

Content not for sharing. For personal academic use only!

Chapter 4. Basic Communication Operations

In most parallel algorithms, processes need to exchange data with other processes. This exchange of data can significantly impact the efficiency of parallel programs by introducing interaction delays during their execution. For instance, recall from [Section 2.5](#) that it takes roughly $t_s + mt_w$ time for a simple exchange of an m -word message between two processes running on different nodes of an interconnection network with cut-through routing. Here t_s is the latency or the startup time for the data transfer and t_w is the per-word transfer time, which is inversely proportional to the available bandwidth between the nodes. Many interactions in practical parallel programs occur in well-defined patterns involving more than two processes. Often either all processes participate together in a single global interaction operation, or subsets of processes participate in interactions local to each subset. These common basic patterns of interprocess interaction or communication are frequently used as building blocks in a variety of parallel algorithms. Proper implementation of these basic communication operations on various parallel architectures is a key to the efficient execution of the parallel algorithms that use them.

In this chapter, we present algorithms to implement some commonly used communication patterns on simple interconnection networks, such as the linear array, two-dimensional mesh, and the hypercube. The choice of these interconnection networks is motivated primarily by pedagogical reasons. For instance, although it is unlikely that large scale parallel computers will be based on the linear array or ring topology, it is important to understand various communication operations in the context of linear arrays because the rows and columns of meshes are linear arrays. Parallel algorithms that perform rowwise or columnwise communication on meshes use linear array algorithms. The algorithms for a number of

like
blocks
in
a
line

problem
under many assumptions
common case
observation
solution sketch

our focus

why linear array

communication operations on a mesh are simple extensions of the corresponding linear array algorithms to two dimensions. Furthermore, parallel algorithms using regular data structures such as arrays often map naturally onto one- or two-dimensional arrays of processes. This too makes it important to study interprocess interaction on a linear array or mesh interconnection network. The hypercube architecture, on the other hand, is interesting because many algorithms with recursive interaction patterns map naturally onto a hypercube topology. Most of these algorithms may perform equally well on interconnection networks other than the hypercube, but it is simpler to visualize their communication patterns on a hypercube.

The algorithms presented in this chapter in the context of simple network topologies are practical and are highly suitable for modern parallel computers, even though most such computers are unlikely to have an interconnection network that exactly matches one of the networks considered in this chapter. The reason is that on a modern parallel computer, the time to transfer data of a certain size between two nodes is often independent of the relative location of the nodes in the interconnection network. This homogeneity is afforded by a variety of firmware and hardware features such as randomized routing algorithms and cut-through routing, etc. Furthermore, the end user usually does not have explicit control over mapping processes onto physical processors. Therefore, we assume that the transfer of m words of data between *any* pair of nodes in an interconnection network incurs a cost of $t_s + mt_w$. On most architectures, this assumption is reasonably accurate as long as a free link is available between the source and destination nodes for the data to traverse. However, if many pairs of nodes are communicating simultaneously, then the messages may take longer. This can happen if the number of messages passing through a cross-section of the network exceeds the cross-section bandwidth ([Section 2.4.4](#)) of the network. In such situations, we need to adjust the value of t_w to reflect the slowdown due to congestion. As discussed in [Section 2.5.1](#), we refer to the adjusted value of t_w as effective t_w . We will make a note in the text when we come across

When links can be congested

Why
hypercube

May be
when we
don't have
contention

Another
assumption

communication operations that may cause congestion on certain networks.

As discussed in [Section 2.5.2](#), the cost of data-sharing among processors in the shared-address-space paradigm can be modeled using the same expression $t_s + mt_w$, usually with different values of t_s and t_w relative to each other as well as relative to the computation speed of the processors of the parallel computer. Therefore, parallel algorithms requiring one or more of the interaction patterns discussed in this chapter can be assumed to incur costs whose expression is close to one derived in the context of message-passing.

**Mole
assum
tions**

In the following sections we describe various communication operations and derive expressions for their time complexity. We assume that the interconnection network supports [cut-through routing](#) ([Section 2.5.1](#)) and that the communication time between any pair of nodes is practically independent of the number of intermediate nodes along the paths between them. We also assume that the communication links are [bidirectional](#); that is, two directly-connected nodes can send messages of size m to each other simultaneously in time $t_s + t_w m$. We assume a [single-port communication](#) model, in which a node can send a message on only one of its links at a time. Similarly, it can receive a message on only one link at a time. However, a node can receive a message while sending another message at the same time on the same or a different link.

Many of the operations described here have [duals](#) and other related operations that we can perform by using procedures very similar to those for the original operations. The dual of a communication operation is the opposite of the original operation and can be performed by reversing the direction and sequence of messages in the original operation. We will mention such operations wherever applicable.

One-to-All Broadcast and All-to-One Reduction

- we ignored propagation delays
- cut-through routing with small packets
⇒ as if we are using circuit switching

D & of dual

Parallel algorithms often require a single process to send identical data to all other processes or to a subset of them. This operation is known as one-to-all broadcast. Initially, only the source process has the data of size m that needs to be broadcast. At the termination of the procedure, there are p copies of the initial data – one belonging to each process. The dual of one-to-all broadcast is all-to-one reduction. In an all-to-one reduction operation, each of the p participating processes starts with a buffer M containing m words. The data from all processes are combined through an associative operator and accumulated at a single destination process into one buffer of size m . Reduction can be used to find the sum, product, maximum, or minimum of sets of numbers – the i th word of the accumulated M is the sum, product, maximum, or minimum of the i th words of each of the original buffers. [Figure 4.1](#) shows one-to-all broadcast and all-to-one reduction among p processes.

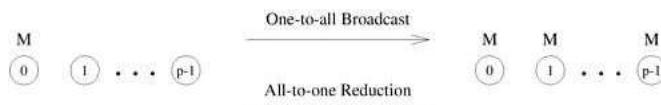


Figure 4.1. One-to-all broadcast and all-to-one reduction.

$(12A) \equiv \text{one-to-all}$

One-to-all broadcast and all-to-one reduction are used in several important parallel algorithms including matrix-vector multiplication, Gaussian elimination, shortest paths, and vector inner product. In the following subsections, we consider the implementation of one-to-all broadcast in detail on a variety of interconnection topologies.

Implementation of $(12A)$ operator:

Ring or Linear Array

① slow →

A naive way to perform one-to-all broadcast is to sequentially send $p - 1$ messages from the source to the other $p - 1$ processes. However, this is inefficient because the source process becomes a bottleneck. Moreover, the communication network is underutilized because only the connection between a single pair of nodes is used at a time. A better broadcast algorithm can be devised using a technique commonly known as recursive doubling. The source process first sends the message to another process. Now both these processes can simultaneously send the message to two other processes that are still waiting for the message. By continuing this

The beautiful thing in this algo is how we select nodes

one to all
initial state
source has m
final state
Each of P
processes has
 m

Real apps using $(12A)$ operat

Cost:

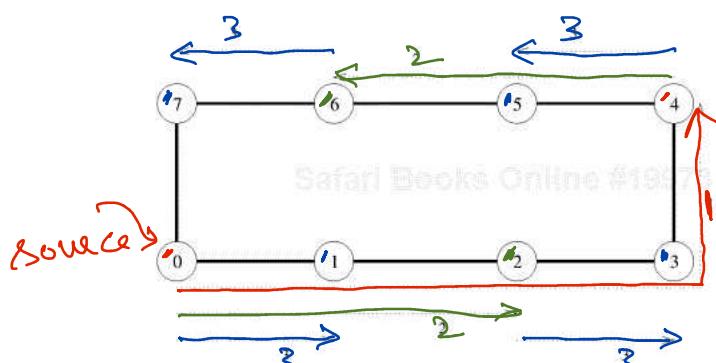
$(t_s + mt_w)(p-1)$

So we have an upper bound with naive algo

select

procedure until all the processes have received the data, the message can be broadcast in $\log p$ steps. *That is a substantial improvement from $(P-1)$ to $\log P$*

The steps in a one-to-all broadcast on an eight-node linear array or ring are shown in [Figure 4.2](#). The nodes are labeled from 0 to 7. Each message transmission step is shown by a numbered, dotted arrow from the source of the message to its destination. Arrows indicating messages sent during the same time step have the same number.



[Figure 4.2](#). One-to-all broadcast on an eight-node ring. Node 0 is the source of the broadcast. Each message transfer step is shown by a numbered, dotted arrow from the source of the message to its destination. The number on an arrow indicates the time step during which the message is transferred.

Note that on a linear array, the destination node to which the message is sent in each step must be carefully chosen. In [Figure 4.2](#), the message is first sent to the farthest node (4) from the source (0). In the second step, the distance between the sending and receiving nodes is halved, and so on. The message recipients are selected in this manner at each step to avoid congestion on the network. For example, if node 0 sent the message to node 1 in the first step and then nodes 0 and 1 attempted to send messages to nodes 2 and 3, respectively, in the second step, the link between nodes 1 and 2 would be congested as it would be a part of the shortest route for both the messages in the second step.

Reduction on a linear array can be performed by simply reversing the direction and the sequence of communication, as shown in [Figure 4.3](#). In the first step, each odd numbered node sends its buffer to the even numbered node just before itself, where the contents of the two buffers are combined into one. After the first step, there are four buffers left to be reduced on nodes 0, 2, 4, and 6, respectively. In the second step, the contents of the buffers on nodes 0 and 2 are accumulated on node 0 and those on

Node labeling will be crucial.

$$\log_2 8 = 3 \text{ steps}$$

for 8 node interconnect

Observations

- ① First send half-way through
- ② All do anti-clock wise for best link utilization
- ③ Network symmetry helped us (*Internet does not have such symmetry!*)

*Example:
Why we chose the nodes the way we did.*

nodes 6 and 4 are accumulated on node 4. Finally, node 4 sends its buffer to node 0, which computes the final result of the reduction.

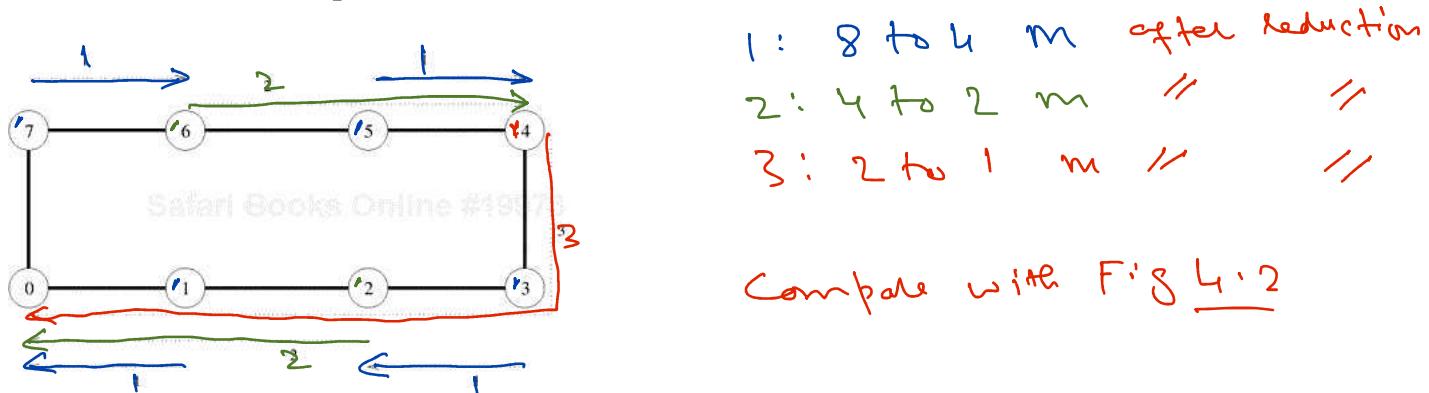


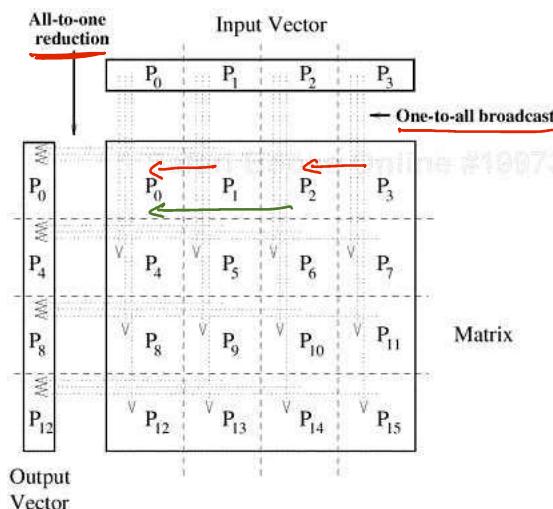
Figure 4.3. Reduction on an eight-node ring with node 0 as the destination of the reduction.

EXAMPLE 4.1 MATRIX-VECTOR MULTIPLICATION

Consider the problem of multiplying an $n \times n$ matrix A with an $n \times 1$ vector x on an $n \times n$ mesh of nodes to yield an $n \times 1$ result vector y . Algorithm 8.1 shows a serial algorithm for this problem.

Figure 4.4 shows one possible mapping of the matrix and the vectors in which each element of the matrix belongs to a different process, and the vector is distributed among the processes in the topmost row of the mesh and the result vector is generated on the leftmost column of processes.

Re call
owner's
Compute
rule!



- ① In parallel, after $\log_2 4 = 2$ steps, all processes have required data
- ② In 2 steps reduction happens

FIGURE 4.4. ONE-TO-ALL BROADCAST AND ALL-TO-ONE REDUCTION IN THE MULTIPLICATION OF A 4×4 MATRIX WITH A 4×1 VECTOR.

Since all the rows of the matrix must be multiplied with the vector, each process needs the element of the vector residing in the topmost process of its column. Hence, before computing the ma-

matrix-vector product, each column of nodes performs a one-to-all broadcast of the vector elements with the topmost process of the column as the source. This is done by treating each column of the $n \times n$ mesh as an n -node linear array, and simultaneously applying the linear array broadcast procedure described previously to all columns.

After the broadcast, each process multiplies its matrix element with the result of the broadcast. Now, each row of processes needs to add its result to generate the corresponding element of the product vector. This is accomplished by performing all-to-one reduction on each row of the process mesh with the first process of each row as the destination of the reduction operation.

For example, P_9 will receive $x[1]$ from P_1 as a result of the broadcast, will multiply it with $A[2, 1]$ and will participate in an all-to-one reduction with P_8 , P_{10} , and P_{11} to accumulate $y[2]$ on P_8 . ▪

Mesh

We can regard each row and column of a square mesh of p nodes as a linear array of \sqrt{p} nodes. So a number of communication algorithms on the mesh are simple extensions of their linear array counterparts. A linear array communication operation can be performed in two phases on a mesh. In the first phase, the operation is performed along one or all rows by treating the rows as linear arrays. In the second phase, the columns are treated similarly.

$$\begin{aligned} & \text{X} \\ & \text{---} \\ & \text{X} = P \\ & \Rightarrow X = \sqrt{P} \end{aligned}$$

Consider the problem of one-to-all broadcast on a two-dimensional square mesh with \sqrt{p} rows and \sqrt{p} columns. First, a one-to-all broadcast is performed from the source to the remaining $(\sqrt{p}-1)$ nodes of the same row. Once all the nodes in a row of the mesh have acquired the data, they initiate a one-to-all broadcast in their respective columns. At the end of the second phase, every node in the mesh has a copy of the initial message. The communication steps for one-to-all broadcast on a mesh are illustrated in Figure 4.5 for $p = 16$, with node 0 at the bottom-left corner as

the source. Steps 1 and 2 correspond to the first phase, and steps 3 and 4 correspond to the second phase.

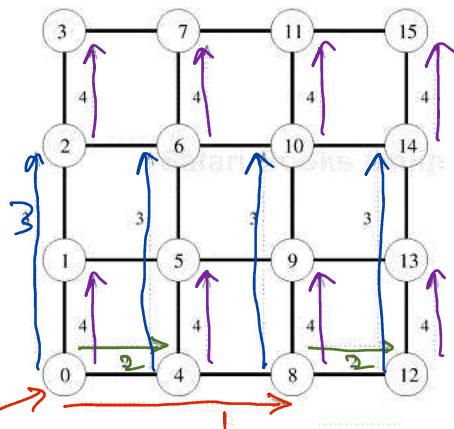


Figure 4.5. One-to-all broadcast on a 16-node mesh.

- ① After \sqrt{P} (2) steps bottom row has data
- ② 2 more steps for each column.

We can use a similar procedure for one-to-all broadcast on a three-dimensional mesh as well. In this case, rows of $p^{1/3}$ nodes in each of the three dimensions of the mesh would be treated as linear arrays. As in the case of a linear array, reduction can be performed on two- and three-dimensional meshes by simply reversing the direction and the order of messages.

Hypercube

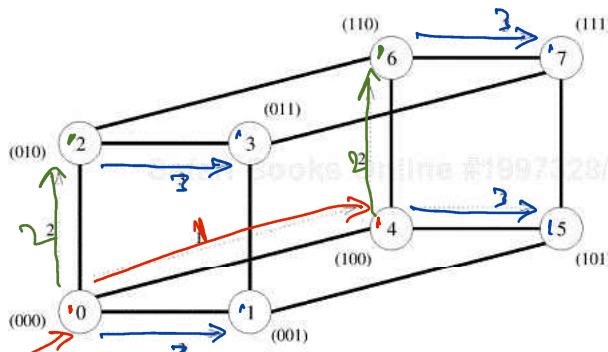
The previous subsection showed that one-to-all broadcast is performed in two phases on a two-dimensional mesh, with the communication taking place along a different dimension in each phase. Similarly, the process is carried out in three phases on a three-dimensional mesh. A hypercube with 2^d nodes can be regarded as a d -dimensional mesh with two nodes in each dimension. Hence, the mesh algorithm can be extended to the hypercube, except that the process is now carried out in d steps – one in each dimension.

Figure 4.6 shows a one-to-all broadcast on an eight-node (three-dimensional) hypercube with node 0 as the source. In this figure, communication starts along the highest dimension (that is, the dimension specified by the most significant bit of the binary representation of a node label) and proceeds along successively lower dimensions in subsequent steps. Note that the source and the destination nodes in three communication

Moving
like linear
array

steps of the algorithm shown in [Figure 4.6](#) are identical to the ones in the broadcast algorithm on a linear array shown in [Figure 4.2](#). However, on a hypercube, the order in which the dimensions are chosen for communication does not affect the outcome of the procedure. [Figure 4.6](#) shows only one such order. Unlike a linear array, the hypercube broadcast would not suffer from congestion if node 0 started out by sending the message to node 1 in the first step, followed by nodes 0 and 1 sending messages to nodes 2 and 3, respectively, and finally nodes 0, 1, 2, and 3 sending messages to nodes 4, 5, 6, and 7, respectively.

That is
be cause
a lot
of links
are available



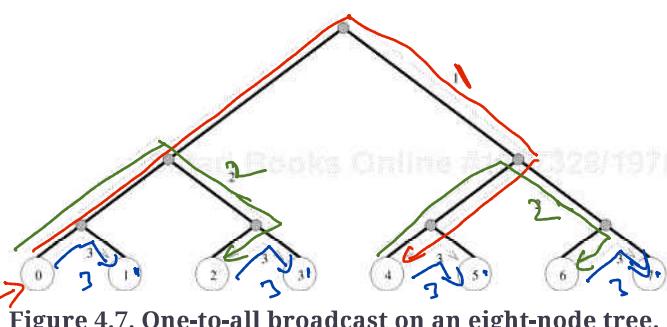
[Figure 4.6](#). One-to-all broadcast on a three-dimensional hypercube. The binary representations of node labels are shown in parentheses.

Balanced Binary Tree

The hypercube algorithm for one-to-all broadcast maps naturally onto a balanced binary tree in which each leaf is a processing node and intermediate nodes serve only as switching units. This is illustrated in [Figure 4.7](#) for eight nodes. In this figure, the communicating nodes have the same labels as in the hypercube algorithm illustrated in [Figure 4.6](#). [Figure 4.7](#) shows that there is no congestion on any of the communication links at any time. The difference between the communication on a hypercube and the tree shown in [Figure 4.7](#) is that there is a different number of switching nodes along different paths on the tree.

Observations

- ① No link Congestion in steps



[Figure 4.7](#). One-to-all broadcast on an eight-node tree.

Detailed Algorithms

This is a major observation

A careful look at Figures 4.2, 4.5, 4.6, and 4.7 would reveal that the basic communication pattern for one-to-all broadcast is identical on all the four interconnection networks considered in this section. We now describe procedures to implement the broadcast and reduction operations. For the sake of simplicity, the algorithms are described here in the context of a hypercube and assume that the number of communicating processes is a power of 2. However, they apply to any network topology, and can be easily extended to work for any number of processes (Problem 4.1).

Algorithm 4.1 shows a one-to-all broadcast procedure on a 2^d -node network when node 0 is the source of the broadcast. The procedure is executed at all the nodes. At any node, the value of $\underline{my_id}$ is the label of that node. Let X be the message to be broadcast, which initially resides at the source node 0. The procedure performs d communication steps, one along each dimension of a hypothetical hypercube. In **Algorithm 4.1**, communication proceeds from the highest to the lowest dimension (although the order in which dimensions are chosen does not matter). The loop counter i indicates the current dimension of the hypercube in which communication is taking place. Only the nodes with zero in the i least significant bits of their labels participate in communication along dimension

$i = 1 = 010$

$i = 2 = 100$

$i = 3 = 110$

$i = 4 = 111$

$i = 5 = 101$

$i = 6 = 011$

$i = 7 = 001$

$i = 8 = 000$

$i = 9 = 001$

$i = 10 = 011$

$i = 11 = 111$

$i = 12 = 110$

$i = 13 = 101$

$i = 14 = 010$

$i = 15 = 001$

$i = 16 = 000$

$i = 17 = 001$

$i = 18 = 011$

$i = 19 = 111$

$i = 20 = 110$

$i = 21 = 101$

$i = 22 = 010$

$i = 23 = 001$

$i = 24 = 000$

$i = 25 = 001$

$i = 26 = 011$

$i = 27 = 111$

$i = 28 = 110$

$i = 29 = 101$

$i = 30 = 010$

$i = 31 = 001$

$i = 32 = 000$

$i = 33 = 001$

$i = 34 = 011$

$i = 35 = 111$

$i = 36 = 110$

$i = 37 = 101$

$i = 38 = 010$

$i = 39 = 001$

$i = 40 = 000$

$i = 41 = 001$

$i = 42 = 011$

$i = 43 = 111$

$i = 44 = 110$

$i = 45 = 101$

$i = 46 = 010$

$i = 47 = 001$

$i = 48 = 000$

$i = 49 = 001$

$i = 50 = 011$

$i = 51 = 111$

$i = 52 = 110$

$i = 53 = 101$

$i = 54 = 010$

$i = 55 = 001$

$i = 56 = 000$

$i = 57 = 001$

$i = 58 = 011$

$i = 59 = 111$

$i = 60 = 110$

$i = 61 = 101$

$i = 62 = 010$

$i = 63 = 001$

$i = 64 = 000$

$i = 65 = 001$

$i = 66 = 011$

$i = 67 = 111$

$i = 68 = 110$

$i = 69 = 101$

$i = 70 = 010$

$i = 71 = 001$

$i = 72 = 000$

$i = 73 = 001$

$i = 74 = 011$

$i = 75 = 111$

$i = 76 = 110$

$i = 77 = 101$

$i = 78 = 010$

$i = 79 = 001$

$i = 80 = 000$

$i = 81 = 001$

$i = 82 = 011$

$i = 83 = 111$

$i = 84 = 110$

$i = 85 = 101$

$i = 86 = 010$

$i = 87 = 001$

$i = 88 = 000$

$i = 89 = 001$

$i = 90 = 011$

$i = 91 = 111$

$i = 92 = 110$

$i = 93 = 101$

$i = 94 = 010$

$i = 95 = 001$

$i = 96 = 000$

$i = 97 = 001$

$i = 98 = 011$

$i = 99 = 111$

$i = 100 = 110$

$i = 101 = 101$

$i = 102 = 010$

$i = 103 = 001$

$i = 104 = 000$

$i = 105 = 001$

$i = 106 = 011$

$i = 107 = 111$

$i = 108 = 110$

$i = 109 = 101$

$i = 110 = 010$

$i = 111 = 001$

$i = 112 = 000$

$i = 113 = 001$

$i = 114 = 011$

$i = 115 = 111$

$i = 116 = 110$

$i = 117 = 101$

$i = 118 = 010$

$i = 119 = 001$

$i = 120 = 000$

$i = 121 = 001$

$i = 122 = 011$

$i = 123 = 111$

$i = 124 = 110$

$i = 125 = 101$

$i = 126 = 010$

$i = 127 = 001$

$i = 128 = 000$

$i = 129 = 001$

$i = 130 = 011$

$i = 131 = 111$

$i = 132 = 110$

$i = 133 = 101$

$i = 134 = 010$

$i = 135 = 001$

$i = 136 = 000$

$i = 137 = 001$

$i = 138 = 011$

$i = 139 = 111$

$i = 140 = 110$

$i = 141 = 101$

$i = 142 = 010$

$i = 143 = 001$

$i = 144 = 000$

$i = 145 = 001$

$i = 146 = 011$

$i = 147 = 111$

$i = 148 = 110$

$i = 149 = 101$

$i = 150 = 010$

$i = 151 = 001$

$i = 152 = 000$

$i = 153 = 001$

$i = 154 = 011$

$i = 155 = 111$

$i = 156 = 110$

$i = 157 = 101$

$i = 158 = 010$

$i = 159 = 001$

$i = 160 = 000$

$i = 161 = 001$

$i = 162 = 011$

$i = 163 = 111$

$i = 164 = 110$

$i = 165 = 101$

$i = 166 = 010$

$i = 167 = 001$

$i = 168 = 000$

$i = 169 = 001$

$i = 170 = 011$

$i = 171 = 111$

$i = 172 = 110$

$i = 173 = 101$

$i = 174 = 010$

$i = 175 = 001$

$i = 176 = 000$

$i = 177 = 001$

$i = 178 = 011$

$i = 179 = 111$

$i = 180 = 110$

$i = 181 = 101$

$i = 182 = 010$

$i = 183 = 001$

$i = 184 = 000$

$i = 185 = 001$

$i = 186 = 011$

$i = 187 = 111$

$i = 188 = 110$

$i = 189 = 101$

$i = 190 = 010$

$i = 191 = 001$

$i = 192 = 000$

$i = 193 = 001$

$i = 194 = 011$

$i = 195 = 111$

$i = 196 = 110$

$i = 197 = 101$

$i = 198 = 010$

$i = 199 = 001$

$i = 200 = 000$

$i = 201 = 001$

$i = 202 = 011$

$i = 203 = 111$

$i = 204 = 110$

$i = 205 = 101$

$i = 206 = 010$

$i = 207 = 001$

$i = 208 = 000$

$i = 209 = 001$

$i = 210 = 011$

$i = 211 = 111$

$i = 212 = 110$

$i = 213 = 101$

$i = 214 = 010$

$i = 215 = 001$

$i = 216 = 000$

$i = 217 = 001$

$i = 218 = 011$

$i = 219 = 111$

$i = 220 = 110$

$i = 221 = 101$

$i = 222 = 010$

$i = 223 = 001$

$i = 224 = 000$

$i = 225 = 001$

$i = 226 = 011$

$i = 227 = 111$

$i = 228 = 110$

$i = 229 = 101$

$i = 230 = 010$

$i = 231 = 001$

$i = 232 = 000$

$i = 233 = 001$

$i = 234 = 011$

$i = 235 = 111$

$i = 236 = 110$

$i = 237 = 101$

$i = 238 = 010$

$i = 239 = 001$

$i = 240 = 000$

$i = 241 = 001$

$i = 242 = 011$

$i = 243 = 111$

$i = 244 = 110$

$i = 245 = 101$

$i = 246 = 010$

$i = 247 = 001$

$i = 248 = 000$

$i = 249 = 001$

$i = 250 = 011$

$i = 251 = 111$

$i = 252 = 110$

$i = 253 = 101$

$i = 254 = 010$

$i = 255 = 001$

$i = 256 = 000$

$i = 257 = 001$

$i = 258 = 011$

$i = 259 = 111$

$i = 260 = 110$

$i = 261 = 101$

$i = 262 = 010$

$i = 263 = 001$

$i = 264 = 000$

$i = 265 = 001$

$i = 266 = 011$

$i = 267 = 111$

$i = 268 = 110$

$i = 269 = 101$

$i = 270 = 010$

$i = 271 = 001$

$i = 272 = 000$

$i = 273 = 001$

$i = 274 = 011$

$i = 275 = 111$

$i = 276 = 110$

$i = 277 = 101$

$i = 278 = 010$

$i = 279 = 001$

$i = 280 = 000$

$i = 281 = 001$

$i = 282 = 011$

$i = 283 = 111$

$i = 284 = 110$

$i = 285 = 101$

$i = 286 = 010$

$i = 287 = 001$

$i = 288 = 000$

$i = 289 = 001$

$i = 290 = 011$

$i = 291 = 111$

$i = 292 = 110$

$i = 293 = 101$

$i = 294 = 010$

$i = 295 = 001$

$i = 296 = 000$

$i = 297 = 001$

$i = 298 = 011$

$i = 299 = 111$

$i = 300 = 110$

$i = 301 = 101$

$i = 302 = 010$

$i = 303 = 001$

$i = 304 = 000$

$i = 305 = 001$

$i = 306 = 011$

$i = 307 = 111$

$i = 308 = 110$

$i = 309 = 101$

$i = 310 = 010$

$i = 311 = 001$

$i = 312 = 000$

$i = 313 = 001$

$i = 314 = 011$

$i = 315 = 111$

$i = 316 = 110$

$i = 317 = 101$

$i = 318 = 010$

$i = 319 = 001$

$i = 320 = 000$

$i = 321 = 001$

$i = 322 = 011$

$i = 323 = 111$

$i = 324 = 110$

$i = 325 = 101$

$i = 326 = 010$

$i = 327 = 001$

$i = 328 = 000$

$i = 329 = 001$

$i = 330 = 011$

$i = 331 = 111$

$i = 332 = 110$

$i = 333 = 101$

$i = 334 = 010$

$i = 335 = 001$

$i = 336 = 000$

$i = 337 = 001$

$i = 338 = 011$

$i = 339 = 111$

$i = 340 = 110$

$i = 341 = 101$

$i = 342 = 010$

$i = 343 = 001$

$i = 344 = 000$

$i = 345 = 001$

$i = 346 = 011$

$i = 347 = 111$

$i = 348 = 110$

$i = 349 = 101$

$i = 350 = 010$

$i = 351 = 001$

$i = 352 = 000$

$i = 353 = 001$

$i = 354 = 011$

$i = 355 = 111$

$i = 356 = 110$

$i = 357 = 101$

$i = 358 = 010$

$i = 359 = 001$

$i = 360 = 000$

$i = 361 = 001$

$i = 362 = 011$

$i = 363 = 111$

$i = 364 = 110$

$i = 365 = 101$

$i = 366 = 010$

$i = 367 = 001$

$i = 368 = 000$

$i = 369 = 001$

$i = 370 = 011$

$i = 371 = 111$

$i = 372 = 110$

$i = 373 = 101$

$i = 374 = 010$

$i = 375 = 001$

$i = 376 = 000$

$i = 377 = 001$

$i = 378 = 011$

$i = 379 = 111$

$i = 380 = 110$

$i = 381 = 101$

$i = 382 = 010$

$i = 383 = 001$

$i = 384 = 000$

$i = 385 = 001$

$i = 386 = 011$

$i = 387 = 111$

$i = 388 = 110$

$i = 389 = 101$

$i = 390 = 010$

$i = 391 = 001$

$i = 392 = 000$

$i = 393 = 001$

$i = 394 = 011$

$i = 395 = 111$

$i = 396 = 110$

$i = 397 = 101$

$i = 398 = 010$

$i = 399 = 001$

$i = 400 = 000$

$i = 401 = 001$

$i = 402 = 011$

$i = 403 = 111$

$i = 404 = 110$

$i = 405 = 101$

$i = 406 = 010$

$i = 407 = 001$

$i = 408 = 000$

$i = 409 = 001$

$i = 410 = 011$

$i = 411 = 111$

$i = 412 = 110$

$i = 413 = 101$

$i = 414 = 010$

$i = 415 = 001$

$i = 416 = 000$

$i = 417 = 001$

$i = 418 = 011$

$i = 419 = 111$

$i = 420 = 110$

$i = 421 = 101$

$i = 422 = 010$

$i = 423 = 001$

$i = 424 = 000$

$i = 425 = 001$

$i = 426 = 011$

$i = 427 = 111$

$i = 428 = 110$

$i = 429 = 101$

$i = 430 = 010$

$i = 431 = 001$

$i = 432 = 000$

$i = 433 = 001$

$i = 434 = 011$

$i = 435 = 111$

$i = 436 = 110$

$i = 437 = 101$

$i = 438 = 010$

$i = 439 = 001$

$i = 440 = 000$

$i = 441 = 001$

$i = 442 = 011$

$i = 443 = 111$

$i = 444 = 110$

$i = 445 = 101$

$i = 446 = 010$

$i = 447 = 001$

$i = 448 = 000$

$i = 449 = 001$

$i = 450 = 011$

$i = 451 = 111$

$i = 452 = 110$

$i = 453 = 101$

$i = 454 = 010$

$i = 455 = 001$

$i = 456 = 000$

$i = 457 = 001$

$i = 458 = 011$

$i = 459 = 111$

$i = 460 = 110$

$i = 461 = 101$

$i = 462 = 010$

$i = 463 = 001$

$i = 464 = 000$

$i = 465 = 001$

$i = 466 = 011$

$i = 467 = 111$

$i = 468 = 110$

$i = 469 = 101$

$i = 470 = 010$

$i = 471 = 001$

$i = 472 = 000$

$i = 473 = 001$

$i = 474 = 011$

$i = 475 = 111$

$i = 476 = 110$

$i = 477 = 101$

$i = 478 = 010$

$i = 479 = 001$

$i = 480 = 000$

$i = 481 = 001$

$i = 482 = 011$

$i = 483 = 111$

$i = 484 = 110$

$i = 485 = 101$

$i = 486 = 010$

$i = 487 = 001$

$i = 488 = 000$

$i = 489 = 001$

$i = 490 = 011$

$i = 491 = 111$

$i = 492 = 110$

$i = 493 = 101$

$i = 494 = 010$

$i = 495 = 001$

$i = 496 = 000$

$i = 497 = 001$

$i = 498 = 011$

$i = 499 = 111$

$i = 500 = 110$

$i = 501 = 101$

$i = 502 = 010$

$i = 503 = 001$

$i = 504 = 000$

$i = 505 = 001$

$i = 506 = 011$

$i = 507 = 111$

$i = 508 = 110$

$i = 509 = 101$

$i = 510 = 010$

$i = 511 = 001$

$i = 512 = 000$

$i = 513 = 001$

$i = 514 = 011$

$i = 515 = 111$

$i = 516 = 110$

$i = 517 = 101$

$i = 518 = 010$

$i = 519 = 001$

$i = 520 = 000$

$i = 521 = 001$

$i = 522 = 011$

$i = 523 = 111$

$i = 524 = 110$

$i = 525 = 101$

$i = 526 = 010$

$i = 527 = 001$

$i = 528 = 000$

$i = 529 = 001$

$i = 530 = 011$

$i = 531 = 111$

$i = 532 = 110$

$i = 533 = 101$

$i = 534 = 010$

$i = 535 = 001$

$i = 536 = 000$

$i = 537 = 001$

$i = 538 = 011$

$i = 539 = 111$

$i = 540 = 110$

$i = 541 = 101$

$i = 542 = 010$

$i = 543 = 001$

$i = 544 = 000$

$i = 545 = 001$

$i = 546 = 011$

$i = 547 = 111$

$i = 548 = 110$

$i = 549 = 101$

$i = 550 = 010$

$i = 551 = 001$

$i = 552 = 000$

$i = 553 = 001$

$i = 554 = 011$

$i = 555 = 111$

$i = 556 = 110$

$i = 557 = 101$

$i = 558 = 010$

$i = 559 = 001$

$i = 560 = 000$

$i = 561 = 001$

$i = 562 = 011$

$i = 563 = 111$

$i = 564 = 110$

$i = 565 = 101$

$i = 566 = 010$

$i = 567 = 001$

$i = 568 = 000$

$i = 569 = 001$

$i = 570 = 011$

$i = 571 = 111$

$i = 572 = 110$

$i = 573 = 101$

$i = 574 = 010$

$i = 575 = 001$ </

(the i least significant bits of *mask* are ones). The AND operation on Line 6 selects only those nodes that have zeros in their i least significant bits.

Participants: senders + receivers

Among the nodes selected for communication along dimension i , the nodes with a zero at bit position i send the data, and the nodes with a one at bit position i receive it. The test to determine the sending and receiving nodes is performed on Line 7. For example, in [Figure 4.6](#), node 0 (000) is the sender and node 4 (100) is the receiver in the iteration corresponding to $i = 2$. Similarly, for $i = 1$, nodes 0 (000) and 4 (100) are senders while nodes 2 (010) and 6 (110) are receivers.

Bit wise AND

$$i=2, i=4$$

000 AND

$$\begin{array}{r} 10 \\ \times 0 \\ \hline 00 \end{array}$$

sender

$$\begin{array}{r} 100 \\ \times 100 \\ \hline 110 \end{array}$$

receiver

\leftarrow) A
drawback

And
solution

001 =
1st node
001
XOR
011
= 0
New
source

000
001
011
011
100
110
111
111
attacker
1 001
0 000
3 011
2 010
5 101
4 100
7 111
5 110

Algorithm 4.1 works only if node 0 is the source of the broadcast. For an arbitrary source, we must relabel the nodes of the hypothetical hypercube by XORing the label of each node with the label of the source node before we apply this procedure. A modified one-to-all broadcast procedure that works for any value of *source* between 0 and $p - 1$ is shown in **Algorithm 4.2**. By performing the XOR operation at Line 3, **Algorithm 4.2** relabels the source node to 0, and relabels the other nodes relative to the source. After this relabeling, the algorithm of **Algorithm 4.1** can be applied to perform the broadcast.

Algorithm 4.3 gives a procedure to perform an all-to-one reduction on a hypothetical d -dimensional hypercube such that the final result is accumulated on node 0. Single node-accumulation is the dual of one-to-all broadcast. Therefore, we obtain the communication pattern required to implement reduction by reversing the order and the direction of messages in one-to-all broadcast. Procedure ALL_TO_ONE_REDUCE(d, my_id, m, X, sum) shown in **Algorithm 4.3** is very similar to procedure ONE_TO_ALL_BC(d, my_id, X) shown in **Algorithm 4.1**. One difference is that the communication in all-to-one reduction proceeds from the lowest to the highest dimension. This change is reflected in the way that variables *mask* and i are manipulated in **Algorithm 4.3**. The criterion for determining the source and the destination among a pair of communicating nodes is also reversed (Line 7). Apart from these differences, procedure ALL_TO_ONE_REDUCE has extra instructions (Lines 13 and 14) to add the

contents of the messages received by a node in each iteration (any associative operation can be used in place of addition).

Example 4.1. One-to-all broadcast of a message X from node 0 of a d -dimensional p -node hypercube ($d = \log p$). AND and XOR are bitwise logical-and and exclusive-or operations, respectively.

```

1.  procedure ONE_TO_ALL_BC( $d$ ,  $my\_id$ ,  $X$ )
2.  begin
3.       $mask := 2^d - 1$ ; /* Set all  $d$  bits of  $mask$  to 1 */
4.      for  $i := d - 1$  downto 0 do /* Outer loop */ i = current dimension
5.           $mask := mask \text{ XOR } 2^i$ ; /* Set bit  $i$  of  $mask$  to 0 */
6.          if ( $my\_id$  AND  $mask$ ) = 0 then /* If lower  $i$  bits of  $my\_id$  are 0
7.              if ( $my\_id$  AND  $2^i$ ) = 0 then
8.                   $msg\_destination := my\_id \text{ XOR } 2^i$ ;
9.                  send  $X$  to  $msg\_destination$ ;
10.             else
11.                  $msg\_source := my\_id \text{ XOR } 2^i$ ;
12.                 receive  $X$  from  $msg\_source$ ;
13.             endelse;
14.         endif;
15.     endfor;
16. end ONE_TO_ALL_BC

```

mask helps identify nodes that will participate in this round.

Example 4.2. One-to-all broadcast of a message X initiated by $source$ on a d -dimensional hypothetical hypercube. The AND and XOR operations are bitwise logical operations.

```

1.  procedure GENERAL_ONE_TO_ALL_BC( $d$ ,  $my\_id$ ,  $source$ ,  $X$ )
2.  begin
3.       $my\_virtual\_id := my\_id \text{ XOR } source$ ;
4.       $mask := 2^d - 1$ ;
5.      for  $i := d - 1$  downto 0 do /* Outer loop */
6.           $mask := mask \text{ XOR } 2^i$ ; /* Set bit  $i$  of  $mask$  to 0 */
7.          if ( $my\_virtual\_id$  AND  $mask$ ) = 0 then
8.              if ( $my\_virtual\_id$  AND  $2^i$ ) = 0 then
9.                   $virtual\_dest := my\_virtual\_id \text{ XOR } 2^i$ ;

```

```

10.         send X to (virtual_dest XOR source);
           /* Convert virtual_dest to the label of the physical destination */
11.         else
12.             virtual_source := my_virtual_id XOR 2i;
13.             receive X from (virtual_source XOR source);
           /* Convert virtual_source to the label of the physical source */
14.         endelse;
15.     endfor;
16. end GENERAL_ONE_TO_ALL_BC

```

Example 4.3. Single-node accumulation on a d -dimensional hypercube. Each node contributes a message X containing m words, and node 0 is the destination of the sum. The AND and XOR operations are bitwise logical operations.

```

1. procedure ALL_TO_ONE_REDUCE(d, my_id, m, X, sum)
2. begin
3.     for j := 0 to m - 1 do sum[j] := X[j];
4.     mask := 0;
5.     for i := 0 to d - 1 do
        /* Select nodes whose lower i bits are 0 */
6.         if (my_id AND mask) = 0 then
7.             if (my_id AND 2i) ≠ 0 then
8.                 msg_destination := my_id XOR 2i;
9.                 send sum to msg_destination;
10.            else
11.                msg_source := my_id XOR 2i;
12.                receive X from msg_source;
13.                for j := 0 to m - 1 do
14.                    sum[j] := sum[j] + X[j];
15.                endelse;
16.            mask := mask XOR 2i; /* Set bit i of mask to 1 */
17.        endfor;
18.    end ALL_TO_ONE_REDUCE

```

Cost Analysis

Analyzing the cost of one-to-all broadcast and all-to-one reduction is fairly straightforward. Assume that p processes participate in the opera-

Caution :-
Be careful
about this
cont. Don't
blindly apply!

tion and the data to be broadcast or reduced contains m words. The broadcast or reduction procedure involves $\log p$ point-to-point simple message transfers, each at a time cost of $t_s + t_w m$. Therefore, the total time taken by the procedure is

Equation 4.1.

$$\underline{T = (t_s + t_w m) \log p}.$$

All-to-All Broadcast and Reduction

All-to-all broadcast is a generalization of one-to-all broadcast in which all p nodes simultaneously initiate a broadcast. A process sends the same m -word message to every other process, but different processes may broadcast different messages. All-to-all broadcast is used in matrix operations, including matrix multiplication and matrix-vector multiplication. The dual of all-to-all broadcast is all-to-all reduction, in which every node is the destination of an all-to-one reduction (Problem 4.8). [Figure 4.8](#) illustrates all-to-all broadcast and all-to-all reduction.

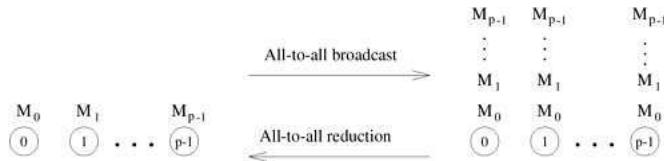


Figure 4.8. All-to-all broadcast and all-to-all reduction.

One way to perform an all-to-all broadcast is to perform p one-to-all broadcasts, one starting at each node. If performed naively, on some architectures this approach may take up to p times as long as a one-to-all broadcast. It is possible to use the communication links in the interconnection network more efficiently by performing all p one-to-all broadcasts simultaneously so that all messages traversing the same path at the same time are concatenated into a single message whose size is the sum of the sizes of individual messages.

The following sections describe all-to-all broadcast on linear array, mesh, and hypercube topologies.

Linear Array and Ring