in a communication step are of the same size in our algorithm. Therefore, the network bandwidth is fully utilized and the total transmission cost is $\tau_d M$. However, the bipartite matching scheme does not consider this case. In the Caterpillar algorithm, the transmission cost in a communication step is dominated by the largest message transferred in that step. Let m_i denote the size of the largest message sent in a communication step i . Note that $\sum_{i=0}^{Q-1} m_i \geq M$. The total start-up cost of the Caterpillar and our algorithm is QT_s since the number of communication steps is the same. On the other hand, the total transmission cost of our algorithm is $\tau_d M$ which is less than that of the Caterpillar algorithm. The Caterpillar algorithm as well as our algorithm perform the schedule and index computation in $O(Q)$ time. However, the schedule and index computation

the theorem in terms of the number of reorganizations required to convert \mathbf{C}_{send} to its desired final form. The total start-up cost is proportional to the number of communication steps. Therefore, it is in proportion to the number of row reorganizations. Also, communication pairs in each communication step communicate messages of the same size. Therefore, source processors have the same total transmission cost without wasting the network bandwidth. The total transmission cost is proportional to the size of an array assigned to each processor.

(i) non all-to-all communication: We know that each row of \mathbf{C}_{send} represents communication step i . The total number of communication steps is L_s . In each step, processor pairs communicate one block per superblock between them. Therefore, data transfer cost can be estimated as follows,

$$\begin{aligned} \text{data transfer cost} &= L_s \cdot \left(T_s + \frac{N}{PL_s} \tau_d \right) \\ &= L_s T_s + \frac{N}{P} \tau_d \end{aligned}$$

(ii) all-to-all communication: From Theorem 3, we need only Q rows in the generalized circulant matrix form to develop the communication schedule table \mathbf{C}_{send} . We know that each row of \mathbf{C}_{send} represents communication step i . The total number of communication steps is L_s . In each step, equal-sized messages are communicated between processor pairs. The message size in the first r rows on \mathbf{C}_{send} is $q + 1$ blocks per superblock, where $r = L_s \bmod Q$ and $q = L_s/Q$. The message size in remaining rows is q blocks per superblock. Therefore, data transfer cost can be estimated as follows,

$$\begin{aligned} \text{data transfer cost} &= QT_s + \{r(q + 1) + (Q - r)q\} \frac{N}{PL_s} \tau_d \\ &= QT_s + L_s \frac{N}{PL_s} \tau_d \\ &= QT_s + \frac{N}{P} \tau_d \end{aligned}$$

□

5 Experimental Results

Our experiments were conducted on the IBM SP2 at the Maui High Performance Computing Center. The algorithms were written in C and MPI. MPI communication primitives were used for interprocessor communication.

Table 5 shows a comparison of the proposed algorithm with the Caterpillar algorithm[18] and the bipartite matching scheme[5] with respect to the data transfer cost and schedule and index com-

same message size, K_1 is a multiple of G_2 . Therefore, for all-to-all communication case, the send communication schedule table is a $Q \times P$ matrix which consists of the first $\{0, 1, \dots, Q - 1\}^{th}$ rows in a $L_s \times P$ generalized circulant matrix. Every G_2^{th} row in the circulant block matrix is folded into its first row. Therefore, equal-sized messages are transferred in a communication step. \square

Corollary 2 *For the redistribution problem $\mathfrak{R}_x(P, K, Q)$, the proposed algorithm minimizes the data transfer cost in both non all-to-all communication case and all-to-all communication case.*

4.4 Data transfer cost

In distributed memory model, the data transfer cost has two parameters; start-up time and transmission time. The start-up time, T_s , is incurred once for each communication event and is independent of the communicated message size. Generally, the start-up time consists of the transfer request and acknowledgment latencies, context switch latency, and latencies for initializing the message header. The unit transmission time, τ_d , is the cost of transferring a message of unit length over the network. The transmission time for a message is proportional to its size. Thus, the data transfer time for sending a message of size m units from one processor to another is modeled as $T_s + m\tau_d$. In this model, a reorganization of the data elements among the processors, in which each processor has m units of data for another processor, also takes $T_s + m\tau_d$ time. This model assumes that there is no node contention. This is ensured by our communication schedules for redistribution. Theorem 4 shows the data transfer cost for our redistribution algorithms using the distributed memory model.

Theorem 4 *Consider an array with N elements initially distributed cyclic(x) on P processors. The array needs to be redistributed to cyclic(Kx) on Q processors. Using our algorithms, the data transfer costs for performing $\mathfrak{R}_x(P, K, Q)$ are (i) $L_s T_s + \frac{N}{P} \tau_d$ in the case of a non all-to-all communication pattern, and (ii) $Q T_s + \frac{N}{P} \tau_d$ in the case of an all-to-all communication pattern.*

Proof: Data transfer cost is considered as a sum of total start-up cost and total transmission cost. Consider our send communication schedule table \mathbf{C}_{send} , which is in generalized circulant matrix form and each row consists of distinct elements. The communication schedule is specified as a sequence of row reorganization on \mathbf{C}_{send} . A row reorganization moves elements to their destination processor specified by \mathbf{C}_{send} . This corresponds to an interprocessor communication event. The number of communication steps is equal to the number of row reorganizations on \mathbf{C}_{send} . We prove

in a $L_s \times P$ generalized circulant matrix as computed by the equations of Theorem 2. Equal-sized messages are transferred in each communication step.

Proof: From Theorem 1, the *dpt* \mathbf{T} is rearranged into a generalized circulant matrix form by column reorganizations. It is computed by Eq. (17) in Theorem 2. The generalized circulant matrix is a $K_1 \times P_1$ circulant block matrix. Each block matrix is a $Q_1 \times G_1$ circulant matrix. We show that every G_2^{th} row in the circulant block matrix consists of the same indices in each column. Therefore, we prove that the send communication schedule table is represented as $Q \times P$ generalized circulant matrix and equal-sized messages are transferred in each communication step.

Consider a generalized circulant matrix computed by Eq. (17) in Theorem 2. The second part of Eq. (17) is related to an offset of the $Q_1 \times G_1$ block matrix. For processor j , the second part is $qP_1 \bmod Q$, where $q = (i_2 - j_2) \bmod Q_1$ and $0 \leq q, i_2 < Q_1$. Since $qP_1 \bmod Q = qP_1 - \{(qP_1/Q) \times Q\} = G_2\{qP_2 - Q_1(qP_1/Q)\}$, $qP_1 \bmod Q$ is a multiple of G_2 . Therefore, the second part of Eq. (17) is an element in *set* $\{0, G_2, 2G_2, \dots, (Q_1 - 1)G_2\}$. Similarly, the first part of Eq. (17) is related to the base value of each $Q_1 \times G_1$ block matrix. A base value is determined by $\{\{n(j_1 - i_1)\} \bmod P_1\} \bmod Q = \{n(j_1 - i_1)\} \bmod P_1$, for processor j and $0 \leq i_1 < K_1$. For processor j , consider the base value of the block matrix at the first row and the block matrix at the G_2^{th} row. For the block matrix at the first row, the base value in Eq. (17) is computed as follows,

$$\begin{aligned} (nj_1) \bmod P_1 &= nj_1 - \{(nj_1)/P_1\}P_1 \\ &= nj_1 - \lambda_1 G_2 \end{aligned}$$

where $\lambda_1 = [(nj_1)/P_1]P_2$ and $P_1 = P_2G_2$. For the block matrix in the G_2^{th} row, the base value in Eq. (17) is also computed as follows,

$$\begin{aligned} n(j_1 - G_2) \bmod P_1 &= n(j_1 - G_2) - \{n(j_1 - G_2)/P_1\}P_1 \\ &= (nj_1) - \{n - [n(j_1 - G_2)/P_1]P_2\}G_2 \\ &= (nj_1) - \lambda_2 G_2 \end{aligned}$$

where $\lambda_2 = n - [n(j_1 - G_1)/P_1]P_2$. Even though their base values are different, the difference of both base values are a multiple of G_2 . When the base value is divided by G_2 , their remainder is the same. Therefore, the $Q_1 \times G_1$ block matrices with these base values consist of the same destination processor indices, because offset values are multiples of G_2 . So, block matrices at every G_2^{th} row in each column consist of the same entries in \mathbf{C}_{send} . In the case of all-to-all communication with

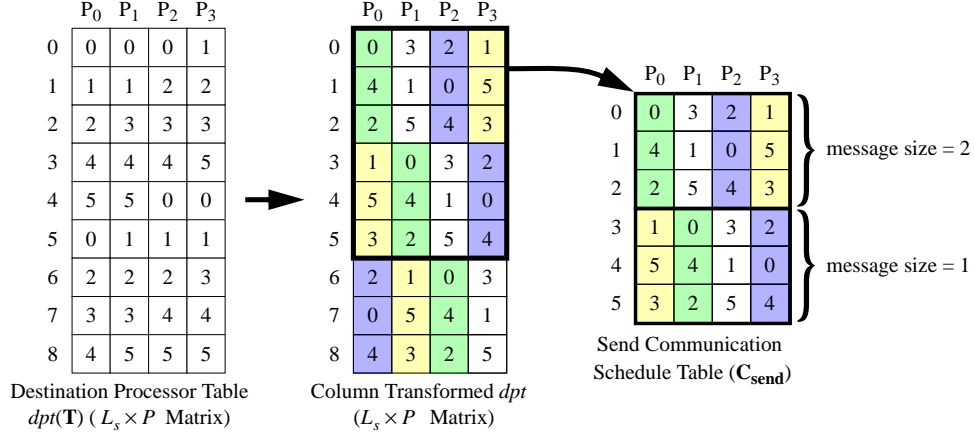


Figure 10: Example illustrating an all-to-all case with different message sizes: $\mathfrak{R}_x(4, 3, 6)$

the same entries but different circular-shifted patterns. These blocks can be folded onto the blocks in their first row. Therefore, the first G_2 rows in the block matrix only are used in determining a send communication schedule table \mathbf{C}_{send} . It is a $Q \times P$ generalized circulant matrix. Since blocks in every G_2^{th} row are folded onto blocks in their first row, for all-to-all communication case with different message sizes, blocks in the first $(K_1 \bmod G_2)$ rows of \mathbf{C}_{send} have size $\lceil \frac{K_1}{G_2} \rceil$, while blocks in the remaining rows have size $\lfloor \frac{K_1}{G_2} \rfloor$.

Figure 10 shows an example of the send communication schedule table of $\mathfrak{R}_x(4, 3, 6)$, generated for all-to-all case with different message sizes. In this example, each processor has more entries than 6 destination processors. The corresponding dpt is a $L_s \times P$ matrix, where $L_s = 9$ and $P = 4$. Applying column reorganizations on it results in a generalized circulant matrix, which can be considered as a $K_1 \times P_1$ block matrix, where $K_1 = 3$ and $P_1 = 4$. Each block is a $Q_1 \times G_1$ matrix, where $Q_1 = 3$ and $G_1 = 1$. The first $G_2 = 2$ rows are used as a send communication schedule table \mathbf{C}_{send} . The 3rd row is folded onto the 1st row. Hence, the message size in the 1st row is 2 and that in the 2nd row is 1. If K_1 is a multiple of G_2 , the message size in every step is the same. Therefore, the network bandwidth is fully utilized by sending equal sized messages in each communication step.

Theorem 3 summarizes the above intuition and shows that the send communication schedule table for all-to-all communication can be obtained as a generalized circulant matrix. Further, it ensures that equal-sized messages are transferred in every communication step.

Theorem 3 *If a redistribution problem $\mathfrak{R}_x(P, K, Q)$ requires all-to-all communication, the send communication schedule table is a $Q \times P$ matrix which consists of the first $\{0, 1, \dots, Q - 1\}^{\text{th}}$ rows*

From Eq. (22) and Eq. (23),

$$X = \{n(j_1 - i_1)\} \bmod P_1 \quad (24)$$

$$Y = \{m(j_1 - i_1)\} \bmod K_1 \quad (25)$$

Eq. (21) becomes $a'_2 = (i_2 - j_2) \bmod Q_1$. By replacing a'_1 and a'_2 with i_1 and i_2 , Eq. (19) can be rewritten as follows

$$\mathbf{D}'_{\text{init}}(i, j) = (XK_1 + i_1)G_1 + \{(i_2 - j_2) \bmod Q_1\}K_1P + j_2$$

Therefore, $\mathbf{C}_{\text{send}}(i, j) = (\mathbf{D}'_{\text{init}}(i, j)/K) \bmod Q$ gives Eq. (17) since $i_1G_1 + j_1 < K$. Similarly, we can prove Eq. (18). \square

The above formulae for computing the communication schedule and index set for redistribution are extremely efficient compared with the methods presented in [5], which use a bipartite matching algorithm. Furthermore, using our formulae, each processor computes only entries which it needs in its send communication schedule table. Hence, the schedule and index set computation can be performed in a distributed way and the total cost of computing the schedule and index set is $O(\max(P, Q))$. Our scheme minimizes the number of communication steps and avoids node contention. In each communication step, equal-sized messages are transferred. Therefore, our scheme minimizes the data transfer cost.

Corollary 1 *In the proposed redistribution algorithm, the costs of index set and communication schedule computations are $O(\max(P, Q))$. The amortized cost to compute a step in the communication schedule and index set computation is $O(1)$.*

4.3 All-to-all communication

The all-to-all communication case arises if $G(= G_1G_2) \leq K$ as stated in Lemma 1, where $G_1 = \gcd(P, K)$ and $G_2 = \gcd(P_1, Q)$. From the first superblock, the *dpt* \mathbf{T} is constructed. The *dpt* \mathbf{T} is a $L_s \times P$ matrix, where $L_s = K_1Q_1$. Since $Q = Q_1G_2$ and $G_2 \leq K_1$, the number of rows in *dpt* \mathbf{T} $L_s \geq Q$. Therefore, each column has more entries than Q destination processors. In each column, several blocks are transferred to the same destination. The column reorganizations as stated in Section 4.2 are applied to the *dpt* \mathbf{T} , which results in a generalized circulant matrix and is a $K_1 \times P_1$ circulant block matrix. Each block is a $Q_1 \times G_1$ submatrix which is also a circulant matrix. In the block matrix, the first G_2 blocks in each column are distinct. Blocks in every G_2^{th} row have

Theorem 2 gives the formulae to compute the individual entries of \mathbf{C}_{send} and \mathbf{D}_{send} efficiently. Referring to Figure 9, in the communication schedule table, \mathbf{C}_{send} , i indices represent the communication steps and j indices represent the source processor indices. Each entry $\mathbf{C}_{\text{send}}(i, j)$ refers to the destination processor to which the j^{th} processor communicates in the i^{th} communication step. In the following theorem, i_1, i_2, j_1 , and j_2 refer to the indices defined in the proof of Theorem 1. (See Eq. (11) and Eq. (13))

Theorem 2 *Given redistribution parameters, K, P , and Q , let $G_1 = \gcd(P, K)$, $P_1 = P/G_1$, $K_1 = K/G_1$, $G_2 = \gcd(P_1, Q)$, and $Q_1 = Q/G_2$. A send communication schedule table \mathbf{C}_{send} and the send data location table \mathbf{D}_{send} in the generalized circulant matrix form can be computed as follows:*

$$\mathbf{C}_{\text{send}}(i, j) = \{ \{n(j_1 - i_1)\} \bmod P_1 + \{(i_2 - j_2) \bmod Q_1\} P_1 \} \bmod Q \quad (17)$$

$$\mathbf{D}_{\text{send}}(i, j) = \{m(j_1 - i_1)\} \bmod K_1 + \{(i_2 - j_2) \bmod Q_1\} K_1 \quad (18)$$

where n and m are solutions of $nK_1 - mP_1 = 1$.

Proof: From the proof of Theorem 1, the $(i, j)^{\text{th}}$ entry of $\mathbf{D}'_{\text{init}}$ is shown as follows.

$$\mathbf{D}'_{\text{init}}(i, j) = \mathbf{D}_{\text{init}1}(a', j) = (a'_1 P_1 + j_1) G_1 + (a'_2 K_1 P + j_2) \quad (19)$$

From Eq. (11) and Eq. (13),

$$i_1 = (a'_1 P_1 + j_1) \bmod K_1 \quad (20)$$

$$i_2 = (a'_2 + j_2) \bmod Q_1 \quad (21)$$

Let $t = (a'_1 P_1 + j_1)$. From Eq. (20),

$$t = X K_1 + i_1 = Y P_1 + j_1$$

where $0 \leq X < P_1$ and $0 \leq Y < P_1$. Hence,

$$X K_1 - Y P_1 = (j_1 - i_1) \quad (22)$$

We have to solve a Diophantine equation to find X and Y . Since $\gcd(K_1, P_1) = 1$, we can find m and n using the Euclid algorithm such that

$$nK_1 - mP_1 = 1 \quad (23)$$

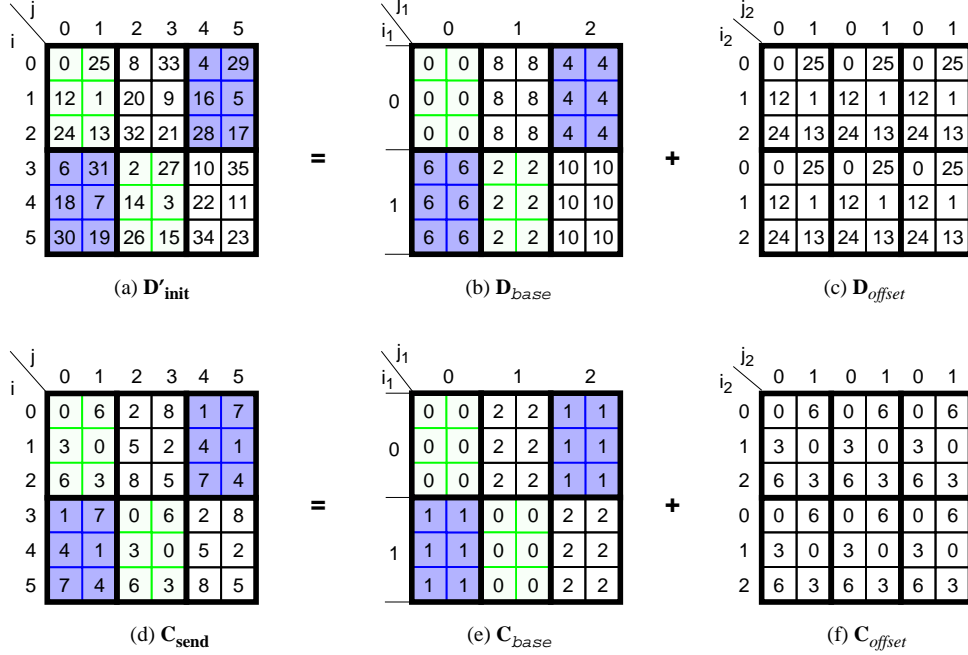


Figure 9: Decomposition of D'_{init} and C_{send}

and offset respectively. The row and column of the base are indexed as i_1 and j_1 respectively. Also, the row and column of the offset are indexed as i_2 and j_2 respectively. Thus, entries of the base matrix is independent of i_2 and j_2 . All entries within $Q_1 \times G_1$ submatrix of C_{base} have the same value, given by $(\frac{(a'_1 P_1 + j_1)}{K_1}) \bmod Q = \frac{(a'_1 P_1 + j_1)}{K_1}$. Similarly, the entries of the offset are independent of i_1 and j_1 . Therefore, all $Q_1 \times G_1$ submatrices of C_{offset} are identical to one another. Each is a $Q_1 \times G_1$ circulant matrix. Figure 9 shows the base and the offset of D'_{init} and C_{send} for $\mathfrak{R}_x(6, 4, 9)$.

With reference to Figure 5, observe that C_{send} helps to specify the row reorganizations that convert D'_{init} to D'_{final} . The initial column reorganizations convert D_{init} to D'_{init} and T to C_{send} . Although Section 3 indicates that the column reorganizations reorganize data within a processor, this operation can be expensive for large array sizes. Instead, the reorganization can be done by maintaining pointers to the elements of the array. Each source processor has a table which points to the data blocks to be packed in a communication step. It is denoted as *send data location table* D_{send} . Each entry of D_{send} is the local block index of the corresponding entry of D'_{init} . Therefore, $D_{send}(i, j) = D'_{init}(i, j)/P$. Each entry of C_{send} , $C_{send}(i, j)$, points to the destination processor of the corresponding entry of D_{send} , $D_{send}(i, j)$. Our scheme computes the schedule and data index set at the same time.

The *receive communication schedule table* C_{recv} and the *receive data location table* D_{recv} can be computed in a similar way. They are not discussed further.

Similarly, we can convert each submatrix to a circulant matrix. Elements within the j_2^{th} column of submatrix $\mathbf{B}_{a'_1, j_1}$ are circularly shifted by j_2 positions. The relationship between a'_2 and i_2 is as follows.

$$i_2 = (a'_2 + j_2) \bmod Q_1 \quad (13)$$

Each reorganized submatrix is shown below.

$$\mathbf{C}_{i_1, j_1} = \left(\frac{(a'_1 P_1 + j_1)}{K_1} + \begin{bmatrix} 0 & (Q_1 - 1)P_1 & \cdots & (Q_1 - G_1 + 1)P_1 \\ P_1 & 0 & \cdots & (Q_1 - G_1 + 2)P_1 \\ \vdots & \vdots & & \vdots \\ (Q_1 - 1)P_1 & (Q_1 - 2)P_1 & \cdots & (Q_1 - G_1)P_1 \end{bmatrix} \right) \bmod Q \quad (14)$$

The resulting matrix is a block-wise circulant matrix and its submatrices are also circulant matrices. Therefore, it is a generalized circulant matrix. This matrix can be used as a send communication schedule table \mathbf{C}_{send} . The same column reorganizations are applied to the $\mathbf{D}_{\text{init1}}$. The reorganized distribution table, $\mathbf{D}'_{\text{init}}$ is in the following form:

$$\mathbf{D}'_{\text{init}} = \begin{bmatrix} \mathbf{D}_{0,0} & \mathbf{D}_{0,1} & \cdots & \mathbf{D}_{0,P_1-1} \\ \mathbf{D}_{1,0} & \mathbf{D}_{1,1} & \cdots & \mathbf{D}_{1,P_1-1} \\ \vdots & \vdots & & \vdots \\ \mathbf{D}_{K_1-1,0} & \mathbf{D}_{K_1-1,1} & \cdots & \mathbf{D}_{K_1-1,P_1-1} \end{bmatrix} \quad (15)$$

$$\mathbf{D}_{i_1, j_1} = (a'_1 P_1 + j_1)G_1 + \begin{bmatrix} 0 & (Q_1 - 1)K_1 P + 1 & \cdots & (Q_1 - G_1 + 1)K_1 P + 1 \\ K_1 P & 1 & \cdots & (Q_1 - G_1 + 2)K_1 P + 1 \\ \vdots & \vdots & & \vdots \\ (Q_1 - 1)K_1 P & (Q_1 - 2)K_1 P + 1 & \cdots & (Q_1 - G_1)K_1 P + G_1 - 1 \end{bmatrix} \quad (16)$$

We will show that every row of \mathbf{C}_{send} has P distinct numbers from the set of $\{0, 1, 2, \dots, Q - 1\}$. Consider the 0^{th} row of \mathbf{C}_{send} . For row $i_1 = 0$, it implies $(a'_1 P_1 + j_1) \bmod K_1 = 0$. Then, $(a'_1 P_1 + j_1) \in \{0, K_1, \dots, (P_1 - 1)K_1\}$ since $0 \leq a'_1 < K_1$ and $0 \leq j_1 < P_1$. Therefore, $\frac{(a'_1 P_1 + j_1)}{K_1} \in \{0, 1, \dots, (P_1 - 1)\}$. In the first row, the second term of Eq. (14) is $([0 \ (Q_1 - 1)P_1 \ \cdots \ (Q_1 - G_1 + 1)P_1]) \bmod Q$. Since $Q_1 = \frac{Q}{\gcd(Q, P_1)}$, $Q_1 P_1 \bmod Q = 0$. Therefore, the second term is $[0 \ Q - P_1 \ \cdots \ Q - (G_1 - 1)P_1]$. By summing the base and offset values, elements in the first row will have a value in range 0 to $P_1 - 1$ or in range $Q - (G_1 - 1)P_1$ to $Q - 1$. Since $P_1 - 1 < Q - (G_1 - 1)P_1$, row 0 of the send communication table \mathbf{C}_{send} consists of P distinct destination processor indices. Since \mathbf{C}_{send} is a generalized circulant matrix, every row of \mathbf{C}_{send} consists of P distinct destination processor indices. \square

The send communication schedule table \mathbf{C}_{send} is $K_1 \times P_1$ block matrix. Each block is $Q_1 \times G_1$ submatrix. The submatrix \mathbf{B}_{i_1, j_1} of Eq. (16) consists of two parts. These can be referred as base

replacing j by $j_1G_1 + j_2$. Now, let us consider a submatrix $\mathbf{A}_{a'_1, j_1}$. Element $\mathbf{D}_{\text{init1}}(a', j)$ is in submatrix $\mathbf{A}_{a'_1, j_1}$, and is located at (a'_2, j_2) within $\mathbf{A}_{a'_1, j_1}$. This matrix can be written as,

$$\mathbf{A}_{a'_1, j_1} = (a'_1P_1 + j_1)G_1 + \begin{bmatrix} 0 & 1 & \cdots & G_1 - 1 \\ K_1P & K_1P + 1 & \cdots & K_1P + G_1 - 1 \\ \vdots & \vdots & & \vdots \\ (Q_1 - 1)K_1P & (Q_1 - 1)K_1P + 1 & \cdots & (Q_1 - 1)K_1P + G_1 - 1 \end{bmatrix} \quad (8)$$

The corresponding dpt \mathbf{T}_1 of $\mathbf{D}_{\text{init1}}$ is obtained by replacing each element by $(\mathbf{D}_{\text{init1}}(a', j)/K) \bmod Q$. This can be represented again in block matrix form as,

$$\mathbf{T}_1 = \begin{bmatrix} \mathbf{B}_{0,0} & \mathbf{B}_{0,1} & \cdots & \mathbf{B}_{0,P_1-1} \\ \mathbf{B}_{1,0} & \mathbf{B}_{1,1} & \cdots & \mathbf{B}_{1,P_1-1} \\ \vdots & \vdots & & \vdots \\ \mathbf{B}_{K_1-1,0} & \mathbf{B}_{K_1-1,1} & \cdots & \mathbf{B}_{K_1-1,P_1-1} \end{bmatrix} \quad (9)$$

The expanded form of a submatrix can be written as,

$$\mathbf{B}_{a'_1, j_1} = \left(\frac{(a'_1P_1 + j_1)}{K_1} + \begin{bmatrix} 0 & 0 & \cdots & 0 \\ P_1 & P_1 & \cdots & P_1 \\ \vdots & \vdots & & \vdots \\ (Q_1 - 1)P_1 & (Q_1 - 1)P_1 & \cdots & (Q_1 - 1)P_1 \end{bmatrix} \right) \bmod Q \quad (10)$$

Thus, we have completed the proof for stage 1. Further, we can see that each block matrix has a base value and a fixed offset matrix.

Stage 2 (\mathbf{T}_1 to \mathbf{C}_{send}): We can transform \mathbf{T}_1 to \mathbf{C}_{send} by reorganizing submatrices within columns of the block matrix and circularly shifting elements within columns of submatrices. First we will show that $Q_1 \times G_1$ block matrices can be reorganized to obtain a block-wise circulant matrix. Next, we show the elements within a submatrix can be converted to circulant matrix.

Now consider the base $b = \frac{(a'_1P_1 + j_1)}{K_1}$ of $\mathbf{B}_{a'_1, j_1}$ in Eq. (10), where $0 \leq b < P_1$. We refer to each collection of K_1 adjacent submatrices as a run in row-major order. Run r contains submatrices whose base value is r . Through these column reorganizations of the block matrix \mathbf{T}_1 , the k^{th} submatrix of a run moves to the k^{th} row in its column, where $0 \leq k < K_1$. Thus, K_1 submatrices of each run are aligned into a diagonal. In this column reorganization of the block matrix, submatrix $\mathbf{B}_{a'_1, j_1}$ in \mathbf{T}_1 moves to \mathbf{C}_{i_1, j_1} in \mathbf{C}_{send} . The relationship between a'_1 and i_1 is as follows.

$$i_1 = (a'_1P_1 + j_1) \bmod K_1 \quad (11)$$

The resultant block matrix is a block-wise circulant matrix as shown as follows.

$$\mathbf{C}_{\text{send}} = \begin{bmatrix} \mathbf{C}_{0,0} & \mathbf{C}_{0,1} & \cdots & \mathbf{C}_{0,P_1-1} \\ \mathbf{C}_{1,0} & \mathbf{C}_{1,1} & \cdots & \mathbf{C}_{1,P_1-1} \\ \vdots & \vdots & & \vdots \\ \mathbf{C}_{K_1-1,0} & \mathbf{C}_{K_1-1,1} & \cdots & \mathbf{C}_{K_1-1,P_1-1} \end{bmatrix} \quad (12)$$

resulting in the intermediate distribution table $\mathbf{D}_{\text{init1}}$. \mathbf{T}_1 and $\mathbf{D}_{\text{init1}}$ can be represented in block matrix form.

2. **Stage 2:** By considering \mathbf{T}_1 as a block matrix, we will prove that by reorganization of submatrices within columns and circular shift of elements within columns of submatrices, we will obtain a generalized circulant matrix, which is our send communication schedule table \mathbf{C}_{send} .

We now show mathematically that these reorganizations can be performed correctly.

Stage 1 (\mathbf{D}_{init} to \mathbf{T}_1): In $\mathfrak{R}_x(P, K, Q)$, the D_i matrix will have L_s rows and P columns. With $G_1 = \gcd(P, K)$, we observe that every K_1^{th} row will have the same modulo value with respect to K , where $K_1 = \frac{K}{G_1}$.

$$\mathbf{D}_{\text{init}}(a, j) - \mathbf{D}_{\text{init}}(a \bmod K_1, j) \equiv (a/K_1)K_1P \equiv 0 \pmod{K} \quad (6)$$

There are exactly Q_1 rows such that $a \bmod K_1 = a_1$, where $0 \leq a_1 < K_1$. The corresponding Q_1 rows in *dpt* \mathbf{T} have the same pattern but their destination processor indices are different. In the first stage, these rows are gathered by moving row a to row $a' = (a \bmod K_1)Q_1 + a/K_1$, where $0 \leq a' < L_s$. In $\mathbf{D}_{\text{init1}}$, for a given a' the global index of an element can be determined by first mapping to original row, the corresponding $a = (a' \bmod Q_1)K_1 + a'/Q_1 = a'_2K_1 + a'_1$, where $0 \leq a'_1 < K_1$ and $0 \leq a'_2 < Q_1$. From this row shuffling, any element of $\mathbf{D}_{\text{init1}}(a', j)$ is then given by:

$$\begin{aligned} \mathbf{D}_{\text{init1}}(a', j) &= \mathbf{D}_{\text{init}}((a' \bmod Q_1)K_1 + a'/Q_1, j) \\ &= \mathbf{D}_{\text{init}}(a'_2K_1 + a'_1, j) \\ &= (a'_2K_1 + a'_1)P + j \end{aligned}$$

With $P_1 = \frac{P}{G_1}$, $\mathbf{D}_{\text{init1}}$ can be considered as $K_1 \times P_1$ block matrix and each submatrix is a $Q_1 \times G_1$ matrix. Let us denote these submatrices by $\mathbf{A}_{a'_1, j_1}$, where $0 \leq a'_1 < K_1$ and $0 \leq j_1 < P_1$. $\mathbf{D}_{\text{init1}}$ can be represented as,

$$\mathbf{D}_{\text{init1}} = \begin{bmatrix} \mathbf{A}_{0,0} & \mathbf{A}_{0,1} & \cdots & \mathbf{A}_{0,P_1-1} \\ \mathbf{A}_{1,0} & \mathbf{A}_{1,1} & \cdots & \mathbf{A}_{1,P_1-1} \\ \vdots & \vdots & & \vdots \\ \mathbf{A}_{K_1-1,0} & \mathbf{A}_{K_1-1,1} & \cdots & \mathbf{A}_{K_1-1,P_1-1} \end{bmatrix} \quad (7)$$

Let $j_1 = j/G_1$ and $j_2 = j \bmod G_1$ for any j . Recall that the global block index of any element is $\mathbf{D}_{\text{init1}}(a', j) = (a'_2K_1 + a'_1)P + j$, which can be written as $(a'_1P_1 + j_1)G_1 + (a'_2K_1P + j_2)$ by

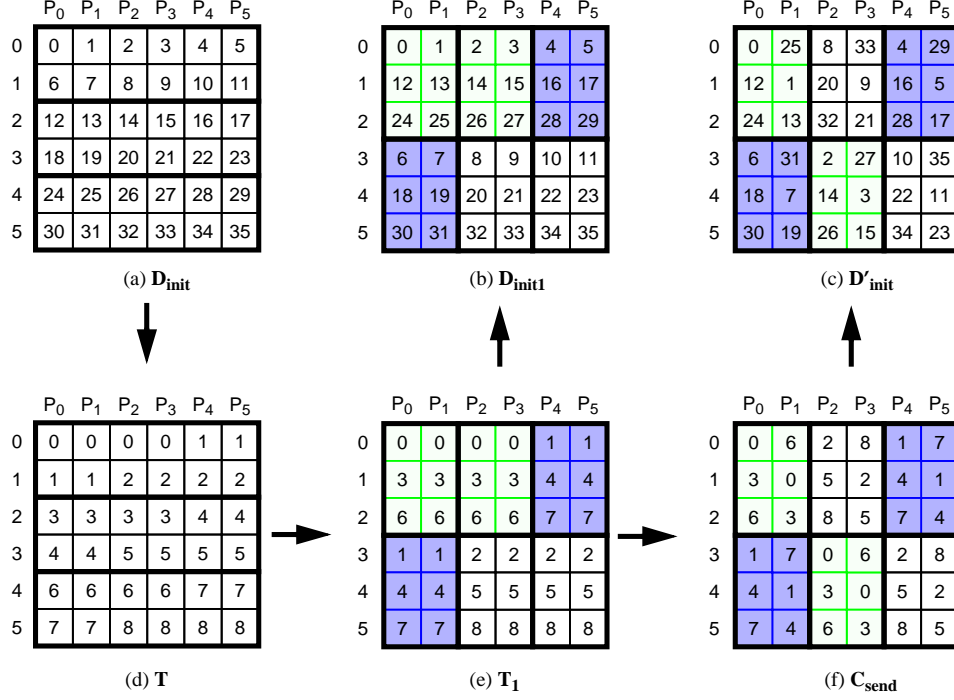


Figure 8: Steps of Column rearrangement

this procedure correctly obtains a generalized circulant matrix which is the send communication schedule table \mathbf{C}_{send} .

Theorem 1 *The initial dpt \mathbf{T} of $\mathfrak{R}_x(P, K, Q)$ with $P \leq Q$ can be reorganized via column reorganizations to a send communication schedule table \mathbf{C}_{send} such that (i) \mathbf{C}_{send} is a generalized circulant matrix and (ii) Every row of \mathbf{C}_{send} has P distinct numbers from the set of $\{0, 1, 2, \dots, Q - 1\}$.*

Proof: In the initial distribution table \mathbf{D}_{init} , each element, $\mathbf{D}_{\text{init}}(a, j)$, is assigned a global block index, $aP + j$, where $0 \leq a < L_s$ and $0 \leq j < P$. The corresponding dpt \mathbf{T} of \mathbf{D}_{init} is obtained by replacing each element by $(\mathbf{D}_{\text{init}}(a, j)/K) \bmod Q$. First, we show that the dpt \mathbf{T} can be converted to a generalized circulant matrix \mathbf{C}_{send} . Then every row of \mathbf{C}_{send} consists of P distinct numbers from the set of $\{0, 1, 2, \dots, Q - 1\}$. Therefore, \mathbf{C}_{send} can be used as a send communication schedule table.

The following are the two major steps in our approach to converting the dpt \mathbf{T} to a generalized circulant matrix.

1. **Stage 1:** We will prove that the dpt \mathbf{T} corresponding to the initial distribution table \mathbf{D}_{init} has some rows with similar patterns. These rows can be brought together in this stage. It results in the intermediate dpt \mathbf{T}_1 . Corresponding operations can be performed on \mathbf{D}_{init}

$(p - Kq) \bmod G \geq K$, none of k satisfies Eq. (5). For example, consider $p = 0$ and $q = Q - 1$ pair which gives $(p - Kq) \bmod G = K$. The right hand side has a maximum value of $K - 1$. Therefore, processor 0 doesn't send any message to processor $Q - 1$. \square

Among these three cases, the case of all-to-all processor communication with the same message size can be optimally scheduled using a trivial round-robin schedule. However, it is non trivial to achieve the same message size between all pairs of nodes in a communication step for all-to-all case with different message sizes. Therefore, we focus on the two cases of redistribution requiring scheduling of non all-to-all communication and all-to-all communication with different message sizes.

4.2 Non all-to-all communication

We first explain how $dpt \mathbf{T}$ is transformed to the send communication schedule table \mathbf{C}_{send} . Given the redistribution parameters P , Q , and K , we get the $L_s \times P$ initial distribution table \mathbf{D}_{init} and its $dpt \mathbf{T}$. Let $G_1 = \gcd(P, K)$, $K_1 = \frac{K}{G_1}$ and $P_1 = \frac{P}{G_1}$. In the dpt , every K_1 row has a similar pattern. It has a different destination processor index. We shuffle the rows such that rows having similar pattern are adjacent resulting in the shuffled $dpt \mathbf{T}_1$. The shuffled $dpt \mathbf{T}_1$ is divided into Q_1 slices in the row direction, $Q_1 = \frac{L_s}{K_1}$. It is divided into P_1 slices in the column direction. Now, $dpt \mathbf{T}_1$ can be considered as a $K_1 \times P_1$ block matrix made of $Q_1 \times G_1$ submatrices. This block matrix is then converted into a generalized circulant matrix by reorganization of submatrices within columns and rotating individual columns within submatrix by appropriate amounts. This generalized circulant matrix can then be used as a communication schedule table \mathbf{C}_{send} . In this procedure, the K identical values in row 0 of the dpt are distributed to K distinct rows, and hence, row 0 has distinct values. Since \mathbf{C}_{send} is a generalized circulant matrix, all rows are distinct and we achieve a conflict-free schedule.

In the above reorganizations, an element is moved within its column. So, it does not incur any interprocessor communication. Figure 8 shows an example where $dpt \mathbf{T}$ of $\mathfrak{R}_x(6, 4, 9)$ is converted to generalized circulant matrix form \mathbf{C}_{send} by column reorganizations. In this example, $L_s = 6$, $G_1 = 2$, $K_1 = 2$, $P_1 = 3$, and $Q_1 = 3$. Figure 8(a) shows the initial distribution table, \mathbf{D}_{init} , and (d) shows the corresponding $dpt \mathbf{T}$. Rows of \mathbf{D}_{init} and \mathbf{T} are shuffled, as shown in Figure 8(b) and (e). Now we can partition the shuffled tables into submatrices of size 3×2 . The diagonalization of submatrices and diagonalization of elements in each submatrix is shown in Figure 8(c) and (f). Figure 8(f) is a generalized circulant matrix, \mathbf{C}_{send} . In the following theorem, we will formally show that

Let $L = lcm(P, KQ)$ and $G = gcd(P, KQ)$. As discussed in Section 3.1, global blocks i and $L+i$ are mapped to the same source processor p and redistributed to the same destination processor q , because L is a least common multiplier of P and KQ . The boundaries for Eq. (3) are

$$\begin{aligned} 0 &\leq p < P \\ 0 &\leq q < Q \\ 0 &\leq k < K \end{aligned} \tag{4}$$

From the redistribution equation, we can classify communication patterns into 3 classes for the redistribution problem $\mathfrak{R}_x(P, K, Q)$ (See [5] for an alternative formulation for the *cyclic(x)* to *cyclic(y)* problem) according to the following Lemma.

Lemma 1 *The communication pattern induced by $\mathfrak{R}_x(P, K, Q)$ requires: (i) non all-to-all communication if $G > K$, (ii) all-to-all communication with the same message size if $K = \alpha G$, where α is an integer greater than 0, and (iii) all-to-all communication with different message sizes if $G < K$ and $K \neq \alpha G$.*

Proof: The redistribution equation, Eq. (3), can be rewritten as

$$p - Kq = (nQK - mP) + k$$

which can be expressed as $p - Kq = \lambda G + k$ because $nQK - mP$ is a multiple of G and λ is an arbitrary integer. For any p and q , if there is at least one block satisfying the above equation, this redistribution requires all-to-all communication. When both sides are divided by G , their remainders are

$$(p - Kq) \bmod G = k \bmod G \tag{5}$$

Assume that $G \leq K$. There is at least one k satisfying Eq. (5), since k is between $[0, K - 1]$. Note that if $K = \alpha G$ for any integer $\alpha > 0$, k is between $[0, \alpha G - 1]$. Thus, for the expression $k \bmod G$ in Eq. (5), α distinct numbers result in the same remainder. Therefore, Eq. (5) has α solutions for any (p, q) pair. This redistribution requires all-to-all communication with the same message size. In the initial distribution table, α blocks in a column are transferred to the same destination processor. If $K > G$ and $K \neq \alpha G$, then it requires all-to-all communication with different message sizes. Assuming that $G > K$, we can find a processor pair (p, q) exchanging no message. When

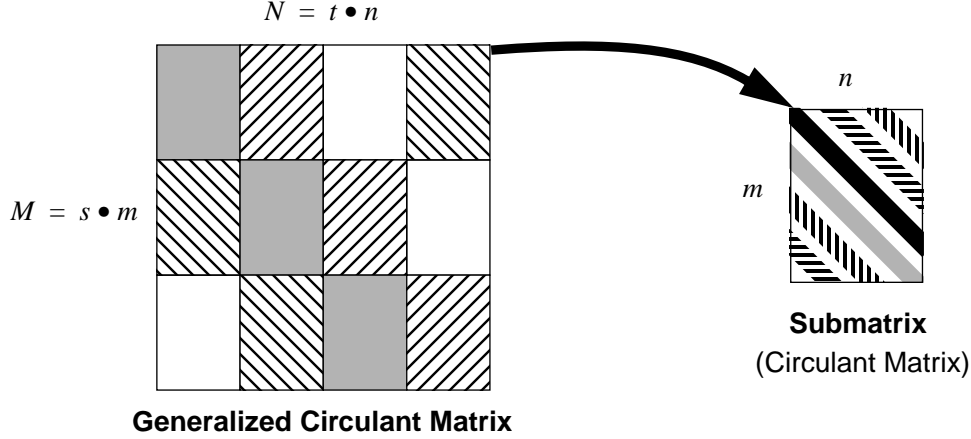


Figure 7: Generalized circulant matrix

4 Efficient Redistribution Algorithms

In this section, we discuss the communication patterns that arise in performing array redistribution, develop the algorithms for our table-based approach to obtain the send schedule, and present correctness proofs of our techniques.

4.1 Communication patterns of redistribution

Consider the movement of the global block i in the redistribution problem $\mathfrak{R}_x(P, K, Q)$. It is the m^{th} local block of source processor p , where $m = i/P$ and $p = i \bmod P$.

$$i = P \times m + p \tag{1}$$

After redistribution, K consecutive global blocks become one block in the new layout. Therefore, global block i is k^{th} block in the i_K^{th} new global block, where $k = i \bmod K$ and $i_K = i/K$. The i_K^{th} new global block is at n^{th} new local block of destination processor q , where $n = i_K/Q$ and $q = i_K \bmod Q$.

$$i = (Q \times n + q) \times K + k \tag{2}$$

We can derive the following redistribution equation.

$$i = P \times m + p = (Q \times n + q) \times K + k \tag{3}$$

For example, consider the 16^{th} global block in $\mathfrak{R}_x(3, 2, 6)$ of Figure 1. Here, $i = 16$, $P = 3$, $p = 1$ and $m = 5$. All block indices start from 0. Also, $K = 2$, $k = 0$, $i_K = 8$, $n = 1$ and $q = 2$.

0	3	2	1
1	0	3	2
2	1	0	3
3	2	1	0

(a) $m = n = 4$

0	3	2	1
1	0	3	2
2	1	0	3

(b) $m = 3, n = 4$

0	3	2
1	0	3
2	1	0
3	2	1

(c) $m = 4, n = 3$

Figure 6: Circulant matrix examples

In the following subsection, we define generalized circulant matrix. Using the generalized circulant matrix, we derive an efficient contention-free communication schedule for $\mathfrak{R}_x(P, K, Q)$.

3.3 Communication scheduling using generalized circulant matrix

Our framework for communication schedule performs local rearrangement of data within each processor as well as interprocessor communication. The local rearrangement of data, which we call column reorganization, results in a send communication schedule table \mathbf{C}_{send} . We will show that for any P , K and Q , the send communication schedule is indeed a generalized circulant matrix which avoids node contention.

Definition 1 *An $m \times n$ matrix is a circulant matrix if it satisfies the following properties:*

1. *If $m \leq n$, row $k =$ row 0 circularly right shifted k times, $0 \leq k < m$.*
2. *If $m > n$, column $l =$ column 0 circularly down shifted l times, $0 \leq l < n$.*

Figure 6 shows several different circulant matrices. Note that the above definition can be extended to block circulant matrices by changing “row” to “row block”.

Definition 2 *An $M \times N$ matrix is a generalized circulant matrix if the matrix can be partitioned into blocks of size $m \times n$, where $M = s \cdot m$ and $N = t \cdot n$, for some $s, t > 0$ such that the block matrix forms a circulant matrix and each block is either a circulant matrix or a generalized circulant matrix.*

Figure 7 illustrates a generalized circulant matrix. There are two observations about the generalized circulant matrix: (i) the s blocks along each block diagonal are identical, and (ii) if all the elements in row 0 are distinct, then in each row all elements are distinct.

We will show that through our approach the destination processor table \mathbf{T} is transformed to a generalized circulant matrix \mathbf{C}_{send} with distinct elements in each row.

Figure 5 shows our table-based framework for redistribution. To convert the initial distribution table \mathbf{D}_{init} to the final distribution table $\mathbf{D}_{\text{final}}$, *dpt* \mathbf{T} can be used. But, the use of \mathbf{T} itself as a communication schedule is not efficient. It leads to node contention, since several processors try to send their data to the same destination processor in a communication step. For example, in Figure 5, during step 0, both source processors 0 and 1 try to communicate with destination processor 0. However, if every row of \mathbf{T} consists of P distinct destination processor indices among $\{0, 1, \dots, Q-1\}$, node contention can be avoided in each communication step. This is the motivation for the column reorganizations.

To eliminate node contention, the *dpt* \mathbf{T} is reorganized by column reorganizations. The reorganized table is called a *send communication schedule table*, \mathbf{C}_{send} . In section 4, we discuss how these reorganizations are performed. \mathbf{C}_{send} is a $L_s \times P$ matrix as well. Each entry of \mathbf{C}_{send} is a destination processor index and each row corresponds to a contention-free communication step. To maintain the correspondence between \mathbf{D}_{init} and \mathbf{T} , the same set of column reorganizations is applied to \mathbf{D}_{init} which results in a distribution table, $\mathbf{D}'_{\text{init}}$ corresponding to \mathbf{C}_{send} . In a communication step, blocks in a row of $\mathbf{D}'_{\text{init}}$ are transferred to their destination processors specified by the corresponding entries in \mathbf{C}_{send} . Referring to Figure 5, in the first communication step, source processors 0, 1 and 2 transfer blocks 0, 4 and 2 to destination processors 0, 2 and 1 respectively as specified by \mathbf{C}_{send} . Such a step is called row reorganization. The distribution table $\mathbf{D}'_{\text{final}}$ corresponding to the received blocks in destination processors is reorganized into the final distribution table $\mathbf{D}_{\text{final}}$ by another set of column reorganizations. (For this example, we do not need this operation.) The received blocks are then stored in the memory locations of the destination processors. The key idea is to choose a \mathbf{C}_{send} such that the required row reorganizations (communication events) can be performed efficiently and it supports easy-to-compute contention-free communication scheduling.

So far, we have discussed a redistribution problem from *cyclic*(x) on P processors to *cyclic*(Kx) on Q processors. A dual relationship exists between the problem from *cyclic*(x) on P processors to *cyclic*(Kx) on Q processors and the problem from *cyclic*(Kx) on Q processors to *cyclic*(x) on P processors. The redistribution from *cyclic*(Kx) on Q processors to *cyclic*(x) on P processors is the redistribution with reverse direction of the redistribution $\mathfrak{R}_x(P, K, Q)$. Its send (receive) communication schedule table is the same as the receive (send) communication schedule table of $\mathfrak{R}_x(P, K, Q)$. Therefore, our scheme for $\mathfrak{R}_x(P, K, Q)$ can be extended to the redistribution problem from *cyclic*(Kx) on Q processors to *cyclic*(x) on P processors.

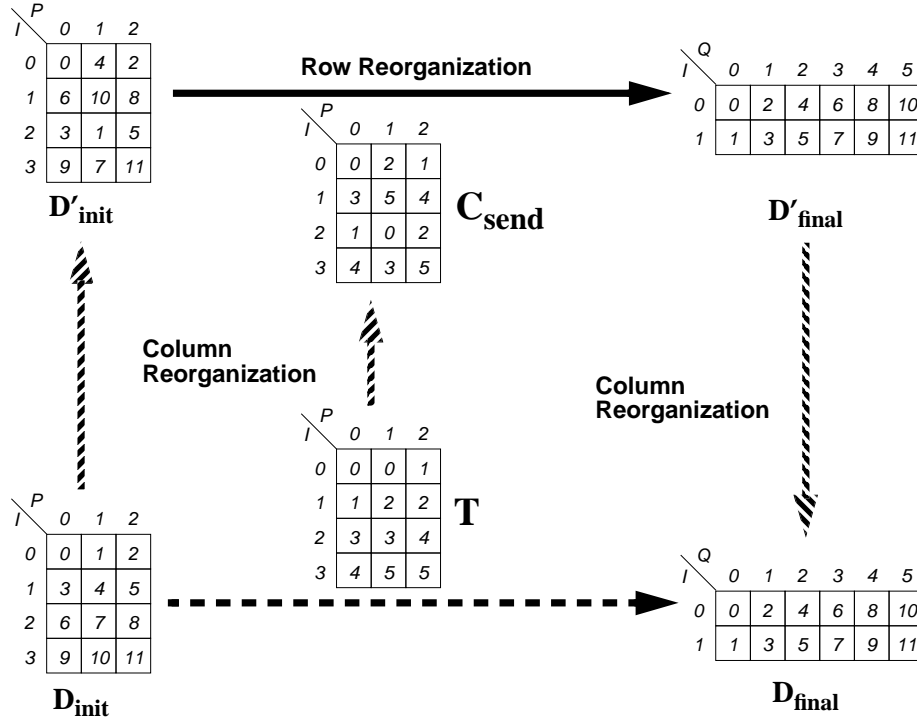


Figure 5: Table conversion process for redistribution

redistribution parameters and its global block index. A send communication events table is constructed by replacing each block index in the initial distribution table with its destination processor index as shown in Figure 4. This is denoted as *destination processor table (dpt)* \mathbf{T} . The $(i, j)^{th}$ entry of \mathbf{T} is the destination processor index of i^{th} local block in source processor j and $0 \leq i < L_s$ i.e. \mathbf{T} considers only one superblock. It is a $L_s \times P$ matrix. Each row corresponds to a communication step. On the other hand, each destination processor has to know from which source processors it receives messages. Source processor index of each block in the final distribution is also determined by the redistribution parameters and its global block index. A receive communication events table can be constructed by replacing each global block index in \mathbf{D}_{final} with its source processor index. In our algorithm, during a communication step, a processor sends data to at most one destination processor. If $Q \geq P$, at most P processors in the destination processor set can receive data and the other destination processors remain idle during that communication step. Therefore, each communication step can have at most P communicating pairs. On the other hand, if $Q < P$, only Q destination processors can receive data at a time. The maximum number of communicating pairs in a communication step is $\min(P, Q)$. Without loss of generality, in the following discussion we assume that $Q \geq P$.

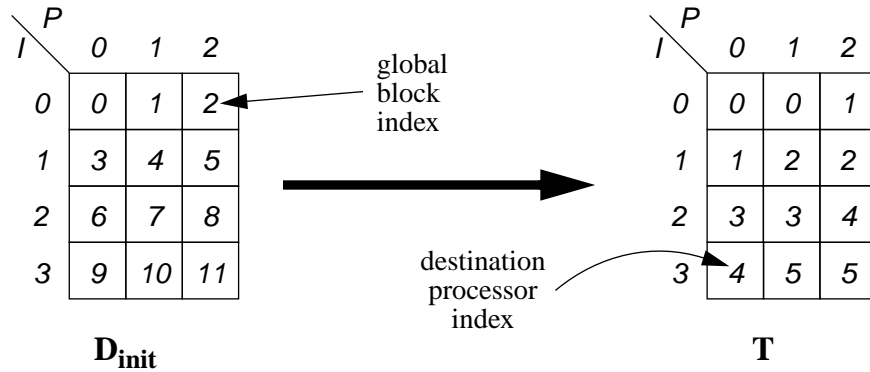


Figure 4: The destination processor table \mathbf{T}

The initial layout can be partitioned into collections of rows. Each collection consists of $L_s = \text{lcm}(P, KQ)/P$ rows. Similarly, the final layout can be partitioned into disjoint collections of rows; each collection having $L_d = \text{lcm}(P, KQ)/Q$ rows. Note that each collection corresponds to a superblock. Blocks, which are located at the same relative position within a superblock, are moved in the same way during the redistribution. These blocks can be transferred in a single communication step. The MPI derived data type can handle these blocks as a single block. Without loss of generality, we will consider only the first superblock in the following to illustrate our algorithm. We refer to the tables representing the indices of the blocks within the first superblock in the initial (final) layout as *initial distribution table* \mathbf{D}_{init} (*final distribution table* $\mathbf{D}_{\text{final}}$). These are shown in Figure 3(c) and (f), respectively. The cyclic redistribution problem essentially involves reorganizing blocks within each superblock from an initial distribution table \mathbf{D}_{init} to a final distribution table $\mathbf{D}_{\text{final}}$.

3.2 A table-based framework for redistribution

Given the redistribution parameters, P , K , and Q , each block's location in \mathbf{D}_{init} and $\mathbf{D}_{\text{final}}$ can be determined. Through redistribution, each block moves from its initial location in \mathbf{D}_{init} to the final location in $\mathbf{D}_{\text{final}}$. Thus, the processor ownership and the local memory location of each block are changed by redistribution. This redistribution can be conceptually considered as a table conversion process from \mathbf{D}_{init} to $\mathbf{D}_{\text{final}}$, which can be decomposed into independent column and row reorganizations. In a column reorganization, blocks are rearranged within a column of the table. This is therefore a local operation within a processor's memory. In a row reorganization, blocks within a row are rearranged. This operation therefore leads to a change in ownership of the blocks, and requires interprocessor communication.

The destination processor of each block in the initial distribution table is determined by the

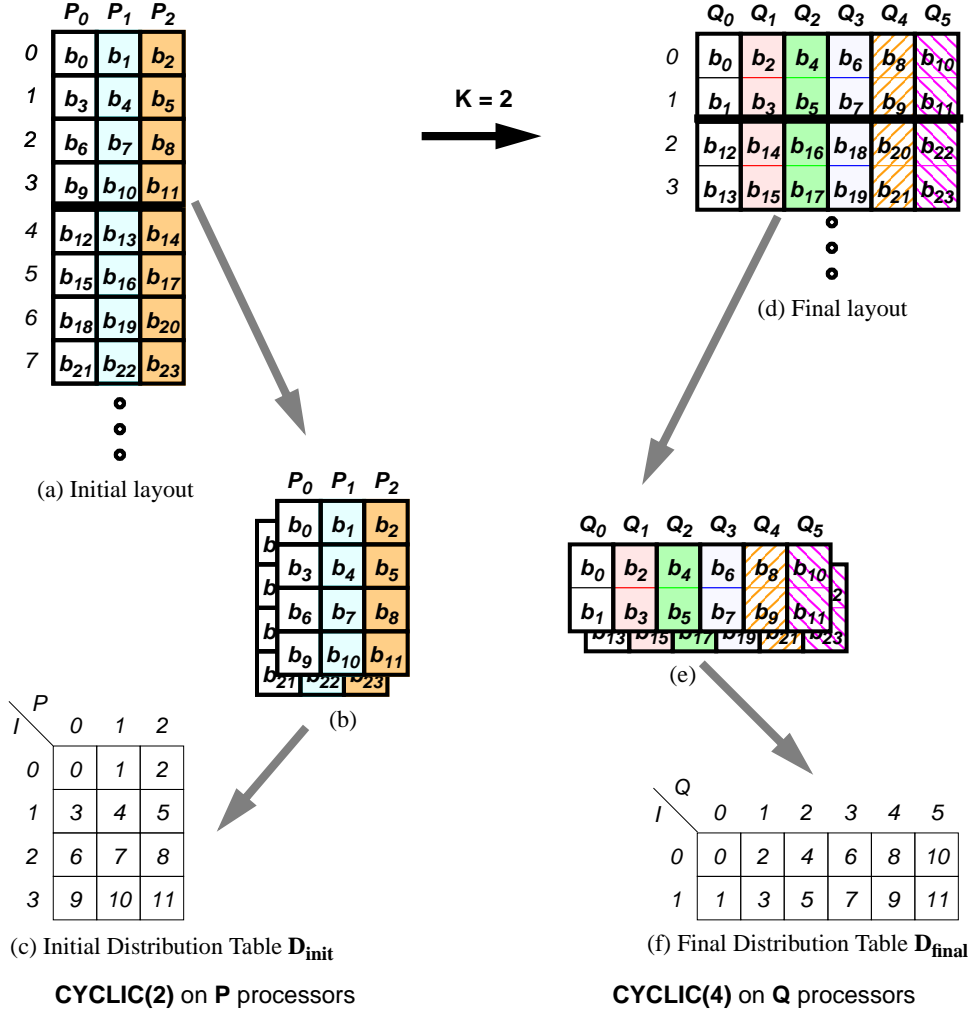


Figure 3: Block-Cyclic Redistribution from $\text{cyclic}(x)$ on P processors to $\text{cyclic}(Kx)$ on Q processors from processor point of view. In this example, $P = 3$, $Q = 6$ and $K = 2$.

same fashion.

From the *processor point of view*, the block-cyclic distribution can be represented by a 2-dimensional table. Each column corresponds to a processor and each row index is a local block index. Each entry in the table is a global block index. Therefore, element (i, j) in the table represents the i^{th} local block of the j^{th} processor. Figure 3 shows the example of $\mathfrak{R}_2(3, 2, 6)$ from the processor point of view. Blocks are distributed on the table in a round-robin fashion. The table corresponding to source processors is denoted as initial layout representing which blocks are initially assigned to which source processors. Similarly, the final layout represents which blocks are assigned to which destination processors. Our problem is to redistribute the blocks from initial layout to final layout. These layouts are shown in Figure 3(a) and (d) respectively.

$cyclic(y)$ on Q processors, where x, y, P , and Q are arbitrary positive integers. Their method to compute the communication schedule uses bipartite matching. They propose two strategies, stepwise strategy and greedy strategy. In stepwise strategy, they try to minimize the number of communication steps but the total data transfer cost is not optimized. In greedy strategy, the total transmission cost is optimized but the number of communication steps is not minimized. Even though they reduced the data transfer time, the time to compute the communication schedule using bipartite matching is significant. The schedule computation cost is $O((P + Q)^4)$. As the number of processors is increased, the schedule computation time can be larger than the data transfer cost.

In [4], Y.C. Chung *et al.* have proposed the index computation method for redistribution from $cyclic(x)$ to $cyclic(y)$ on the same processor set. They proposed the basic-cycle calculation (denoted as BCC) technique which is the closed forms for source/destination processor indices of array elements. These closed forms are useful to efficiently determine the communication sets of a basic-cycle. They did not consider the communication schedule in this technique. Therefore, the node contention problems exist on the redistribution communications.

3 Our Approach to Redistribution

In this section, we present our approach to block-cyclic redistribution problem. In subsection 3.1, we discuss two views of redistribution and illustrate the concept of a superblock. In the following subsection, we explain our table-based framework for redistribution using the destination processor table and column and row reorganizations. In subsection 3.3, we discuss the generalized circulant matrix formalism which allows us to compute communication schedule efficiently.

3.1 Array and processor points of view

Figure 1 shows $\mathfrak{R}_2(3, 2, 6)$ from the *array point of view*. The elements of the array are shown along a single horizontal axis. The processor indices are marked above each block. For the redistribution $\mathfrak{R}_x(P, K, Q)$, a periodicity can be found in the block movement pattern. For example, in Figure 1, b_0, b_3, b_6 , and b_9 , which are initially assigned to P_0 , are moved to Q_0, Q_1, Q_3 , and Q_4 respectively. After that, b_{12} in P_0 is moved to Q_0 again. We find that the communication pattern between b_0 and b_{11} is repeated on other blocks. Such a collection of blocks is called as a superblock. The period of this block movement pattern is $lcm(P, KQ)$, and is the size of the superblock. In Figure 1, superblock size is $lcm(3, 2 \times 6) = 12$. In the next superblock, blocks b_{12} to b_{23} are moved in the

processor sets as well as multidimensional arrays. However, PITFALLS does not consider communication scheduling during redistribution. Their main contribution is to determine the elements to be exchanged.

Other previous studies [4, 8, 10, 22, 23, 24] consider redistribution from $cyclic(x)$ to $cyclic(Kx)$ on the same set of processors. This is a classical redistribution problem. Techniques to compute index sets are proposed in [8, 22, 23]. However, explicit scheduling of the communication to minimize overall redistribution time is not considered.

In [22, 23], Choudhary *et al.* present efficient index computation algorithms for the special case when $P \bmod K \equiv 0$. They also consider redistribution from $cyclic(x)$ to $cyclic(y)$ on the same processor set, for arbitrary x and y . Although it is possible to explicitly calculate the destination and source processor of each element of the local array, such a scheme is expensive for use in practice as potential node contentions occur in each communication step. In [23], *gcd* and *lcm* methods have been proposed. These are two phase algorithms where $cyclic(x)$ is first redistributed to $cyclic(m)$, followed by a redistribution from $cyclic(m)$ to $cyclic(y)$. Here, m can be *gcd* or *lcm* of x and y . In [23], it is shown that multidimensional arrays can be redistributed by applying these algorithms to each dimension of the array separately.

In [8], Ni *et al.* provide new logical processor numbers (*lpids*) for the target distribution, so as to minimize the amount of data to be communicated during redistribution. Data blocks which have the same *lpids* across source and target distributions, need not be moved. However, index set computation becomes complicated. This approach is not applicable when the number of source and target processors are different.

Sadayappan *et al.* [10] and Walker *et al.* [24] have proposed algorithms which reduce communication overheads. In [24], a K step schedule is derived for $cyclic(x)$ to $cyclic(Kx)$ redistribution on the same processor set. In each step, processors exchange data in a contention free manner: each processor sends its data to exactly one processor and receives data from exactly one processor. Therefore contention at the destination nodes does not occur. A similar communication schedule is shown in [10]. Although the communication schedule presented in [24] is based on modular arithmetic and that in [10] is based on tensor products, the resulting communication schedules are similar.

In [5], Desprez *et al.* propose a general solution for block cyclic redistribution on arbitrary source and destination processor sets. They can handle redistribution from $cyclic(x)$ on P processors to

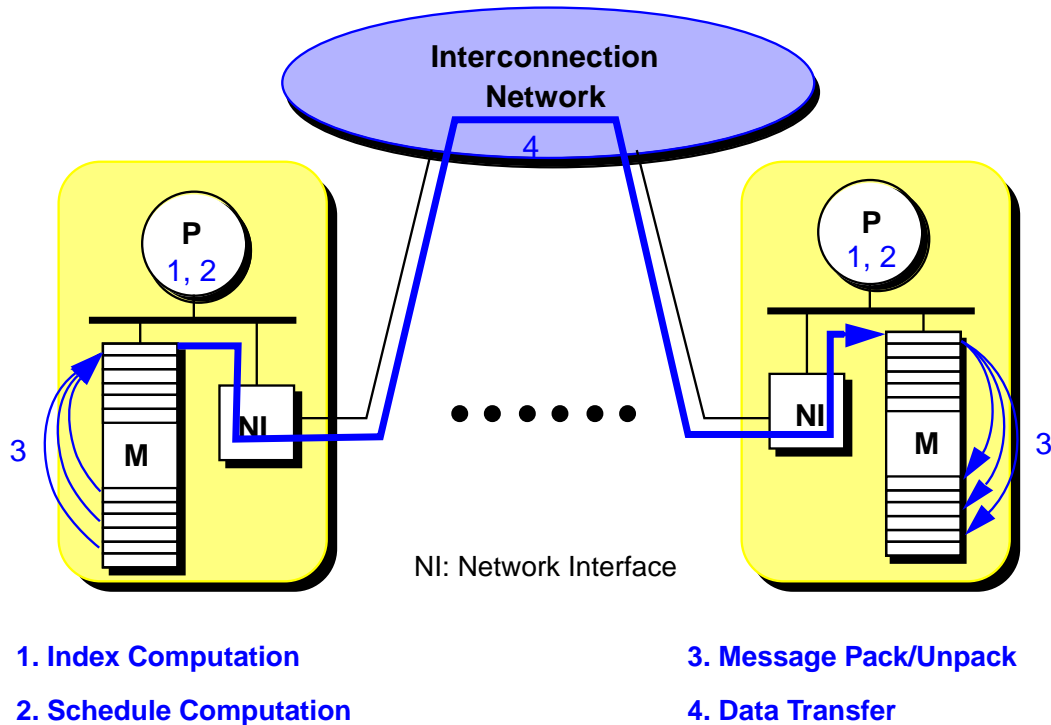


Figure 2: Steps in performing redistribution

the buffer. The time to perform this data gathering at the sender is the message packing cost. Similarly, at the receiving side, each message is to be unpacked and data words need to be stored in appropriate memory locations.

Data Transfer Cost: The data transfer cost for each communication step consists of start-up cost and transmission cost. The start-up cost is incurred by software overheads in each communication operation. The total start-up cost can be reduced by minimizing the number of message transfer steps. The transmission cost is incurred in transferring the bits over the network and depends on the network bandwidth.

2.3 Related work

The array redistribution problem has been the focus of several research efforts [4, 5, 8, 10, 19, 22, 23, 24]. Many of these efforts have targeted efficient implementation of high level compiler directives such as the `REDISTRIBUTE` directive in HPF.

In [19], Banerjee *et al.* use a line segment formalism called PITFALLS to represent the *cyclic(x)* distribution. The array elements that map to a processor are represented as a set of stride line segments. For every pair of processors, the array elements whose indices are in the intersection of the respective line sets are exchanged. Their technique handles arbitrary source and destination

Table 4: An example of all-to-all communication with different message size. $P = 8$, $Q = 10$, and $K = 6$.

Source \ Destination	0	1	2	3	4	5	6	7	8	9
0	2	1	2	1	2	1	2	1	2	1
1	2	1	2	1	2	1	2	1	2	1
2	1	2	1	2	1	2	1	2	1	2
3	1	2	1	2	1	2	1	2	1	2
4	2	1	2	1	2	1	2	1	2	1
5	2	1	2	1	2	1	2	1	2	1
6	1	2	1	2	1	2	1	2	1	2
7	1	2	1	2	1	2	1	2	1	2

communication with different message sizes. This is shown in Table 4 where $N = 120x$, $P = 8$, $Q = 10$, and $K = 6$. Here half the messages are of size $1x$ while others are of size $2x$.

2.2 Cost of performing redistribution

We briefly explain the four costs associated with data redistribution:

Index Computation Cost: Each source processor determines the destination processor indices of the array elements that belong to it and computes the local memory location (local index) of each array element. This local index is used to pack array elements into a message. Similarly, each destination processor determines the source processor indices of received messages and computes their local indices to find out the location where the received message is to be stored. The total time to compute these indices is denoted as index computation cost.

Schedule Computation Cost: The communication schedule specifies a collection of sender-receiver pairs for each communication step. Since in each communication step, a processor can send at most one message and a processor can receive at most one message, careful scheduling is required to avoid contention while minimizing the number of communication steps. Time to compute this communication schedule can be significant. Reducing this cost is an important criteria in performing run-time redistribution.

Message Packing/Unpacking Cost: At each sender, a message consists of words from different memory locations which need to be gathered in a buffer in the sending node. Typically, this requires a memory read and a memory write operation to gather the data to form a compact message in

Table 2: An example of non all-to-all communication. $P = 6$, $Q = 10$ and $K = 3$.

Source \ Destination	0	1	2	3	4	5	6	7	8	9
0	1		1		1		1		1	
1	1		1		1		1		1	
2	1		1		1		1		1	
3		1		1		1		1		1
4		1		1		1		1		1
5		1		1		1		1		1

Table 3: An example of all-to-all communication with same message size. $P = 6$, $Q = 10$, and $K = 4$.

Source \ Destination	0	1	2	3	4	5	6	7	8	9
0	2	2	2	2	2	2	2	2	2	2
1	2	2	2	2	2	2	2	2	2	2
2	2	2	2	2	2	2	2	2	2	2
3	2	2	2	2	2	2	2	2	2	2
4	2	2	2	2	2	2	2	2	2	2
5	2	2	2	2	2	2	2	2	2	2

P processors in a round-robin fashion. Block b_i is assigned to processor j , P_j , where $j = i \bmod P$.

In this paper, we study the problem of redistributing from $cyclic(x)$ on P processors to $cyclic(Kx)$ on Q processors, which is denoted as $\mathfrak{R}_x(P, K, Q)$. Figure 1 (c) shows $\mathfrak{R}_2(3, 2, 6)$. Here, initially pairs of consecutive elements of the array are distributed over $P = 3$ source processors. Then, the block size is doubled and the new blocks (of size 4) are redistributed over $Q = 6$ destination processors.

In performing $\mathfrak{R}_x(P, K, Q)$, three classes of communication patterns between source and destination processors arise. One case is the non *all-to-all* communication: every source processor communicates with some of the destination processors. This case is shown in Table 2, where $N = 30x$, $P = 6$, $Q = 10$ and $K = 3$. Note that the messages are of equal size ($1x$). The second case is the *all-to-all* communication with the same message size. In Table 3, all the source processors have messages of the same size ($2x$) and each source processor communicates with all the destination processors. Here, $N = 120x$, $P = 6$, $Q = 10$ and $K = 4$. The last case is the *all-to-all*

of any scheme that generates an optimal communication schedule. Thus, our scheme minimizes the total time for data redistribution. This makes our scheme attractive for run-time as well as compile-time data redistribution.

Our techniques can be used for implementing scalable redistribution libraries, for implementing the REDISTRIBUTE directive in HPF [11], and for developing parallel algorithms for supercomputer applications. In particular, these techniques lead to efficient distributed corner turn operation, a key communication kernel needed in parallelizing signal processing applications [13, 21, 25].

Our redistribution scheme has been implemented using MPI and C. It can be easily ported to various HPC platforms. We have performed several experiments to illustrate the improved performance compared with the state-of-the-art. The experiments were performed to determine the data transfer, schedule and index computation costs. In one of these experiments, we used 64 processors on an IBM SP2 which were partitioned into 28 source processors and 36 destination processors. The expansion factor was set to 14. The array size was varied from 2.26 Mbytes to 56.4 Mbytes. Compared with the Caterpillar algorithm, our data transfer times were lower. The ratio of the data transfer time of our algorithm to that of the Caterpillar algorithm was between 49.2% and 53.7%. Although the resulting data transfer time using the bipartite matching scheme and that of our algorithm will be the same, the schedule computation time for the bipartite matching scheme will be significantly larger than our scheme.

The rest of this paper is organized as follows. Section 2 defines the array redistribution problem and reviews prior work. Section 3 explains our table-based framework. It also discusses the generalized circulant matrix formalism for deriving conflict free communication schedules. Section 4 explains our redistribution algorithm and index computation in detail and gives the proofs of their correctness and their complexity. Section 5 reports our experimental results on the IBM SP-2. Concluding remarks are made in Section 6.

2 Background and Related Work

2.1 Problem definition

The block-cyclic distribution, $cyclic(x)$, of an array is defined as follows: given an array with N elements, P processors, and a block size x , the array elements are partitioned into contiguous blocks of x elements each. The i^{th} block, b_i , consists of array elements whose indices vary from ix to $(i + 1)x - 1$, where i is a global block index and $0 \leq i < \frac{N}{x}$. These blocks are distributed onto

directly compare our scheme with the bipartite matching scheme in [5] as they specifically excluded this case. If we do apply their bipartite graph formulation and find a series of matchings, we will either obtain non-optimal number of steps with equal sized messages in each step or optimal number of steps (trivially using Caterpillar scheme) with unequal sized messages in some steps.

In the other interesting case where non-all-to-all communication is performed, each source processor communicates with the same number of destination processors. Thus, each source processor will have exactly one message to any of its destination processors. Therefore, the bipartite graph for the non-all-to-all communication case will be regular. The approach in [5] reduces to finding a series of maximum matchings in a regular bipartite graph with unit edge weights. Although we know the optimal number of matchings to obtain a communication schedule, it is still non trivial to exploit the above property to derive the schedule. After the first maximum matching is found by the stepwise approach in [5], the regularity is lost as the number of source and destination processors are not the same. To the best of our knowledge, the fastest scheme to find a maximum matching in a unit weight bipartite graph has $O(V^{1/2}E)$ complexity, where V is the number of nodes and E is the number of edges [6, 17, 20]. For our problem, since there can be $lcm(P, KQ)$ such matchings, the complexity to compute the schedule will be $O((P + Q)^{3.5})$. Although this complexity is better compared with the general case, it is still expensive to compute these graph matchings. On a state-of-the-art workstation, the schedule computation time using an efficient bipartite matching algorithm is on the order of several *msecs* for P and Q around 40 (See Table 6). As we will show in our experimental results section, the schedule computation cost can be up to 50% of the data transfer cost for P and Q around 40 and data array size of few Mbytes.

In this paper, we propose a novel and efficient algorithm for data redistribution from *cyclic*(x) on P processors to *cyclic*(Kx) on Q processors. Our algorithm uses *optimal* number of communication steps and *fully* utilizes the network bandwidth in each step. The communication schedule is determined using a generalized circulant matrix framework. The communication schedule of each node is computed in parallel. The schedule computation cost is $O(max(P, Q))$. Our implementations show that the schedule computation time is in the range of 100's of $\mu secs$ when P and Q are in the range 50-100. Each processor computes its own index set and its communication schedule using only a set of equations derived from our generalized circulant matrix formulation. Our experimental results show that the schedule computation time is negligible compared with the data transfer cost for array sizes of interest. The message packing/unpacking cost is the same as that

Table 1: Comparison of various schemes for array redistribution

	Key Features	
	Schedule & Index Computation	Communication
PITFALLS [19]	<ul style="list-style-type: none"> No communication scheduling Index computation using a line segment formalism 	<ul style="list-style-type: none"> Node contention occurs Does not minimize the transmission cost
BCC [4]	<ul style="list-style-type: none"> No communication scheduling Efficient index computation Source and destination sets are same 	<ul style="list-style-type: none"> Node contention occurs Does not minimize the transmission cost
Caterpillar [18]	<ul style="list-style-type: none"> Simple scheduling algorithm Index computation by scanning the array segments 	<ul style="list-style-type: none"> No node contention Does not minimize the transmission cost and the number of communication steps
Bipartite Matching Scheme [5]	<ul style="list-style-type: none"> Large schedule computation overhead Schedule computation time: $O((P+Q)^{3.5})$ 	<ul style="list-style-type: none"> No node contention Stepwise strategy: minimizes the number of communication steps Greedy strategy: minimizes the transmission cost
Our Scheme	<ul style="list-style-type: none"> Fast schedule and index computations Schedule computation time: $O(\max(P,Q))$ 	<ul style="list-style-type: none"> No node contention Minimizes the number of communication steps and the data transfer cost

transfer cost.

In [5] the general problem of data redistribution from $cyclic(x)$ on P processors to $cyclic(y)$ on Q processors is elegantly formulated as a bipartite graph matching problem. The technique is to construct a bipartite graph with P source nodes and Q destination nodes with edges representing the amount of data communication between pairs of processors. A series of maximal matchings is found and each matching represents a conflict free communication and the set of matchings represent the communication schedule. In the experimental section of [5] the schedule computation time is *not* included in the total data redistribution time. However, the schedule computation time can be significant as the complexity of bipartite graph matching procedure used is $O((P+Q)^4)$ for the general case of the problem. Further, the authors consider only those scenarios that result in non-all-to-all communication among processor sets.

In this paper, we consider the problem of rearranging data from $cyclic(x)$ on P processors to $cyclic(Kx)$ on Q processors. For this class of redistribution problem there are three different cases: non all-to-all communication, all-to-all communication with different message sizes, and all-to-all communication with the same message size. The all-to-all communication with the same message size is a trivial one as there is no need to perform any scheduling and the Caterpillar scheme will be optimal. However, for the other two cases, our technique provides significant performance improvement. For the case of all-to-all communication with different message sizes, we cannot

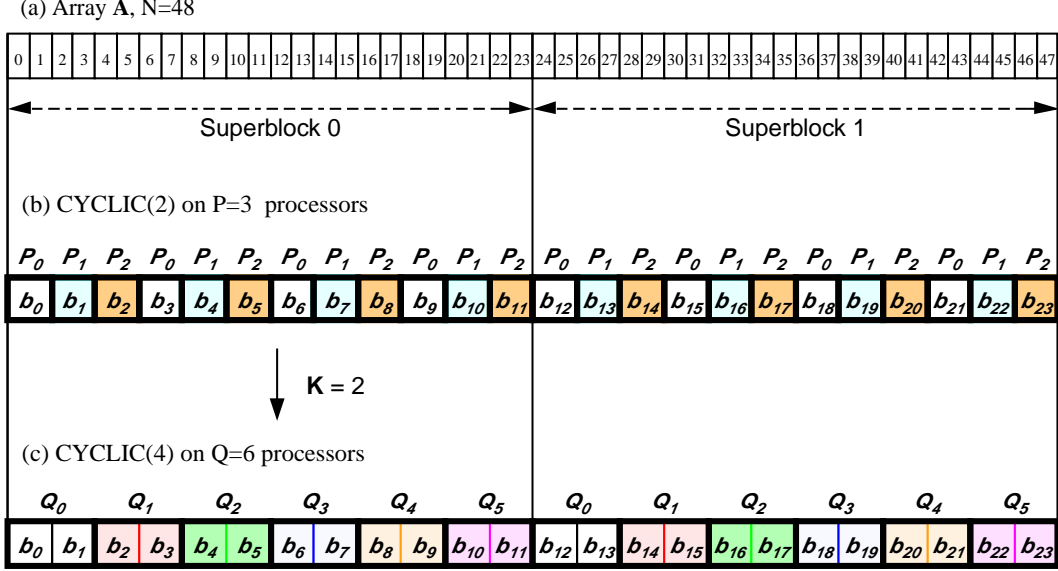


Figure 1: Block-Cyclic Redistribution from array point of view: (a) the array of elements, (b) $cyclic(x)$ on P processors, (c) from $cyclic(x)$ on P processors to $cyclic(Kx)$ on Q processors. In this example, $x = 2$, $P = 3$, $Q = 6$ and $K = 2$.

processor will unpack a received message and place the data blocks in appropriate local memory locations. Finally, the transfer of data blocks over the network incurs a fixed start-up cost plus a transmission cost proportional to the size of a message. All these four steps contribute to the total array redistribution cost.

Table 1 summarizes the key features of the well known data distribution algorithms in the literature. All of the known algorithms ignore one or more of the above costs. Some schemes focus only on efficient index set computation and completely ignore scheduling the communication events [4, 8, 19, 22, 23]. Based on the index of a block, these schemes focus on finding its destination processor and generating messages for the same destination. Communication scheduling is not considered. These lead to node contention in performing the communication. This in turn leads to higher data transfer costs as some nodes incur additional delays. Other schemes eliminate node contention by explicitly scheduling the communication events [10, 12, 24]. Although the schemes in [10, 12, 24] have an efficient scheduling algorithm, these were designed for data redistribution on the same processor set. For redistribution between different processor sets, the Caterpillar algorithm was proposed in [18]. It uses a simple round robin schedule to avoid node contention. However, this algorithm does not fully utilize the network bandwidth i.e., the size of the data sent by the nodes in a communication step varies from node to node. This leads to increased data

1 Introduction

Many High Performance Computing (HPC) applications, including scientific computing and signal processing, consist of several stages [13, 21, 25]. Examples of such applications include the multi-dimensional Fast Fourier Transform, the Alternative Direction Implicit (ADI) method for solving two-dimensional diffusion equation, and linear algebra solvers. As the program execution proceeds from one computational stage to another, the data access patterns and the number of processors required for exploiting the parallelism in the application may change. These changes usually cause the data distribution in a stage to be unsuitable for the subsequent stage. Therefore, data redistribution is needed to between two subsequent stages to reduce the number of remote accesses. Since the parameters of redistribution are generally unknown at compile time, run-time data redistribution is necessary. However, the cost of redistribution can offset the performance benefits that can be achieved by the redistribution. Therefore, run-time redistribution must be implemented efficiently to ensure overall performance improvement.

Array data are typically distributed in a *block-cyclic* pattern onto a given set of processors. The block-cyclic distribution with block size x is denoted as *cyclic*(x). A block contains x consecutive array elements. Blocks are assigned to processors in a round-robin fashion. Other distribution patterns, *cyclic* and *block* distribution, are special cases of the block-cyclic distribution. In general, the block-cyclic array redistribution problem is to reorganize an array from one block-cyclic distribution to another, *i.e.*, from *cyclic*(x) to *cyclic*(y). An important case of this problem is a redistribution from *cyclic*(x) to *cyclic*(Kx) which arises in many scientific and signal processing applications. This type of data redistribution can occur within a same processor set, or between different processor sets. An example of block-cyclic redistribution is shown in Figure 1.

Data redistribution from a given initial layout to a final layout consists of four major steps - index computation, communication schedule, message packing and unpacking, and finally data transfer. With a given initial and final data layout, each processor determines the local indices of the data blocks that belong in its memory initially and at the end of redistribution. This calculation from the global indices to local indices of data blocks is called index computation. Data blocks in a source processor need to be moved to destination processor depending on the required final data layout. The parallel communications of data blocks among processors will be determined by a communication schedule. For efficient communication of data blocks, a source processor should pack all blocks destined to a destination processor in one message. Likewise, a destination

Efficient Algorithms for Block-Cyclic Array Redistribution between Processor Sets

Neungsoo Park* and Viktor K. Prasanna*
Department of EE-Systems, EEB-200C
University of Southern California
Los Angeles, CA 90089-2562
{neungsoo + prasanna}@halcyon.usc.edu
<http://ceng.usc.edu/~prasanna/>

C. S. Raghavendra
The Aerospace Corporation
P.O. Box 92957
Los Angeles, CA 90009-2957
raghu@rush.aero.org

Revised March 22, 1999

Abstract

Run-time array redistribution is necessary to enhance the performance of parallel programs on distributed memory supercomputers. In this paper, we present an efficient algorithm for array redistribution from *cyclic*(x) on P processors to *cyclic*(Kx) on Q processors. The algorithm reduces the overall time for communication by considering the data transfer, communication schedule, and index computation costs. The proposed algorithm is based on a generalized circulant matrix formalism. Our algorithm generates a schedule that minimizes the number of communication steps and eliminates node contention in each communication step. The network bandwidth is fully utilized by ensuring that equal-sized messages are transferred in each communication step. Furthermore, the time to compute the schedule and the index sets is significantly smaller. It takes $O(\max(P, Q))$ time and is less than 1% of the data transfer time. In comparison, the schedule computation time using the stat-of-the-art scheme (which is based on the bipartite matching scheme) is 10 to 50% of the data transfer time for similar problem sizes. Therefore, our proposed algorithm is suitable for run-time array redistribution. To evaluate the performance of our scheme, we have implemented the algorithm using C and MPI on the IBM SP2. Results show that our algorithm performs better than the previous algorithms with respect to the total redistribution time which includes the time for data transfer, schedule, and index computation.

*This work was supported in part by the US DoD High Performance Computing Modernization Office and Rome Laboratory, Air Force Materiel Command, USAF, under agreement number F30602-97-2-0016. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright annotation thereon.