

Sorting and Routing on the Array with Reconfigurable Optical Buses

S. Rajasekaran¹ and S. Sahni²

Department of CIS, University of Florida

Abstract. In this paper we present efficient algorithms for sorting and packet routing on the AROB (Array with Reconfigurable Optical Buses) model.

1 Introduction

An Array with Reconfigurable Optical Buses (AROB) [18] is essentially an $m \times n$ reconfigurable mesh in which the buses are implemented using optical technology. This model has attracted the attention of many researchers in the recent past owing to its promise in superior practical performance.

A 4×4 reconfigurable mesh is shown in figure 1. The switch in each processor can be used to connect together subsets of the four bus segments connected to the processor. Reconfigurable meshes that use electronic buses have been studied extensively. Various models such as the RN [3], RMESH [7], PARBUS [9], M_r [19], RMBM [25], and DMBC [24] have been proposed and studied.

Reconfigurable meshes with optical buses have been less extensively studied. In the AROB model of [18], the allowable switch settings of the processors are the same as those in the RN model of [3]. These are shown in figure 2. A bus link connects two adjacent processors x and y and has two associated wave guides. One of the wave guides permits an optical signal to travel from x to y and the other permits signal movement from y to x . By setting processor switches, bus links are connected together to form disjoint buses. On each bus, we need to specify which orientation of the waveguide on each link of the bus is to be

¹Research of this author is supported, in part, by an NSF Grant CCR-9596065

²Research of this author is supported, in part, by the Army Research Office under grant DAA H04-95-1-

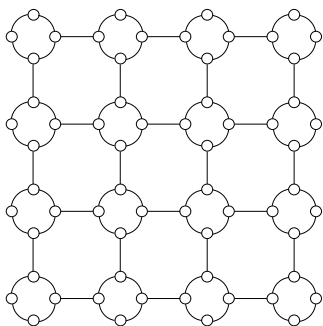


Figure 1: A 4×4 Reconfigurable Mesh

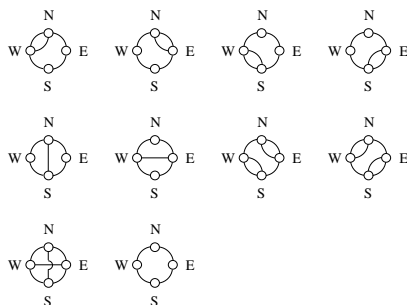


Figure 2: Possible Switch Connections

used. The resulting directed graph that represents the bus should be a directed chain. The root of this chain is the bus ‘leader’. The length of a bus is the number of links on the chain representing that bus. The position of any processor on a bus is its distance from the bus leader. The time needed to transmit a message on a bus is referred to as one cycle. A cycle is divided into slots of duration τ and each slot can carry a different optical signal. τ is the time needed for an optical pulse to move down one bus link. Pavel and Akl [18] have argued that for reasonable size meshes (say up to 1000×1000), the number of slots in a cycle may be assumed to be n for an $n \times n$ mesh. Further, the duration of a cycle may be assumed constant and comparable to the time for a CPU operation.

To assist a processor in determining which slot to use, each processor has a slot counter.

These counters may be started at the beginning of a cycle. The bus leader initiates a light pulse at this time (i.e., it writes a one to the bus). The counter at each processor stops when the light pulse reaches that processor. This special timing mechanism does not require any bus read operation. The terminal counter value is the distance of the processor from the bus leader. Note that because a processor can read/write from/to its bus during only one slot of a cycle, it cannot poll the up to n light pulses moving through it in one cycle. An additional AROB feature that facilitates the development of algorithms is the delay unit at each processor. This permits a processor to introduce a one time slot delay in the light pulses passing through it.

A linear AROB (LAROB) is a $1 \times n$ AROB [18].

In Section 2 we provide some preliminaries and survey known results in the area of AROBs. Sections 3 and 4 are devoted to the problems of sorting and packet routing, respectively. In Section 5 we provide our conclusions and list some open problems.

2 Preliminaries

In this section we provide some preliminary facts and results that will be employed in the paper.

2.1 Problem Definitions

Given a sequence of numbers, say, k_1, k_2, \dots, k_n , the problem of sorting is to rearrange them in nondecreasing order.

In any fixed connection network, a single step of interprocessor communication can be thought of as a packet routing task. The problem of routing can be stated as follows: There is a packet of information at each node that is destined for some other node. Send all the packets to their correct destinations as quickly as possible making sure that at most one packet crosses any edge at any time. Packet routing is equivalent to the random access write operation first defined by Nassimi and Sahni [14]. The *run time* of any packet routing algorithm is defined to be the time taken by the last packet to reach its destination. The *queue size* is the maximum number of packets that any processor will have to store during the algorithm.

The problem of *partial permutation routing* is the task of routing where at most one packet originates from any node and at most one packet is destined for any node. Any routing problem where at most h packets originate from any node and at most h packets are destined for any node will be called *$h - h$ routing* or *h -relations* [26].

2.2 Previous Results and Extensions

In [18], the AROB model has been defined. Similar models have been employed before as well (see e.g., [17]). A related model known as the *Optical Communication Parallel Computer (OCPC)* has also been defined in the literature (see e.g., [1], [5], [26], [6]). In an OCPC any processor can communicate with any other processor in one unit of time, provided there are no conflicts. If more than one processors try to send a message to the same processor, no message reaches the intended destination.

In [18], algorithms for such problems as prefix computation, routing on a linear array, matrix multiplication, etc. have been given for the AROB. On the other hand, [17] considers the problem of selection on a mesh with optical buses.

Lemma 2.1 *Consider an n processor LAROB. If each processor has a bit, then the prefix sums of these bits can be computed in $O(1)$ cycles [18].*

The idea behind the above the algorithm is as follows: Processor 1 initiates a light pulse in time slot one of a cycle if its bit is zero and in slot two otherwise. All processors start their counters at the start of the cycle and also set their delay units to introduce a one slot delay in case the processor's bit is one. A processor's counter is turned off when the light pulse initiated by processor 1 reaches it. By using the terminal counter value, its data bit and its distance from processor 1, each processor can compute its prefix sum value.

The above algorithm can be extended to show the following [18]:

Lemma 2.2 *The addition of $n \log n$ -bit numbers can be performed in $O(1)$ cycles on a $\log n \times n$ AROB.*

A constant time algorithm for prefix sums on a 2D AROB can be found in [18]. In particular, they show:

Lemma 2.3 *If there is a bit at each node of a $\sqrt{n} \times \sqrt{n}$ AROB, we can compute the prefix sums of these bits in $O(1)$ cycles.*

This algorithm uses the algorithm of [16] to compute the prefix sums of an integer sequence. The maximum bus length employed by this algorithm is $3\sqrt{n}$.

We next show that one can reduce the bus length for the prefix sums problem to \sqrt{n} (which is a factor of 3 improvement over [18]’s algorithm). Furthermore, our algorithm is simpler.

Lemma 2.4 *Prefix sums of bits on a $\sqrt{n} \times \sqrt{n}$ AROB can be computed in $O(1)$ cycles, keeping the maximum bus length as \sqrt{n} .*

Proof. We proceed as in [18]. Say we are interested in computing the prefix sums in row major order. Using the algorithm of lemma 2.1, we can compute the prefix sums along each row in $O(1)$ time. At the end of this step, each processor in the last column has the sum of all 1’s in the corresponding row. Now the original problem of computing prefix sums reduces to computing prefix sums of \sqrt{n} numbers where each number is at most \sqrt{n} (i.e. each number is an $O(\log n)$ -bit number). Next we describe how to compute the prefix sums of \sqrt{n} $O(\log n)$ -bit numbers in $O(1)$ time on a $\sqrt{n} \times \sqrt{n}$ AROB. (In [18], this is done using the algorithm of [16].)

Step 1: Group the numbers with $\log n$ numbers in each group. Allocating a subarray of size $\log n \times \log n$, compute the sum of numbers in each group. For each group, first reduce the problem to that of adding $\log n$ $O(\log \log n)$ -bit numbers (by summing the corresponding bits using lemma 2.1) and then apply lemma 2.2. Maximum bus length is $O(\log n)$.

Step 2: Compute prefix sums of the $\frac{\sqrt{n}}{\log n}$ group sums. This can be done using lemma 2.2 as follows. For each prefix sum (there are $\frac{\sqrt{n}}{\log n}$ of them) allocate a subarray of size $\log n \times \frac{\sqrt{n}}{\log n}$. Broadcast the appropriate numbers to each subarray and within each subarray add the numbers using lemma 2.2. Here also, the maximum bus length is \sqrt{n} .

Step 3: Compute prefix sums local to each group of $\log n$ numbers. This step is similar to step 2. Each group will get a subarray of size $\log n \times \log^2 n$. Maximum bus length is $O(\log n)$.

Clearly, the above algorithm runs in $O(1)$ time. \square

Lemma 2.5 *In a LAROB of size n any permutation can be routed in $O(1)$ cycles [18].*

The time it takes for a packet to move from one processor to the next is assumed to be τ . Consider any permutation to be routed. Let the processors be numbered $1, 2, \dots, n$ starting from left. There is a time slot assigned to each processor for reading from (and writing into) the bus. Let the reading time slot for processor i be $2i$. Processor 1 creates a ‘time slot’ for each packet that moves one edge per τ time. The time slots created will be in the order of the processors, i.e., the first time slot is meant for processor 1, time slot 2 is meant for processor 2, and so on. If processor p has a message for processor q , p will write this message at time $t + (p+q)\tau$, where t is the start time. Clearly, this algorithm terminates in 2 cycles (or $2n\tau$ time).

The above lemma can be strengthened as follows:

Lemma 2.6 *Let \mathcal{L} be a LAROB of size n . Consider a routing problem where $O(1)$ packets originate from any node and $O(1)$ packets are destined for any node. This problem can also be solved in $O(1)$ cycles.*

Proof. Let c (resp. d) be an upper bound on the number of packets destined for (originating from) any node. It suffices to consider the case $d = 1$, since that algorithm can be repeated d times to take care of the general case. There will be $2c$ runs of the algorithm given above for lemma 2.5. The above algorithm has the property that the message read by any processor q is the last message written in time slot q . If more than one processors wrote in time slot q , they can determine if their message was read by q or not by reversing the routing process. This way, in every two executions of the above routing algorithm, a processor receives one message destined for it. \square

A further extension of the above ideas leads to the following [18]:

Lemma 2.7 *Say there are k elements arbitrarily distributed (at most one per processor) in a two dimensional AROB of size $\sqrt{n} \times \sqrt{n}$. We would like to ‘compact’ them in the first $\lceil \frac{k}{\sqrt{n}} \rceil$ rows. This problem can be solved in $O(1)$ cycles.*

The algorithm for the above problem figures out a unique address for each element and then routes the elements using greedy paths. There is no possibility of a collision.

Many $O(1)$ time algorithms are known for sorting on the reconfigurable mesh (e.g., [9],[15]):

Lemma 2.8 *Sorting of n numbers can be performed in $O(1)$ time on a reconfigurable mesh of size $n \times n$.*

The same algorithm runs on the AROB preserving the run time.

Routing on the OCPC. Several packet routing algorithms for the OCPC model can be found in the literature. Anderson and Miller have shown that a special case of $\log n$ -relations on an n -node OCPC can be routed in $O(\log n)$ time [1]. Also, [26] and [5] have presented efficient algorithms for h -relations. An algorithm for arbitrary h -relations with a run time of $O(h + \log \log n)$ has been given by Goldberg, Jerrum, Leighton, and Rao [6]. Recently, Rajasekaran and Sahni [22] have presented an $O(h)$ time (for any h) algorithm for h -relations routing on the AROB model. All of these algorithms are randomized.

2.3 New Results

In this paper we present a sorting algorithm that can sort n general keys in $O(1)$ time on an AROB of size $n^\epsilon \times n$ for any constant $\epsilon > 0$. We also point out that this algorithm is optimal. We also present a sorting algorithm that can sort n k -bit numbers in $O(k)$ time on a LAROB of size n . Notice that such an algorithm cannot be devised even on the CRCW PRAM.

An important class of routing problems known as *BPC permutations* have been proven to be widely applicable in many applications of interest [13]. We present algorithms for routing *BPC permutations* in $O(1)$ cycles on an $n \times n$ AROB. In addition, we give a deterministic algorithm for h -relations that runs in $O(h \log n)$ cycles on a $\sqrt{n} \times \sqrt{n}$ AROB as well as on an n -node LAROB.

3 Sorting on the AROB

In this section we present optimal algorithms for sorting both general and integer keys on the two dimensional AROB. The general sorting algorithm sorts n numbers in an AROB of size $n^\epsilon \times n$ for any constant $\epsilon > 0$. The run time is $O(1)$ cycles and hence the algorithm is optimal in view of the following lower bound:

Lemma 3.1 *Sorting of n numbers needs $\Omega\left(\frac{\log n}{\log(1+\frac{p}{n})}\right)$ time using p parallel comparison tree processors [4].*

This lower bound implies that if sorting of n numbers has to be done in $O(1)$ time, then there must be $\Omega(n^{1+\epsilon})$ processors, for some constant $\epsilon > 0$. In [4], the lower bound has been proven for the parallel comparison tree model of Valiant. Since a parallel comparison tree can simulate an AROB step per step, the same lower bound applies to the AROB as well.

Our integer sorting algorithm runs on a LAROB of size n and can sort n k -bit numbers in $O(k)$ time. Notice that even on the CRCW PRAM model, such an algorithm cannot be devised in view of the lower bound result of Beame and Hastad [2].

3.1 General Sorting

Let k_1, k_2, \dots, k_n be the n given numbers. Think of these numbers of as forming a matrix M with $r = n^{2/3}$ rows and $s = n^{1/3}$ columns. We employ the column sort algorithm of Leighton [11]. There are 7 steps in the algorithm:

Algorithm Sort

1. Sort the columns of M in increasing order;
2. Transpose the matrix preserving the dimension as $r \times s$. In particular, pick the elements in column major order and fill the rows in row major order;
3. Sort each column in increasing order;
4. Rearrange the numbers applying the reverse of the permutation employed in step 2;
5. Sort the columns in a way that adjacent columns are sorted in reverse order;

6. Apply two steps of odd-even transposition sort to the rows. Specifically, in the first step perform a comparison-exchange between processors $2i + 1$ and $2i$, for $i = 0, 1, \dots$ and in the second step perform a comparison-exchange between processors $2i$ and $2i + 1$, for $i = 1, 2, \dots$; and
7. Sort each column in increasing order. At the end of this step, it can be shown that, the numbers will be sorted in column major order.

Implementation on the AROB. The n given numbers will be stored in the first row of the $n^\epsilon \times n$ AROB one key per processor. At any given time each key will know which row and which column of the matrix M it belongs to. Whenever we need to sort the columns, we will make sure that the numbers belonging to the same column will be found in successive processors.

On a LAROB of size n note that any permutation can be performed in $O(1)$ time. This means that steps 2 and 4 can be performed in $O(1)$ time. Step 6 can be performed in $O(1)$ time as well as follows: Rearrange the numbers such that elements in the same row are in successive processors and apply two steps of the odd-even transposition sort. After this, move the keys to where they came from.

Next we describe how we implement steps 1, 3, 5, and 7. We first assume that we have an AROB of size $n^{2/3} \times n$. Later we will indicate how to reduce the size to $n^\epsilon \times n$ for any $\epsilon > 0$.

Partition the AROB into $n^{1/3}$ parts each of size $n^{2/3} \times n^{2/3}$, each part corresponding to a column of M . Rearrange the n given numbers such that the first column of M is in the first $n^{2/3}$ processors of row 1; the second column is in the next $n^{2/3}$ processors of the first row; and so on. Now sort the numbers in each part (i.e., each column of M) using lemma 2.8. This can be done in $O(1)$ time. This implies that steps 1, 3, 5, and 7 of column-sort can be performed in $O(1)$ time. Therefore it follows that n numbers can be sorted in $O(1)$ time on an AROB of size $n^{2/3} \times n$.

We can reduce the size of the AROB to $n^{4/9} \times n$ as follows: We still use Leighton's sort with $r = n^{2/3}$ and $s = n^{1/3}$. In steps 1, 3, 5, and 7, each part of $n^{2/3}$ numbers will be sorted using an AROB of size $n^{4/9} \times n^{2/3}$. This is done using the AROB algorithm above.

In a similar way we can reduce the size to $n^{8/27} \times n$, $n^{16/81} \times n$, and so on. Thus we get the following theorem:

Theorem 3.1 *We can sort n numbers in $O(1)$ cycles using an AROB of size $n^\epsilon \times n$, where ϵ is any constant > 0 .*

3.2 Integer Sorting

In this subsection we present an algorithm for sorting n k -bit numbers in $O(1)$ time on a LAROB with n processors. This algorithm makes use of the idea of *radix sorting* and lemmas 2.1 and 2.5.

Radix Sorting. The idea is captured by the following lemma:

Lemma 3.2 *If n numbers in the range $[0, R]$ can be stable sorted using P processors in time T , then we can also stable sort n numbers in the range $[0, R^c]$ in $O(T)$ time using P processors, c being any constant.*

A sorting algorithm is said to be *stable* if equal keys remain in the same relative order in the output as they were in the input.

The algorithm proceeds as follows: There are k stages. In stage i we sort the numbers with respect to their i th LSBs. To be more specific, in the first stage we sort the numbers with respect to their LSBs. In the next stage, we apply a sort in the resultant sequence with respect to the next LSBs, and so on. Thus there will be k stages in the algorithm.

Each stage can be performed in $O(1)$ time as follows: Notice that each stage is nothing but sorting n 1-bit numbers. Perform a prefix sums operation for the zeros in the input. Do the same for the 1's in the input. Using these two sums, each processor can determine the position of its data in the sorted list. Permute the data to complete the sort for the stage. Since the prefix sums as well as the permutation take $O(1)$ time each, each stage takes $O(1)$ time as well (c.f. lemmas 2.1 and 2.5.)

Thus we have proven the following:

Theorem 3.2 *A LAROB with n processing elements can sort n k -bit numbers in $O(k)$ cycles.*

Realize that no PRAM algorithm can achieve the above performance, since the lower bound theorem of [2] implies that sorting of n bits on the CRCW PRAM will need $\Omega(\frac{\log n}{\log \log n})$ time, given only a polynomial number of processors.

4 Packet Routing

Packet routing is a fundamental problem of parallel computing since algorithms for packet routing can be used as mechanisms for interprocessor communication. In this section we present efficient algorithms for packet routing on the AROB.

For the OCPC model several routing algorithms are known: 1) Anderson and Miller have shown that a special case of $\log n$ -relations on an n -node OCPC can be routed in $O(\log n)$ time [1]; 2) Valiant extended their algorithm to show that any h -relation can be achieved in time $O(h + \log n)$ [26]; 3) Geréb-Graus and Tsantilas have given a simple algorithm that can route any h -relation in time $O(h + \log n \log \log n)$ [5]; 4) A more complicated algorithm with a run time of $O(h + \log \log n)$ has been given by Goldberg, Jerrum, Leighton, and Rao [6]. Recently, Rajasekaran and Sahni [22] show how to perform h -relations routing in $O(1)$ cycles on a LAROB (for any h). The stated time bounds hold with high probability.

There is a crucial difference between the OCPC model and the AROB model. On the OCPC model if more than one messages are sent to some processor π at the same time, none of them reaches π . On the other hand, under the same scenario, one of the messages will reach π on the AROB model. Also, operations such as prefix sums (limited to integers of certain magnitude) and compaction can be performed in $O(1)$ time on the AROB model and not on the OCPC model.

4.1 BPC Permutations

Definition [BPC Permutations][13]: In a network \mathcal{N} with N nodes, $N!$ permutations are possible. An important subset of these permutations is called *Bit Permute and Complement (BPC)*. Assume that N is an integral power of 2 (i.e., $N = 2^p$ for some integer p). Each node of \mathcal{N} can be labeled with a p -bit sequence $a_{p-1}, a_{p-2}, \dots, a_1, a_0$. Any BPC permutation, π can be described with a vector $(\pi_{p-1}, \pi_{p-2}, \dots, \pi_1, \pi_0)$, where $\pi_i \in \{\pm 0, \pm 1, \dots, \pm(p-1)\}$ and $(|\pi_{p-1}|, |\pi_{p-2}|, \dots, |\pi_1|, |\pi_0|)$ is a permutation of $(0, 1, 2, \dots, p-1)$. Under this permutation, if

the origin of a packet is $a_{p-1}, a_{p-2}, \dots, a_1, a_0$ then its destination will be $d_{p-1}, d_{p-2}, \dots, d_1, d_0$, where $d_{|\pi_i|} = a_i$ if π_i is non negative and $d_{|\pi_i|} = \bar{a}_i$ otherwise.

Consider a network \mathcal{N} with 8 nodes, for example. Let the permutation under concern be $\pi = (-2, 0, 1)$. Under this permutation a packet originating from the node (a_2, a_1, a_0) will be destined for (\bar{a}_2, a_0, a_1) . For instance a packet from node $(1, 1, 0)$ will have $(0, 0, 1)$ as its destination.

Many important permutations such as matrix transpose, bit reversal, perfect shuffle, etc. belong to the class of BPC permutations. For example the perfect shuffle can be characterized with the vector $(0, p-1, p-2, \dots, 2, 1)$. The vector $(p/2-1, p/2-2, \dots, 1, 0, p-1, p-2, \dots, p/2)$ corresponds to matrix transpose. Bit reversal is described by $(0, 1, 2, \dots, p-1)$, and so on. The number of permutations in the BPC class is $2^p p!$ where the network size is $N = 2^p$.

4.2 Matrix Transpose Routing

In this section we show how to perform matrix transpose routing in $O(1)$ cycles on an $n \times n$ AROB. Recall that the matrix transpose permutation is a BPC permutation characterized by the vector $(p/2-1, p/2-2, \dots, 1, 0, p, p-1, \dots, p/2)$. If a packet originates at the node (i, j) ($1 \leq i, j \leq n$) in the 2D mesh, then its destination is (j, i) .

To perform this routing in $O(1)$ cycles, we use a two phase algorithm. In the first phase, we connect the switches of nodes as shown in Figure 4.2(a). The buses formed are along the diagonals. The switch connection for any processor (i, j) is SE if $i+j$ is odd and it is NW if $i+j$ is even. Note that every node of the mesh is on at most one bus. Also, if the node (i, j) is on some bus, then (j, i) will also be on the same bus. Perform routing along each bus thus formed. All the packets that are on some bus or the other would have been successfully routed to their correct destinations at the end of this phase. Packets that have not been processed in the first phase are handled in the second phase.

In the second phase connect the switches as shown in Figure 4.2(b) and perform routing along the buses. At the end, every packet would have reached its correct destination. Since routing along a LAROB can be performed in $O(1)$ cycles (c.f. Lemma 2.6), matrix transpose can also be completed in $O(1)$ cycles. Thus we get the following Lemma:

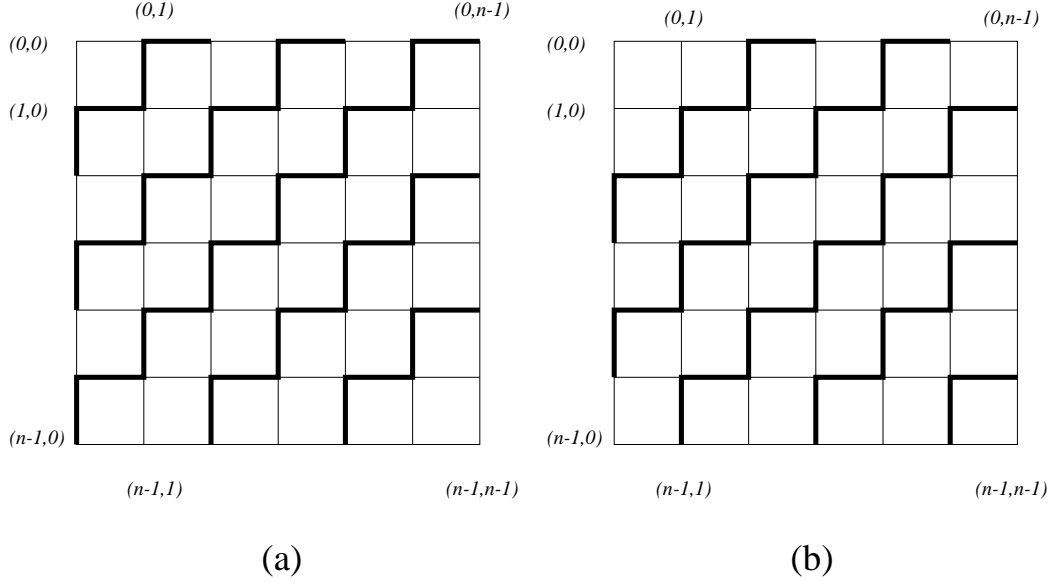


Figure 3: Matrix transpose

Lemma 4.1 *Matrix transpose can be performed on an $n \times n$ AROB in $O(1)$ cycles. \square*

The above Lemma will prove helpful in devising routing algorithms for general BPC permutations.

4.3 BPC Permutation Routing

In this section we present our $O(1)$ cycles algorithm for general BPC permutations. The basic idea is to decompose any BPC permutation into a sequence of five permutations. The first and the fourth permutations in the decomposition are such that they correspond to data movements only along the columns. The second and the fifth permutations are along the rows. Using Lemma 2.6, these four permutations can be performed in $O(1)$ cycles. On the other hand, the third permutation is a matrix transpose which can also be realized in $O(1)$ cycles (c.f. Lemma 4.1). More details follow. Assume w.l.o.g. that p is even.

Let $\pi = (\pi_{p-1}, \pi_{p-2}, \dots, \pi_{p/2}, \pi_{p/2-1}, \dots, \pi_1, \pi_0)$ be the BPC permutation to be routed. Denote the first (second) half of the vector as P_1 (resp. P_0). In other words, $P_1 = \pi_{p-1}, \pi_{p-2}, \dots, \pi_{p/2}$ and $P_0 = \pi_{p/2-1}, \dots, \pi_1, \pi_0$. Let k be the number of symbols in the sequence P_1 that belong to $\{0, 1, 2, \dots, p/2 - 1\}$ and let $\pi_{i_1}, \pi_{i_2}, \dots, \pi_{i_k}$ be these symbols.

Clearly, the number of symbols in P_0 that belong to $\{p/2, p/2 + 1, \dots, p - 1\}$ will also be k . Let these symbols be $\pi_{j_1}, \pi_{j_2}, \dots, \pi_{j_k}$.

Let $A = (a_{p-1}, a_{p-2}, \dots, a_1, a_0)$ be any node and q be the packet originating from A . The five phases in the algorithm are described below:

- **Phase I.** Let $A = A_1, A_0$ where A_1 is the first half of A and A_0 is the second half of A . Also let $A'_1 = a_{i_1}, a_{i_2}, \dots, a_{i_k}$ and $A'_0 = a_{j_1}, a_{j_2}, \dots, a_{j_k}$. Route the packet from A to $A' = A_1 - A'_1, A'_1, A_0$. Here $A_1 - A'_1$ is the subsequence that results from A_1 after eliminating symbols of A'_1 . This task involves routing along the columns. As a part of this phase take care of complements in A_1 , if any. I.e., if $a_i \in A_1$ and π_i is negative, then a_i will appear in A' as \bar{a}_i (in an appropriate place).
- **Phase II.** Now route the packet q from its current location to $A_1 - A'_1, A'_1, A_0 - A'_0, A'_0$. This task involves routing along rows. Any complements in A_0 can also be taken care of.
- **Phase III.** Partition the mesh into submeshes of size $2^k \times 2^k$ and employ matrix transpose within the submeshes. As a result q will reach the node $A_1 - A'_1, A'_1, A_0 - A'_0, A'_0$.
- **Phase IV.** This phase is analogous to phase I and here routing is done along the columns. At the end, q will be in a node the first half of whose label is the same as the first half of the final destination of q . In other words, q will be in its destination row.
- **Phase V.** Finally, a routing step along the rows takes q to its desired destination.

Example. Let \mathcal{N} be a network with $N = 2^8$ nodes. Let the BPC permutation under concern be $(6, -3, -4, 1, 0, -2, 5, 7)$. If the packet q originates from $(a_7, a_6, a_5, a_4, a_3, a_2, a_1, a_0)$, $A_1 = a_7, a_6, a_5, a_4$; $A_0 = a_3, a_2, a_1, a_0$; $A'_1 = a_6, a_4$; and $A'_0 = a_1, a_0$.

In the first phase, q goes to $a_7, \bar{a}_5, \bar{a}_6, a_4, a_3, a_2, a_1, a_0$. In the second phase it goes to $a_7, \bar{a}_5, \bar{a}_6, a_4, a_3, \bar{a}_2, a_1, a_0$. After the matrix transpose in phase III, q will reach $a_7, \bar{a}_5, a_1, a_0, a_3, \bar{a}_2, \bar{a}_6, a_4$.

Phase IV takes q to $a_0, a_7, a_1, \bar{a}_5, a_3, \bar{a}_2, \bar{a}_6, a_4$. And finally at the end of phase V, q will reach $a_0, a_7, a_1, \bar{a}_5, \bar{a}_6, \bar{a}_2, a_4, a_3$. \square

Phases I, II, IV, and V take $O(1)$ cycles each (c.f. Lemma 2.6). Phase III takes $O(1)$ cycles as well in accordance with Lemma 4.1. Thus the whole algorithm takes $O(1)$ cycles. Note that the routing we perform in each phase is a permutation, i.e., there are no conflicts among the packets. We get the following theorem:

Theorem 4.1 *Any BPC permutation can be completed in $O(1)$ cycles on a 2D AROB. \square*

4.4 Routing on a LAROB

Let \mathcal{L} be a LAROB with n processors. We are interested in routing an arbitrary h -relation. Notice that a special case where $h = O(1)$ has already been considered in lemma 2.6.

We look at some special cases of routing before dealing with the general case.

Problem 1. In an n -node network there are at most k packets at any node. Let N be the total number of packets. The problem is to do a load balancing, i.e., to rearrange the packets such that each node has at most $\lceil \frac{N}{n} \rceil$ packets.

Lemma 4.2 *Problem 1 can be solved in $O(k)$ cycles on an n -node LAROB. The same problem can also be solved in $O(k)$ cycles on a $\sqrt{n} \times \sqrt{n}$ AROB.*

Proof. We give the proof for a LAROB. The same proof can be extended to a 2D AROB also. Let the processors order their packets from 1 to k . Perform a prefix sums computation for the first packets of all the processors (using lemma 2.1) and compute a unique address for each such packet. Route the first packets. Prefix takes $O(1)$ cycles and so does the routing. Likewise process the second packets, the third packets, and so on. One should make sure, for example, that if q is the number of first packets, then, the second packets will be routed to nodes starting from $q + 1$. Total time is clearly $O(k)$. \square

Problem 2. There are at most ℓ packets originating from any node of a network \mathcal{N} . Also, at most k packets are destined for any node. Route the packets.

Lemma 4.3 *Problem 2 can be solved on an n -node LAROB in $O(k\ell)$ cycles.*

Proof. The packets are routed in cycles. In any cycle, a processor will choose one of its remaining packets (if any) and try to send it. It may not be successful in one attempt. If it

succeeds, it takes up the next packet; otherwise it will try to send the same packet. A packet will not reach its destination only if there is a conflict. Thus a packet can meet with failure in at most $k - 1$ cycles. This in turn means that every processor will be able to transmit all of its ℓ packets in ℓk cycles or less. \square

4.5 h -Relations Routing

We can also perform deterministic routing in an efficient manner on the AROB:

Lemma 4.4 *Any partial permutation can be routed in $O(\log n)$ time on a $\sqrt{n} \times \sqrt{n}$ AROB.*

Proof. Sort the packets into ascending order of destinations. For this, the destinations are mapped into a single number using the row major mapping scheme. The sorted packets are in processors $1, 2, \dots$ (in row major order). This sort can be accomplished in $O(\log n)$ time using a binary radix sort and two applications of lemma 4.2 to accomplish the sort on each bit (notice that the load balancing scheme of lemma 4.2 is equivalent to a stable sort of bits). Following the sort, no two packets in the same column have the same row as their destination (as there are $\sqrt{n} - 1$ packets between them). So, we may use lemma 2.5 to route packets in each column to their destination rows. Following this, no two packets in the same row have the same column as their destination. So lemma 2.5 may be used again to route packets in the same row to their destination columns. The work done following the sort takes $O(1)$ cycles. So, the overall number of cycles is $O(\log n)$. \square

An extension of the above result can also be proven:

Lemma 4.5 *Any h -relation can be routed in $O(h \log n)$ cycles on a 2D AROB of size $\sqrt{n} \times \sqrt{n}$.*

Proof. Sort the packets into nondescending order of destination. Following the sort, the packets are in the first few processors (in row major order), h packets to a processor. This is accomplished using a binary radix sort on the row major index of the packet's destination processor.

When sorting on bit k of this index, we first concentrate the packets with bit k equal to 0, h packets to a processor and then concentrate those with bit k equal to 1. The process

for each bit value is similar. Consider the case of packets with bit k equal to 0. Call these packets *selected packets*. A processor may have up to h selected packets.

The selected packets in each processor are combined to form a ‘superpacket’ of size at most h . The superpackets are compacted into processors $1, 2, \dots$ (in row major order) using the 2D compaction algorithm of [18]. Since the superpacket size is $O(h)$, this takes $O(h)$ time. The superpackets are now decomposed into the original packets. The original packets are to be further compacted so that we have h packets to a processor. Each packet in a processor is assigned a level number corresponding to its order in the processor. Level numbers are in the range 1 to h .

Prefix sums for the level i packets, $1 \leq i \leq h$ are computed using the 2D prefix sum algorithm given in section 2. The rank $r(i, j)$ of a level i packet in processor j is $\sum_{k=1}^h ps(k, j-1) + (i-1)$, where $ps(k, j-1)$ is the prefix sum of the level k packet in processor $j-1$. The processor $P(i, j)$ to which this packet is to be routed is $\lfloor r(i, j)/h \rfloor$. Furthermore, this packet will be the $round(i, j) = r(i, j) \bmod h + 1$ -th packet in this processor. Since the number of packets in each row is at most $h\sqrt{n}$, no two packets (i, j) and (k, l) , where j and l are processors in the same row, have $column(P(i, j)) = column(P(k, l))$ and $(row(P(i, j)) \neq row(P(k, l)) \text{ or } round(i, j) = round(k, l))$. As a result, the compaction may be completed as below:

Step 1: Perform h rounds of row permutation routing on each row. In round k , packets (i, j) with $round(i, j) = k$ are routed to the processor in column $column(P(i, j))$.

Step 2: Perform h rounds of column permutation routing. In round k , packets (i, j) with $round(i, j) = k$ are routed to the processor in row $row(P(i, j))$.

The radix sort described above takes $O(h \log n)$ time. To complete the h -relation we perform h rounds of column and row permutations. In round i , the level i packets in each column are first routed to the correct row using a column permutation. There can be no collision as for a collision the number of packets destined to the same row needs to be $> h\sqrt{n}$. Next, the level i packets are routed to the correct column using row permutations. Again, collisions are not possible. \square

5 Conclusions

In this paper we have presented efficient algorithms for sorting and packet routing on the AROB. We have considered both integer sorting and general sorting problems. Our general sorting algorithm is optimal. An interesting open problem is if there exists an $O(h)$ deterministic routing algorithm for the 2D AROB. Also, can our integer sorting algorithm be improved?

References

- [1] R.J. Anderson and G.L. Miller, Optical Communication for Pointer Based Algorithms, Technical Report CRI-88-14, Computer Science Department, University of Southern California, 1988.
- [2] P. Beame and J. Hastad, Optimal Bounds for Decision Problems on the CRCW PRAM, Journal of the ACM, 36(3), 1989, pp. 643-670.
- [3] Y. Ben-Asher, D. Peleg, R. Ramaswami, and A. Schuster, The Power of Reconfiguration, Journal of Parallel and Distributed Computing, 1991, pp. 139-153.
- [4] R. Bopanna, A Lower Bound for Sorting on the Parallel Comparison Tree, Information Processing Letters, 1989.
- [5] M. Geréb-Graus and T. Tsantilas, Efficient Optical Communication in Parallel Computers, Symposium on Parallel Algorithms and Architectures, 1992, pp. 41-48.
- [6] L. Goldberg, M. Jerrum, T. Leighton, and S. Rao, A Doubly-Logarithmic Communication Algorithm for the Completely Connected Optical Communication Parallel Computer, Proc. Symposium on Parallel Algorithms and Architectures, 1993.
- [7] E. Hao, P.D. McKenzie and Q.F. Stout, Selection on the Reconfigurable Mesh, Proc. Frontiers of Massively Parallel Computation, 1992.
- [8] E. Horowitz and S. Sahni, *Fundamentals of Computer Algorithms*, Computer Science Press, 1978.

- [9] J. Jang and V.K. Prasanna, An Optimal Sorting Algorithm on Reconfigurable Mesh, Proc. International Parallel Processing Symposium, 1992, pp. 130-137.
- [10] J. Jenq and S. Sahni, Reconfigurable Mesh Algorithms for Image Shrinking, Expanding, Clustering, and Template Matching, Proc. International Parallel Processing Symposium, 1991, pp. 208-215.
- [11] T. Leighton, Tight Bounds on the Complexity of Parallel Sorting, IEEE Transactions on Computers, C-34(4), 1985, pp. 344-354.
- [12] R. Miller, V.K. Prasanna-Kumar, D. Reisis and Q.F. Stout, Meshes with Reconfigurable Buses, in Proc. 5th MIT Conference on Advanced Research in VLSI, 1988, pp. 163-178.
- [13] D. Nassimi and S. Sahni, An Optimal Routing Algorithm for Mesh-Connected Parallel Computers, Journal of the ACM, 27(1), 1980, pp. 6-29.
- [14] D. Nassimi and S. Sahni, A Self-Routing Benes Network and Parallel Permutation Algorithms, IEEE Transactions on Computers, C-30(5), 1981, pp. 332-340.
- [15] M. Nigam and S. Sahni, Sorting n Numbers on $n \times n$ Reconfigurable Meshes with Buses, Proc. International Parallel Processing Symposium, 1993, pp. 174-181.
- [16] S. Olariu, J.L. Schwing and J. Zhang, Integer Problems on Reconfigurable Meshes, with Applications, Proc. 1991 Allerton Conference, 4, 1991, pp. 821-830.
- [17] Y. Pan, Order Statistics on Optically Interconnected Multiprocessor Systems, Proc. First International Workshop on Massively Parallel Processing Using Optical Interconnections, 1994, pp. 162-169.
- [18] S. Pavel and S.G. Akl, Matrix Operations using Arrays with Reconfigurable Optical Buses, manuscript, 1995.
- [19] S. Rajasekaran, Meshes with Fixed and Reconfigurable Buses: Packet Routing, Sorting and Selection, Proc. First Annual European Symposium on Algorithms, Springer-Verlag Lecture Notes in Computer Science 726, 1993, pp. 309-320.

- [20] S. Rajasekaran, Sorting and Selection on Interconnection Networks, to appear in Proc. *DIMACS Workshop on Interconnection Networks and Mapping and Scheduling Parallel Computation*, 1995.
- [21] S. Rajasekaran and J.H. Reif, Derivation of Randomized Sorting and Selection Algorithms, in *Parallel Algorithm Derivation and Program Transformation*, Edited by R. Paige, J.H. Reif, and R. Wachter, Kluwer Academic Publishers, 1993, pp. 187-205.
- [22] S. Rajasekaran and S. Sahni, Sorting, Selection, and Routing on the Array with Reconfigurable Optical Buses, submitted to *IEEE Transactions on Parallel and Distributed Systems*, 1995.
- [23] S. Rao and T. Tsantilas, Optical Interprocessor Communication Protocols, Proc. Workshop on Massively Parallel Processing Using Optical Interconnections, 1994, pp. 266-274.
- [24] S. Sahni, Data Manipulation on the Distributed Memory Bus Computer, to appear in *Parallel Processing Letters*, 1995.
- [25] R.K. Thiruchelvan, J.L. Trahan, and R. Vaidyanathan, Sorting on Reconfigurable Multiple Bus Machines, Proc. International Parallel Processing Symposium, 1994.
- [26] L.G. Valiant, General Purpose Parallel Architectures, in *Handbook of Theoretical Computer Science: Vol. A* (J. van Leeuwen, ed.), North Holland, 1990.
- [27] L.G. Valiant and G.J. Brebner, Universal Schemes for Parallel Communication, Proc. 13th Annual ACM Symposium on Theory of Computing, 1981, pp. 263-277.