

1995

# Many-to-many Personalized Communication with Bounded Traffic

Sanjay Ranka

*Syracuse University, School of Computer and Information Science, ranka@top.cis.syr.edu*

Ravi V. Shankar

*Syracuse University, School of Computer and Information Science, rshankar@top.cis.syr.edu*

Khaled A. Alsabti

*Syracuse University, School of Computer and Information Science, kaalsabt@top.cis.syr.edu*

Follow this and additional works at: [https://surface.syr.edu/lcsmith\\_other](https://surface.syr.edu/lcsmith_other)

 Part of the [Computer Sciences Commons](#)

---

## Recommended Citation

Ranka, Sanjay; Shankar, Ravi V.; and Alsabti, Khaled A., "Many-to-many Personalized Communication with Bounded Traffic" (1995). *College of Engineering and Computer Science - Former Departments, Centers, Institutes and Projects*. 15. [https://surface.syr.edu/lcsmith\\_other/15](https://surface.syr.edu/lcsmith_other/15)

This Article is brought to you for free and open access by the College of Engineering and Computer Science at SURFACE. It has been accepted for inclusion in College of Engineering and Computer Science - Former Departments, Centers, Institutes and Projects by an authorized administrator of SURFACE. For more information, please contact [surface@syr.edu](mailto:surface@syr.edu).

# Many-to-many Personalized Communication with Bounded Traffic

Sanjay Ranka      Ravi V. Shankar      Khaled A. Alsabti  
School of Computer and Information Science  
Syracuse University, Syracuse, NY 13244-4100  
e-mail: ranka, rshankar, kaalsabt@top.cis.syr.edu

## Abstract

This paper presents solutions for the problem of many-to-many personalized communication, with bounded incoming and outgoing traffic, on a distributed memory parallel machine. We present a two-stage algorithm that decomposes the many-to-many communication with possibly high variance in message size into two communications with low message size variance. The algorithm is deterministic and takes time  $2t\mu$  (+ lower order terms) when  $t \geq O(p^2 + p\tau/\mu)$ . Here  $t$  is the maximum outgoing or incoming traffic at any processor,  $\tau$  is the startup overhead and  $\mu$  is the inverse of the data transfer rate. Optimality is achieved when the traffic is large, a condition that is usually satisfied in practice on coarse-grained architectures. The algorithm was implemented on the Connection Machine CM-5. The implementation used the low latency communication primitives (active messages) available on the CM-5, but the algorithm as such is architecture-independent. An alternate single-stage algorithm using distributed random scheduling for the CM-5 was implemented and the performance of the two algorithms were compared.

## 1 Basic Communication Primitives

Communication between processors on a parallel machine can generally be described as  $x$ -to- $y$  communication where  $x$  and  $y$  can be substituted by *one*, *all*, or *many*. “Communication” implies processors sending and receiving messages:  $x$  being *one*, *all*, and *many* respectively, indicates that only one of the  $p$  processors sends data, that all processors send data, and that only some processors send data. Similarly,  $y$  being *one*, *all*, and *many* indicate that from each of the senders, one, all, and many processors receive data respectively. Communication can be further distinguished as a *broadcast/accumulation* or as *personalized communication*. For example, one-to-all

communication could be either a one-to-all broadcast (*single-node broadcast*) where a single processor sends out the same message to all processors, or a one-to-all personalized communication (*single-node scatter*) where a single processor sends out different messages to each processor. This classification is fairly standard in the literature. See, for instance, [5]. Algorithms for performing broadcasts are comparatively easier than those for performing personalized communication. All further discussion in this paper deals with personalized communication. We also restrict ourselves to *collective* communication, where there are multiple senders and multiple receivers.

## 2 Collective Communication Parameters

Any type of communication in a machine with  $p$  processors can be represented using a communication matrix, a  $p \times p$  matrix  $M$  where the addresses of the sending and receiving processors are used as row and column indices. The matrix entry  $m_{ij}$  denotes the size of the message being sent by processor  $P_i$  to processor  $P_j$ . The rows of the matrix are called send vectors and the columns are called receive vectors. The outgoing traffic  $r_i$  is the sum of the sizes of the messages being sent by processor  $P_i$ , while the incoming traffic  $c_j$  is the sum of the sizes of the messages being received by processor  $P_j$ . The outgoing traffic bound  $r$  is the maximum outgoing traffic at any processor, and the incoming traffic bound  $c$  is the maximum incoming traffic at any processor. The overall traffic bound  $t$  is the maximum traffic, incoming or outgoing, at any processor.

$$r_i = \sum_j m_{ij} \quad c_j = \sum_i m_{ij}$$

$$r = \max_i r_i \quad c = \max_j c_j \quad t = \max(r, c)$$

The fan-out  $f_i$  is the number of messages sent by processor  $P_i$ , while the fan-in  $g_j$  is the number of mes-

|       | $P_0$ | $P_1$ | $P_2$ | $P_3$ | $P_4$ | $P_5$ | $P_6$ | $P_7$ |    |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|----|
| $P_0$ |       | 3     | 1     |       | 2     | 1     | 2     | 1     | 10 |
| $P_1$ | 1     |       | 2     | 2     | 1     | 1     | 3     |       | 10 |
| $P_2$ | 4     | 1     |       | 2     |       | 2     |       | 1     | 10 |
| $P_3$ |       | 2     |       |       |       | 3     | 1     | 4     | 10 |
| $P_4$ | 3     |       | 4     |       |       | 2     | 1     |       | 10 |
| $P_5$ | 1     | 2     | 1     | 2     |       |       |       | 4     | 10 |
| $P_6$ |       | 2     | 1     |       | 7     |       |       |       | 10 |
| $P_7$ | 1     |       | 1     | 4     |       | 1     | 3     |       | 10 |
|       | 10    | 10    | 10    | 10    | 10    | 10    | 10    | 10    |    |

Figure 1: A matrix illustrating all-to-many communication with equal traffic

|       | $P_0$ | $P_1$ | $P_2$ | $P_3$ | $P_4$ | $P_5$ | $P_6$ | $P_7$ |    |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|----|
| $P_0$ |       | 3     | 1     |       |       | 1     |       | 1     | 6  |
| $P_1$ | 1     |       | 2     | 2     |       | 1     | 3     |       | 9  |
| $P_2$ | 4     | 1     |       | 2     |       | 2     |       | 1     | 10 |
| $P_3$ |       | 2     |       |       |       | 3     | 1     | 4     | 10 |
| $P_4$ | 3     |       | 4     |       |       | 2     | 1     |       | 10 |
| $P_5$ | 1     | 2     | 1     | 2     |       |       |       | 4     | 10 |
| $P_6$ |       |       |       |       |       |       |       |       | 0  |
| $P_7$ | 1     |       | 1     | 4     |       | 1     |       |       | 7  |
|       | 10    | 8     | 9     | 10    | 0     | 10    | 5     | 10    |    |

Figure 2: A matrix illustrating many-to-many communication with bounded traffic

sages received by processor  $P_j$ . The fan-out bound  $f$  is the maximum fan-out at any processor, and the fan-in bound  $g$  is the maximum fan-in at any processor. The overall fan-in/fan-out bound  $h$  is the maximum number of messages, being sent or received, at any processor.

$$f_i = \sum_j \text{sgn}(m_{ij}) \quad g_j = \sum_i \text{sgn}(m_{ij})$$

$$f = \max_i f_i \quad g = \max_j g_j \quad h = \text{maximum}(f, g)$$

The  $\text{sgn}$  function returns +1,0,-1 depending on whether its argument is positive, zero, or negative.

The relation between the parameters just defined and the different kinds of collective communication is as follows. If  $f_i = p$  for all  $i$  ( $0 \leq i < p$ ), the communication is *all-to-all*. This also implies that  $g_j = p$  for all  $j$  ( $0 \leq j < p$ ). If  $f_i > 0$  for all  $i$ , the communication is *all-to-many*. If  $f_i = 0$  or  $f_i = p$  for each  $i$ , and  $g_j > 0$  for each  $j$ , the communication is *many-to-all*. The general case, where  $f_i \geq 0$  for all  $i$  is *many-to-many* communication.

Collective communication can be further classified based on the sizes of the messages being sent and received. Messages could be *uniform* (of the same size) or *non-uniform* (of different sizes). The variance in message size is an important factor that affects the performance of an algorithm for collective communication. Most algorithms presented in the literature deal only with all-to-all communication with uniform message sizes. In many-to-many communication with bounded traffic  $t$ , message size could vary between 0 and  $t$ . Collective communication with bounded traffic is illustrated in the communication matrices shown in figures 1 and 2. The entry beyond the right margin of row  $P_i$  is the outgoing traffic  $r_i$ , while the entry below column  $P_j$  gives the incoming traffic  $c_j$ . Figure 1 illustrates the case of equal incoming and outgoing traffic at each processor, a special case which will be considered in the description of the algorithms.

Many-to-many communication with an overall traffic bound of  $t$ , cannot be done in time less than  $O(t)$ . The various algorithms presented in this paper take  $O(t)$  time and are optimal under specified conditions. Many-to-many communication with bounded incoming and outgoing traffic appears in a wide variety of parallel algorithms such as matrix transpose on a rectangular grid, in the final phase of sorting algorithms like sample sort, in transformations between any two distributions (like block, cyclic, and block-cyclic) that distribute data equally among all processors, etc. We are using them for performing dynamic permutations [9] and for dealing with highly irregular data accesses involving hot-spots [10] on coarse-grained parallel machines. A detailed version [8] of this paper extends the many-to-many personalized communication algorithm to deal with differing incoming and outgoing traffic bounds.

### 3 CM-5 System Overview

The Connection machine Model CM-5 [11] is a synchronized MIMD distributed-memory parallel machine available in configurations of 32 to 1024 processing nodes. Each node contains a 33 MHz SPARC microprocessor with 32 megabytes of memory, and is rated at 22 Mips and 5 Mflops. Four optional floating-point vector units can be added to each node, and this increases the node's peak performance to 128 Mips and 128 Mflops.

The CM-5 interconnection network has three components: a data network, a control network, and a diagnostic network. The data network has a fat-tree topology and provides high-performance data communication between the system components. The net-

work has a peak bandwidth of about 5 megabytes per second for node-to-node communication. However, if the destination is within the same cluster of 4 or 16 nodes in the fat-tree, a peak bandwidth of 20 megabytes per second and 10 megabytes per second, respectively, can be achieved [11]. The control network handles operations requiring the cooperation of many or all processors. This includes broadcasting, combining, and global operations. The diagnostic network helps in the detection and isolation of errors throughout the system. Both, the control network and the diagnostic network, have a binary tree topology.

Our implementations were performed on a 32-node CM-5 using active messages for low latency communication. Each 20-byte active message packet can carry up to 16 bytes of payload. Sending and receiving a single-packet active message on the CM-5 takes 1.6  $\mu$ s and 1.7  $\mu$ s respectively [4]. We used the CMMD message passing library and CMAML (the CMMD active messages layer) [12]. Two other implementations of active messages on the CM-5 exist: the original CMAM library [4] from UC Berkeley and the Strata library from MIT [2].

The time taken to send a message from one node on the CM-5 to another can be modeled as  $O(\tau + \mu M)$ , where  $\tau$  is the overhead,  $\mu$  is the transfer rate and  $M$  is the size of the message. As mentioned earlier, the value of  $\mu$  depends on whether the destination belongs to a specific subgroup and whether other nodes are sending messages. For our complexity analysis we will assume that  $\tau$  and  $\mu$  are constant, independent of the congestion and distance between two nodes.

#### 4 Collective Communication with Low Message Size Variance

The simplest version of collective communication involves all processors exchanging messages of the same size  $s$ . This is all-to-all personalized communication with uniform messages. Under these conditions, a linear permutation algorithm [1] can be used to perform the communication. A linear permutation algorithm goes through  $p - 1$  iterations, and in iteration  $k$  processor  $P_i$  ( $0 \leq i < p, 0 < k < p$ ) exchanges data with processor  $P_{i \oplus k}$  ( $\oplus$  is the bitwise exclusive OR operator). The time complexity of linear permutation is  $O(sp)$ .

A slightly modified linear permutation algorithm can be used when the messages are not uniform but exhibit only a small variation in size. Here, processors no longer send messages of exactly the same length.

---

#### Linear Permutation

```

For all processors  $P_i, 0 \leq i \leq p - 1$ , in parallel do
    Generate receive vector  $recv$  from the send vectors  $send$  in all the processors;
    for  $k = 1$  to  $p - 1$  do
         $j = i \oplus k$ ;
        if  $send^j > 0$  then  $P_i$  sends a message of size  $send^j$  to  $P_j$ 
        if  $recv^j > 0$  then  $P_i$  receives a message of size  $recv^j$  from  $P_j$ 
        Barrier synchronize with all processors;
    endfor

```

---

Figure 3: Modified Linear Permutation Algorithm

Instead they exchange send and receive vectors, and exchange only messages of the required lengths. The algorithm is shown in figure 3, where  $send^j$  is the size of the message sent to processor  $P_j$  (from  $P_i$ ) and  $recv^j$  is the size of the message received from processor  $P_j$  (by  $P_i$ ). The implicit synchronization in the linear permutation algorithm is replaced by an explicit barrier synchronization, and the algorithm retains the deterministic time complexity of  $O(sp)$  where  $s$  is the upper bound on the sizes of the messages exchanged. The barrier also prevents the communication network from getting congested and this has been shown to improve performance [3].

#### 5 Collective Communication with High Message Size Variance

Dealing with communication in which message sizes show a large variation is a difficult problem. A linear permutation algorithm could take as much as  $O(tp)$  time. Sorting messages by size is not guaranteed to improve performance either. We use a distributed random scheduling algorithm using spin locks to deal with such a situation. The distributed scheduling algorithm [13] was chosen over other graph based techniques because its low overhead enables scheduling to be done dynamically.

The algorithm is presented in figure 4. Each processor maintains a status bit that indicates whether the processor is busy or free. Processors which have messages to send perform a test-and-set operation to determine whether the intended destination is free.

---

## Distributed Random Scheduling

For all processors  $P_i$ ,  $0 \leq i \leq p - 1$ , *in parallel do*

- Generate receive vector *recv* from the send vectors *send* in all the processors;
- Pre-allocate receiving buffers according to receive vector *recv*;
- Repeat*
  - Select a destination node from send vector *send*, use active messages to test-and-set destination node's *busy\_lock*;
  - If the destination node is free to receive message,
    - Send message to the destination node;
    - Upon completion, reset destination node's *busy\_lock* to free;
    - Reset the corresponding entry in send vector *send*;
- Until* send vector *send* is empty
- Wait until all incoming messages arrive at their proper buffers.

---

Figure 4: Distributed Random Scheduling Algorithm

If the destination is free, its status bit is set to busy, and data is transferred as a single message. If the destination is busy, the sending processor tries another intended destination using the same procedure. The test-and-set inquiry operation is shown in figure 5.

We re-implemented the distributed scheduling algorithm using active messages on the CM-5. Two improvements were incorporated into the new implementation, which also replaced the interrupts in the earlier implementation with polling. First, a busy destination processor when replying to the sender of an inquiry gives a measure of how busy it is. The sender notes down this measure and makes sure that the destination will not be disturbed for this measure of time. If the sender receives busy signals from all the intended destinations, it goes to sleep for the amount of time indicated by the minimum of the measures returned by the destinations. The second improvement allowed busy destination processors to give the go-ahead for a new message transfer when the current message transfer is about to get over.

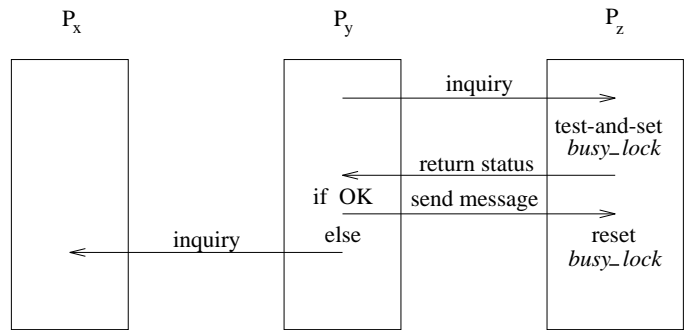


Figure 5: The inquiry operation in Distributed Scheduling

## 6 Two-stage Algorithm for Bounded Traffic Collective Communication

We have developed a two-stage algorithm that decomposes collective communication with high message size variance, into two collective communication stages with low message size variance. In the general case, the fan-out and fan-in at each processor is less than or equal to  $p$  and the traffic bound is  $t$ . Results are given separately for the equal traffic case, where the incoming traffic and the outgoing traffic at each processor is exactly equal to the overall traffic bound  $t$ . Each processor takes on three roles in this two-stage algorithm. First, each processor  $P_i$  for which the fan-out  $f_i$  is non-zero acts as a *source* processor, sending out data during the first stage. Second, each processor participates as an *intermediary*, receiving data during the first stage, and sending data during the second stage. Third, each processor  $P_j$  for which the fan-in  $f_j$  is non-zero acts as a *destination* processor, receiving data during the second stage. The organization of data in the source, intermediate and destination processors is shown in figure 6.

### 6.1 The First Stage

#### Local pre-processing

In source processor  $P_i$  ( $0 \leq i < p$ ) let  $a_{i0}, a_{i1}, \dots, a_{i(p-1)}$  be the number of elements being sent to destination processors  $P_0, P_1, \dots, P_{p-1}$  respectively. In stage 1, each of the  $a_{ij}$  elements is divided into  $p$  parts (each of size either  $\lceil a_{ij}/p \rceil$  or  $\lfloor a_{ij}/p \rfloor$ ) to be sent to processors  $P_0$  to  $P_{p-1}$ .<sup>1</sup> At the end of stage 1, processor  $P_k$  acting

---

<sup>1</sup>In reality, this is only  $p - 1$  messages, since one of the messages is to be sent to the sending processor itself. Our implementations take this into account, but this paper, for the sake of simplicity, continues to refer to  $p$  as the number of messages being sent out.

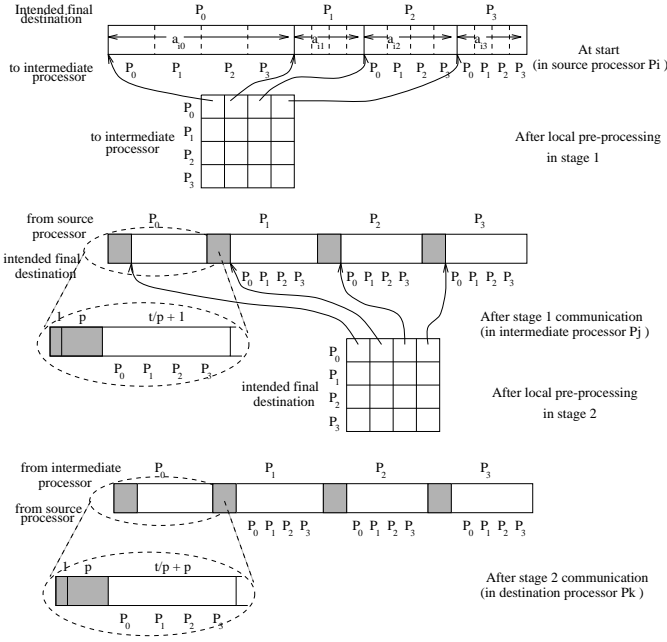


Figure 6: Organization of Data in the Two-stage Algorithm

as an intermediary could receive messages of size up to  $t/p + p$ , since

$$\sum_{i=0}^{p-1} \lceil a_{ik}/p \rceil \leq \sum_{i=0}^{p-1} a_{ik}/p + p \leq c_k/p + p \leq t/p + p$$

The lower bound for the message size is 0, unless we are dealing with the equal traffic case, when the lower bound becomes  $t/p - p$ , since

$$\sum_{i=0}^{p-1} \lfloor a_{ik}/p \rfloor \geq \sum_{i=0}^{p-1} a_{ik}/p - p \geq c_k/p - p \geq t/p - p$$

Our goal is to achieve communication with low variance in message sizes for both stages. A simple change in the algorithm, as described below, can achieve the balance we desire for the first stage.

At any processor  $P_i$ , when dividing the  $a_{ij}$  elements into  $p$  messages for sending, the last  $a_{ij} \bmod p$  elements are assigned to the  $p$  messages in a round robin fashion. This ensures that each intermediate processor  $P_k$  receives messages of size no more than  $\lceil t/p \rceil$ . In the equal traffic case, the message sizes can vary by only one element since the smallest message size is  $\lfloor t/p \rfloor$ .

Figure 7 gives the details of algorithm used for local pre-processing in stage 1. The overall time required is  $O(p^2)$ .

```

procedure Stg1pp(sendl, send_msg_start, send_msg_len)
/* This is code that runs in every processor.
* sendl[0..P-1] is the send vector
* index j gives destination proc #, index k gives intermediate proc #
*
* send_msg_start[0..P-1][0..P-1] gives the index of the element from the
* input array marking the start of each of the P parts of the P
* messages sent out from this source processor;
*
* send_msg_len[0..P-1][0..P] gives the length of those parts; in
* particular, the entry send_msg_len[0..P-1][0] gives the total
* length of messages to each intermediate processor
*/
begin
  for j := 0 to P-1 do
    for k := 0 to P-1 do
      send_msg_len[k][j+1] := sendl[j] div P;
      /* (sendl[j] div P) is the # of elements originally meant for
      processor j now being sent to every intermediate processor */

k := 0;
  for j := 0 to P-1 do
    for x := 1 to (sendl[j] mod P) do
      /* (sendl[j] mod P) is # of elements meant for destination proc
      j that could not be divided equally among the intermediate procs */
      begin
        send_msg_len[k][j+1] := send_msg_len[k][j+1] + 1;
        k := (k+1) mod P;
      end;

data_ptr := 0;
  for j := 0 to P-1 do
    for k := 0 to P-1 do
      begin
        send_msg_start[k][j] := data_ptr;
        data_ptr := data_ptr + send_msg_len[k][j+1];
        send_msg_len[k][0] := send_msg_len[k][0] + send_msg_len[k][j+1];
        /* send_msg_len[k][0] is current message size for intern. proc k */
      end
    end;
  end;
end;

```

Figure 7: Local Pre-Processing in Stage 1

```

procedure Stg2pp (send_msg_start, send_msg_len)
/* send_msg_start[0..P-1][0..P-1] gives the index of the element from the
* input array marking the start of each of the P parts of the P
* messages sent out from this intermediate processor;
*
* send_msg_len[0..P-1][0..P] gives the length of those parts; in
* particular, the entry send_msg_len[0..P-1][0] gives the total
* length of messages to each destination processor
*/
begin
  /* initializing the total length of each message */
  for i := 0 to P-1 do
    send_msg_len[i][0] := 0;

  for i := 0 to P-1 do
    begin
      start := (T/P)*i + 1 + P; /* start position of current message */
      for j := 0 to P-1 do /* T is the traffic bound */
        begin
          total := data_int[start*i+j+1]; /* length of current sub-message */
          send_msg_start[j][i] := start;
          start := start + total;
          send_msg_len[j][0] := send_msg_len[j][0] + total;
          send_msg_len[j][i] := total;
        end
      end
    end
  end;
end;

```

Figure 8: Local Pre-Processing in Stage 2

## Communication

In an initial version of the implementation, local reshuffling was done at the source processors in order to get all the data elements being sent to the same intermediate processor into contiguous memory locations. Such reshuffling gets prohibitively expensive when  $t$  is large. Our current implementation requires that the communication routines take as arguments pointers to  $p$  memory locations in the source processor and  $p$  associated lengths for each message being sent, as shown in figure 6. Note that this does not increase the communication startup latency by a factor of  $p$ .

In the equal traffic case, since the communication is balanced with message lengths differing by just one element, modified linear permutation works best. In the general case, distributed scheduling for the first stage's communication may perform better but linear permutation gives an upper bound on the time taken for communication. A maximum of  $p$  messages of length under  $\lceil t/p \rceil$  may need to be sent. In addition, each of these messages has to be padded with  $p$  lengths (and the sum of these  $p$  lengths, see figure 6) to help the intermediate processor determine the destination processor for each message portion. With linear permutation, the communication time required is  $O(p(\tau + \mu(t/p + p)))$ . If  $t/p < p$ , padding with  $p$  lengths can be replaced by padding with  $t/p$  lengths, making the communication time required  $O(p(\tau + \mu t/p))$ .

## 6.2 The Second Stage

### Local pre-processing

At the intermediate processors, each of which receives  $p$  messages, local pre-processing is done as preparation for the second stage. An initial implementation performed reshuffling. Our current implementation sets up a two arrays of size  $p$  containing pointers and lengths for each message sent out in the second stage. Since a maximum of  $p$  messages could be sent out, this takes  $O(p^2)$  time. Figure 8 gives the steps used for local pre-processing in stage 2.

### Communication

Messages sent out in stage 2 could be of size up to  $t/p + p$ . In the general case, the lower bound on message size is 0, but in the equal traffic case, message size cannot be lower than  $t/p - p$ . Lowering the variance in message size, as was done in stage 1, is not as easy any more. The total size of the messages received at

a destination processor is upper bounded by  $t + p^2$ . In practice, a random reshuffling of messages at the source processor, as explained below, could reduce the mean length of messages reaching a destination processor. The upper bound on the communication time required in stage 2 is  $O(p(\tau + \mu(t/p + p)))$ .

## 6.3 Modifying stage 1 to ease stage 2

The main purpose of the first stage was to spread out data leaving the source processors evenly among the intermediate processors. The intended intermediate processor numbers for the  $p$  messages leaving a source processor can be shuffled randomly within groups of messages of size  $\lceil t/p \rceil$  and  $\lfloor t/p \rfloor$ , without affecting the algorithm. This would still preserve the upper bound derived earlier for total number of data elements sent or received in the first stage. The stage 1 communication now needs to include an extra array of length  $p$  tagged on to each outgoing message. This array gives the permutation that was performed locally before the send and is needed at the destination processors since the  $p$  parts of a message reaching a destination processor must be put back together in order to complete the collective communication. This random reshuffling of messages reduces the expected length of the messages in stage 2 (see [8] for details).

## 6.4 Analysis of time complexity:

The local pre-processing needed for the two-stage algorithm takes  $O(p^2)$  time. The two communication stages take  $O(p(\tau + \mu(t/p + p)))$  time. Thus the two-stage algorithm has a deterministic time complexity of  $O(p\tau + \mu(t + p^2))$ . The algorithm takes time  $O(t)$  and is optimal when  $t \geq O(p^2 + p\tau/\mu)$ . In the case where every  $a_{ij}$  is a multiple of  $p$ , that is, if the message sent by any source processor to any destination processor is a multiple of  $p$ , optimality is achieved when  $t \geq O(p\tau/\mu)$ . The last condition is included to ensure that the startup overhead does not dominate the communication time.

An algorithm for many-to-many communication based on sorting can provide a better time complexity in the general case. Since the destination processors are numbers from a fixed range, local sorting done using a radix-sort takes just  $O(t)$  time. Data movement between processors can be achieved using an adaptation of rotate-sort [6]. Such a combination was used to perform fixed permutations in [7]. The rotate-sort and radix-sort combination performs many-to-many communication in  $O(t)$  time when  $t \geq O(p\tau/\mu)$ . The constants associated with this complexity are, however, much higher than the constants in the two-stage

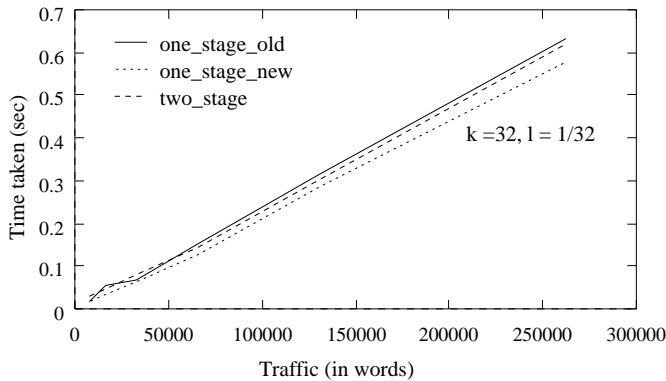


Figure 9: One of the good performances of the two-stage algorithm

algorithm's time complexity.

## 7 Performance Results

The two-stage algorithm and the single-stage algorithm were implemented on the CM-5 using the CMMD message passing library with CMAML active message routines. Communication matrices were generated such that message sizes were non-uniform while the traffic was bounded. Three parameters were used to control the kind of matrix that was generated. The fan-out parameter  $k$  specified the number of processors that each processor communicates with ( $k \leq p$ ). The sum of the messages being sent out and received at each processor was fixed at  $t$ , the traffic parameter. A parameter  $l$  was used to control the non-uniformity of messages sent out by the processors. It was used as follows: Of the  $k$  processors receiving messages from a single processor, the fraction  $lt$  of the traffic reached  $(1-l)k$  processors, while the remaining  $(1-l)t$  traffic reached  $lk$  processors.

Sample values of  $k$  and  $l$  were chosen to highlight a best-case and a worst-case performance of the two-stage algorithm among the trials that were conducted. Figure 9 illustrates the best case in which the two-stage algorithm performed as well as the single-stage algorithms, even out-performing the single-stage algorithm without the improvements. In this trial  $k$  and  $l$  were fixed at 32 and  $1/32$  respectively, which indicates that 1 out of 32 processors received  $31/32$  of the total traffic, while the other 31 processors received in total  $1/32$  of the traffic. It was a trial in which the messages were highly non-uniform in size. Figure 10 illustrates a worst case for the two-stage algorithm. Both the single-stage algorithms out-performed the two-stage one. In this trial  $k$  and  $l$  were fixed at 2 and  $1/2$

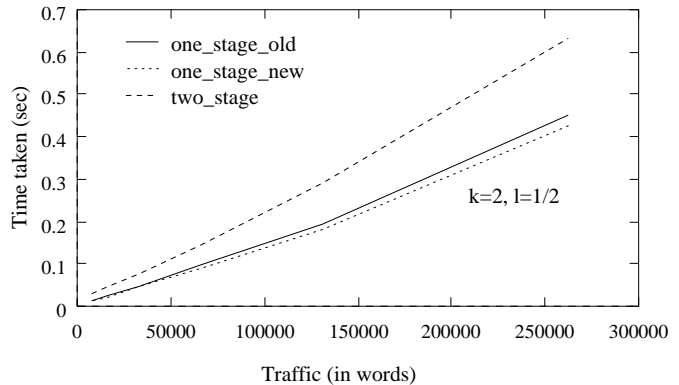


Figure 10: One of the bad performances of the two-stage algorithm

respectively. This indicates that only 2 processors receive data from a single processor, and both of them receive exactly the same amount of traffic. It was a trial in which the messages were uniform in size. The two-stage algorithm's performance remained roughly close to its best-case performance, but the single-stage algorithm's performance improved considerably.

Although the single-stage algorithms were consistent in performing better than the two stage algorithm, they exhibited a much larger variance in the time taken. The two-stage algorithm timings were within a factor of 1.5 times the single-stage readings. It should be noted that the two-stage algorithm is fairly architecture-independent, while the single-stage algorithm (particularly the one with the improvements) is architecture-dependent. The latter is also highly dependent on the availability of low latency communication primitives.

## 8 Conclusions

We have presented a variety of solutions for the problem of many-to-many personalized communication with bounded traffic on a distributed memory parallel machine. A two-stage algorithm that takes time  $O(t)$  when  $t \geq O(p^2 + p\tau/\mu)$  was presented. An algorithm using sorting can improve the result to  $O(t)$  time for  $t \geq O(p\tau/\mu)$ , but the associated constants make this algorithm less desirable for implementation. A single-stage algorithm using distributed random scheduling was implemented and compared with an implementation of the two-stage algorithm. The distributed scheduling algorithm performed better on the CM-5, but this result is not expected to apply to other architectures. Besides, the single-stage algorithm is not deterministic, and that makes it difficult



to ascertain its time complexity.

We have shown that many-to-many personalized communication with non-uniform messages can be performed using two stages of all-to-all personalized communication with uniform messages. Thus, the performance of the two-stage algorithm is roughly half that of an all-to-all personalized communication with the same amount of traffic. The latter problem has been widely investigated in the literature for a variety of interconnection networks (meshes, hypercubes, etc), message passing strategies (wormhole routing, store-and-forward routing, etc), single-port vs. multi-port communication. This makes the two-stage decomposition method useful for a wide variety of architectures. We are currently investigating the performance of these algorithms on other parallel architectures (Intel Paragon, iPSC 860, and the IBM SP1).

## References

- [1] Shahid H. Bokhari. Complete Exchange on the iPSC/860. ICASE Technical Report No. 91-4, NASA Langley Research Center, January 1991.
- [2] Eric A. Brewer and Robert Blumofe, Strata: A Multi-Layer Communications Library, MIT Laboratory of Computer Science Technical Report, February 1994.
- [3] Eric A. Brewer, Bradley C. Kuszmaul, How to Get Good Performance from the CM-5 Data Network, Proceedings of the 8th International Parallel Processing Symposium, April 1994.
- [4] T. von Eicken, D.E. Culler, S.C. Goldstein, K.E.Schauser. Active Messages: a mechanism for integrated communication and computation. Proceedings of the ISCA '92, Gold Coast, Australia, May 1992.
- [5] Vipin Kumar, Ananth Grama, Anshul Gupta, George Karypis. *Introduction to Parallel Computing: Design and Analysis of Algorithms*, Benjamin-Cummings, 1994.
- [6] J. Marberg, E.Gafni. Sorting in Constant Number of Row and Column Phases on a Mesh. *Algorithmica*, Vol.3, pp.561-572, 1988.
- [7] Victor K. Prasanna, Cho-Li Wang, Scalable Data Parallel Object Recognition using Geometric Hashing on the CM-5. Scalable High Performance Computing Conference, SHPCC, 1994.
- [8] Ravi V. Shankar, Khaled A. Alsabti, Sanjay Ranka. The Transportation Primitive, CIS Technical Report, Syracuse University, August 1994.
- [9] Ravi V. Shankar, Sanjay Ranka. Random Data Accesses on a Coarse-Grained Parallel Machine - I. One-to-one Mappings, CIS Technical Report, Syracuse University, October 1994.
- [10] Ravi V. Shankar, Sanjay Ranka. Random Data Accesses on a Coarse-Grained Parallel Machine - II. One-to-many and Many-to-one Mappings, CIS Technical Report, Syracuse University, October 1994.
- [11] Thinking Machines Corporation. *The Connection Machine CM-5 Technical Summary*, October 1991.
- [12] Thinking Machines Corporation. *CMMD Reference Manual Version 3.0*, October 1991.
- [13] Jhy-chun Wang, Tseng-Hui Lin, Sanjay Ranka. Distributed Scheduling of Unstructured Collective Communication on the CM-5. Hawaii International Conference on System Sciences, 1993.