



# Introduction to Parallel Computing

George Karypis

Basic Communication Operations



# Outline

- Importance of Collective Communication Operations
- One-to-All Broadcast
- All-to-One Reduction
- All-to-All Broadcast & Reduction
- All-Reduce & Prefix-Sum
- Scatter and Gather
- All-to-All Personalized



# Collective Communication Operations

- They represent regular communication patterns that are performed by parallel algorithms.
  - Collective: Involve groups of processors
- Used extensively in most data-parallel algorithms.
- The parallel efficiency of these algorithms depends on efficient implementation of these operations.
- They are equally applicable to distributed and shared address space architectures
- Most parallel libraries provide functions to perform them
- They are extremely useful for “getting started” in parallel processing!

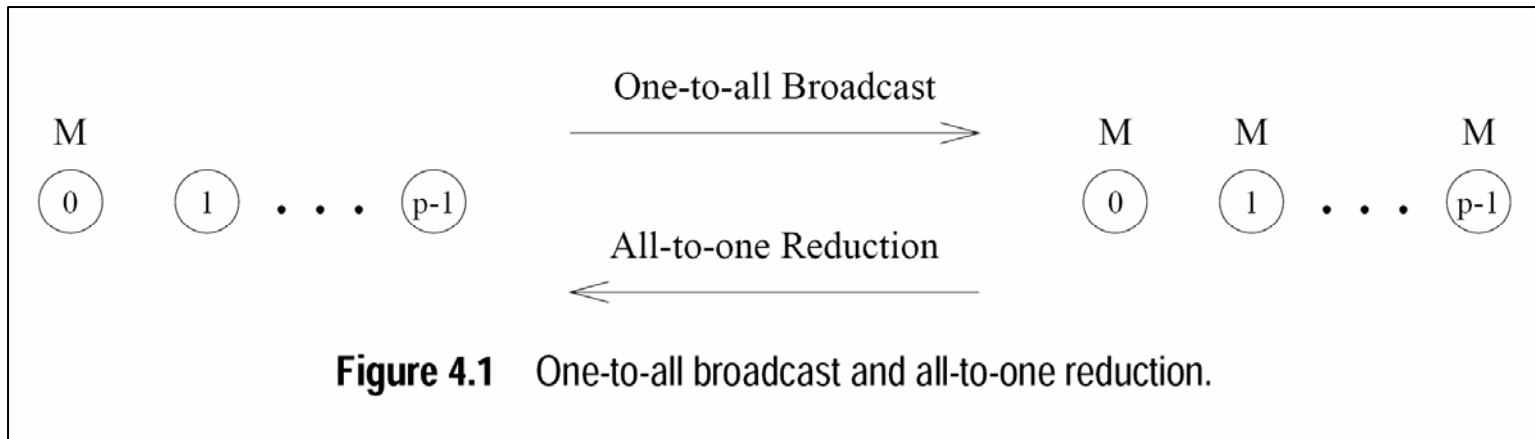


# MPI Names

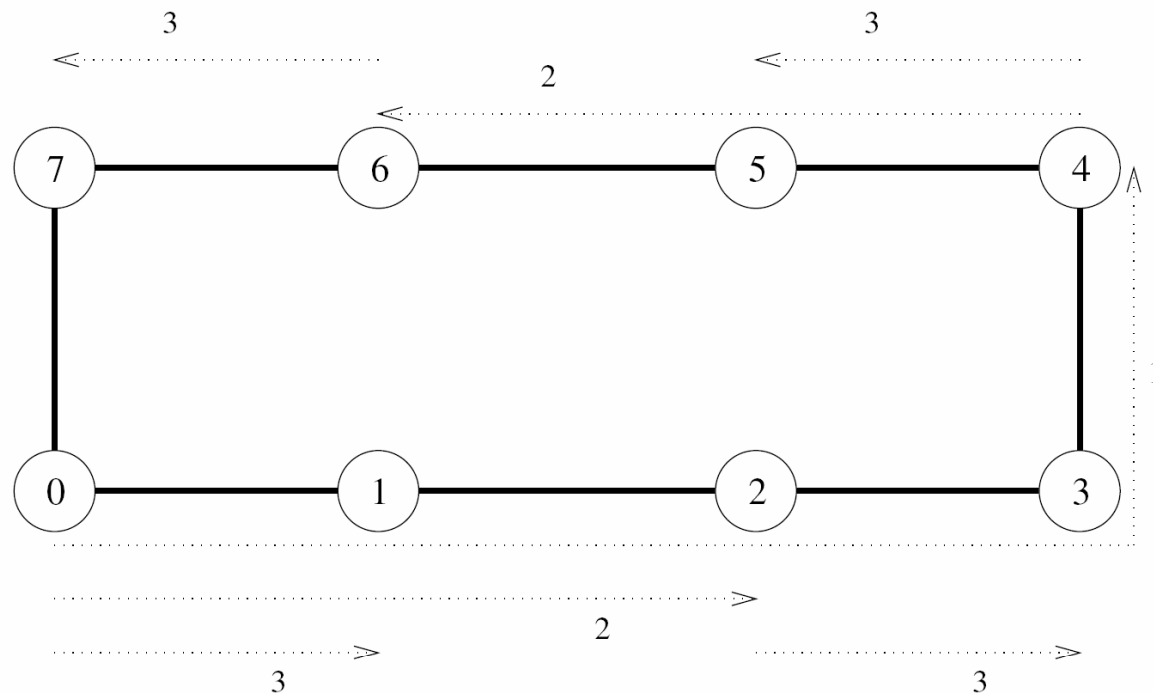
**Table 4.2** MPI names of the various operations discussed in this chapter.

Operation	MPI Name
One-to-all broadcast	MPI_Bcast
All-to-one reduction	MPI_Reduce
All-to-all broadcast	MPI_Allgather
All-to-all reduction	MPI_Reduce_scatter
All-reduce	MPI_Allreduce
Gather	MPI_Gather
Scatter	MPI_Scatter
All-to-all personalized	MPI_Alltoall

# One-to-All Broadcast & All-to-One Reduction

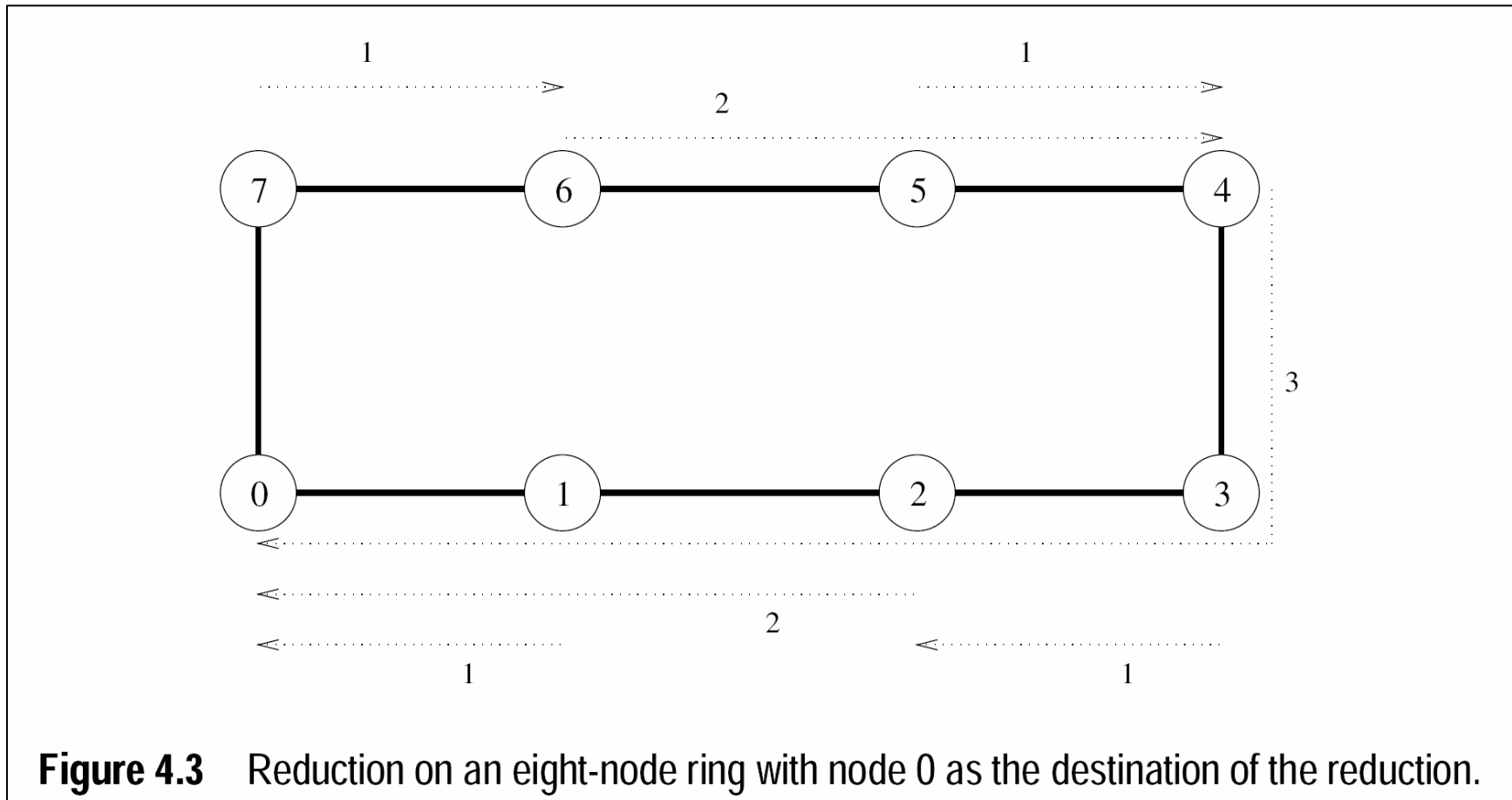


# Broadcast on a Ring Algorithm



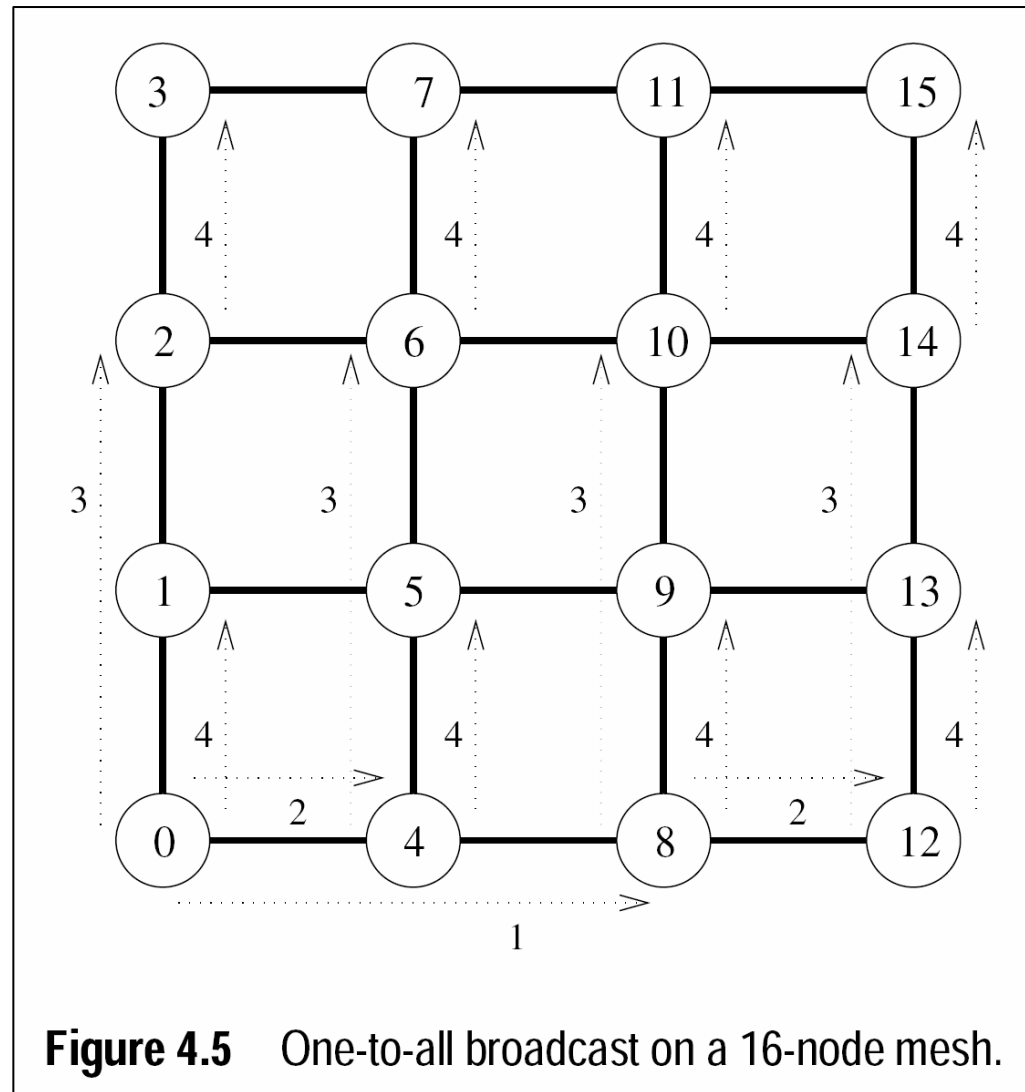
**Figure 4.2** One-to-all broadcast on an eight-node ring. Node 0 is the source of the broadcast. Each message transfer step is shown by a numbered, dotted arrow from the source of the message to its destination. The number on an arrow indicates the time step during which the message is transferred.

# Reduction on a Ring Algorithm



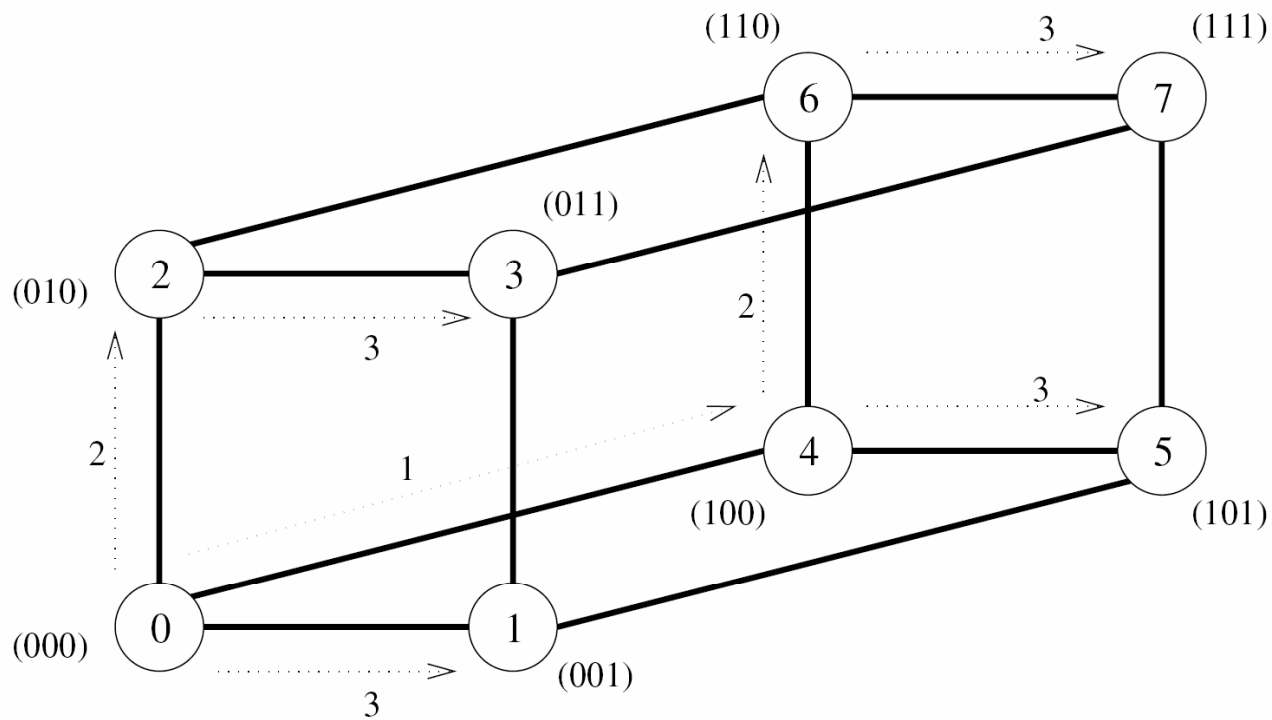
**Figure 4.3** Reduction on an eight-node ring with node 0 as the destination of the reduction.

# Broadcast on a Mesh





# Broadcast on a Hypercube



**Figure 4.6** One-to-all broadcast on a three-dimensional hypercube. The binary representations of node labels are shown in parentheses.

# Code for the Broadcast Source: Root

```
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 */
5.          mask := mask 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 XOR  $2^i$ ;
9.                  send X to msg_destination;
10.             else
11.                 msg_source := my_id XOR  $2^i$ ;
12.                 receive X from msg_source;
13.             endelse;
14.         endif;
15.     endfor;
16. end ONE_TO_ALL_BC
```

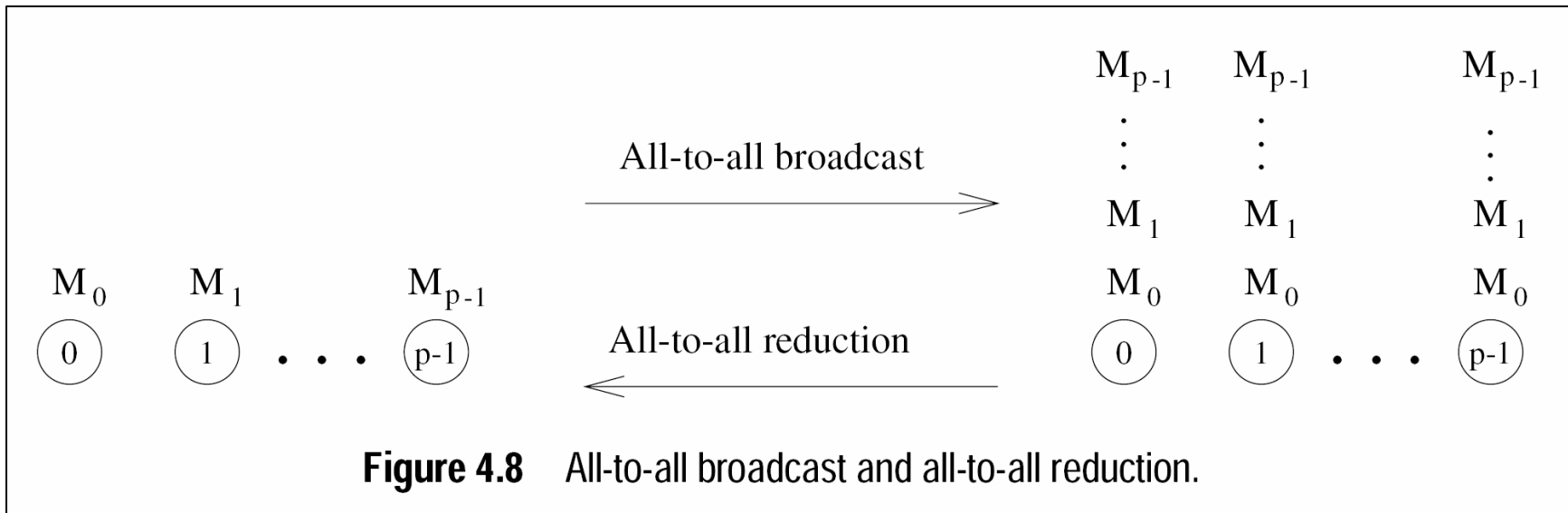
**Algorithm 4.1** One-to-all broadcast of a message  $X$  from node 0 of a  $d$ -dimensional  $p$ -node hypercube ( $d = \log p$ ). AND and XOR are bitwise logical-and and exclusive-or operations, respectively.

# Code for Broadcast Arbitrary Source

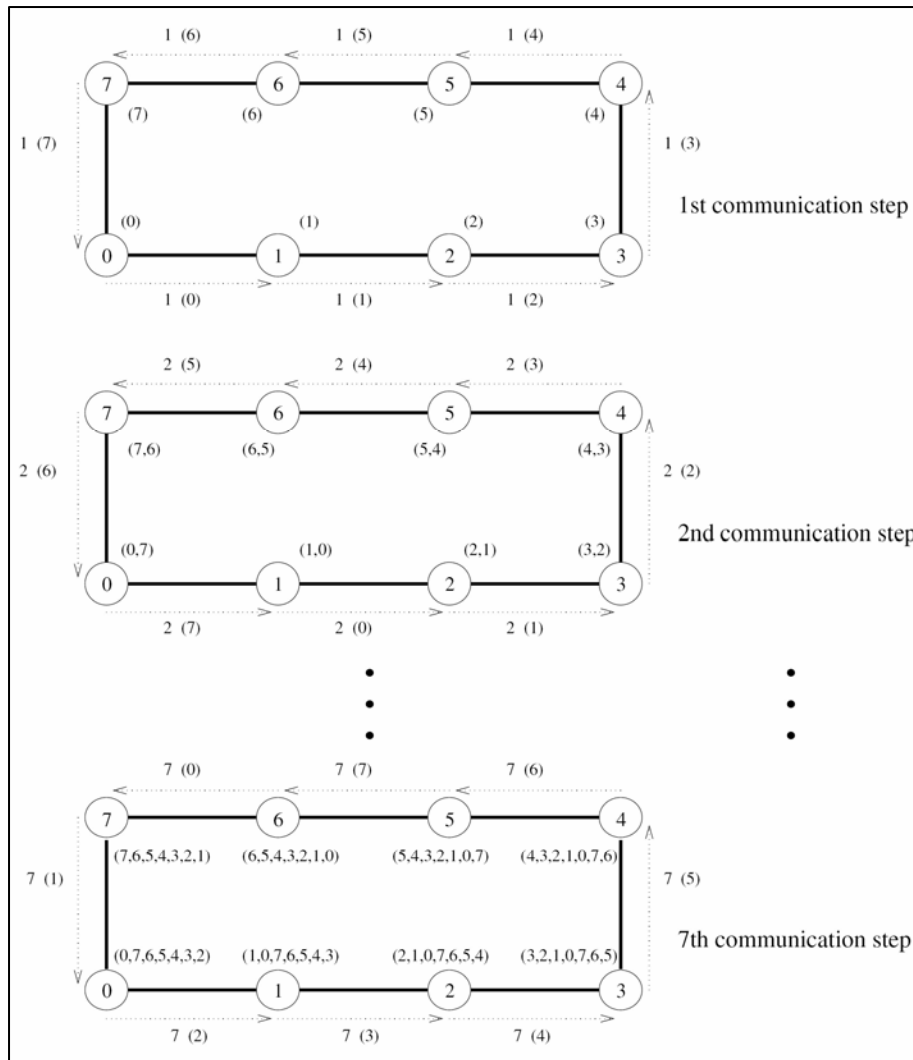
```
1.  procedure GENERAL_ONE_TO_ALL_BC(d, my_id, source, X)
2.  begin
3.    my_virtual_id := my_id XOR source;
4.    mask :=  $2^d - 1$ ;
5.    for i := d - 1 downto 0 do    /* Outer loop */
6.      mask := mask XOR  $2^i$ ;    /* Set bit i of mask to 0 */
7.      if (my_virtual_id AND mask) = 0 then
8.        if (my_virtual_id AND  $2^i$ ) = 0 then
9.          virtual_dest := my_virtual_id 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  $2^i$ ;
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
```

**Algorithm 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.

# All-to-All Broadcast & Reduction



# All-to-All Broadcast for Ring



```

1. procedure ALL_TO_ALL_BC_RING(my_id, my_msg, p, result)
2. begin
3.   left := (my_id - 1) mod p;
4.   right := (my_id + 1) mod p;
5.   result := my_msg;
6.   msg := result;
7.   for i := 1 to p - 1 do
8.     send msg to right;
9.     receive msg from left;
10.    result := result ∪ msg;
11.  endfor;
12. end ALL_TO_ALL_BC_RING
  
```

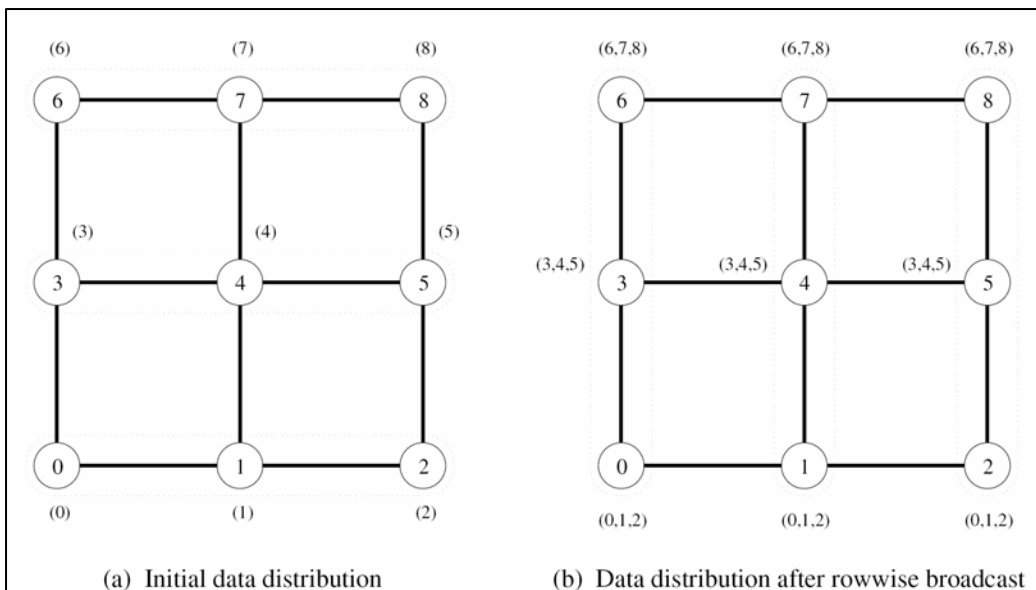
**Algorithm 4.4** All-to-all broadcast on a  $p$ -node ring.

```

1. procedure ALL_TO_ALL_RED_RING(my_id, my_msg, p, result)
2. begin
3.   left := (my_id - 1) mod p;
4.   right := (my_id + 1) mod p;
5.   recv := 0;
6.   for i := 1 to p - 1 do
7.     j := (my_id + i) mod p;
8.     temp := msg[j] + recv;
9.     send temp to left;
10.    receive recv from right;
11.  endfor;
12.  result := msg[my_id] + recv;
13. end ALL_TO_ALL_RED_RING
  
```

**Algorithm 4.5** All-to-all reduction on a  $p$ -node ring.

# All-to-All Broadcast on a Mesh



**Figure 4.10** All-to-all broadcast on a  $3 \times 3$  mesh. The groups of nodes communicating with each other in each phase are enclosed by dotted boundaries. By the end of the second phase, all nodes get  $(0,1,2,3,4,5,6,7)$  (that is, a message from each node).

```

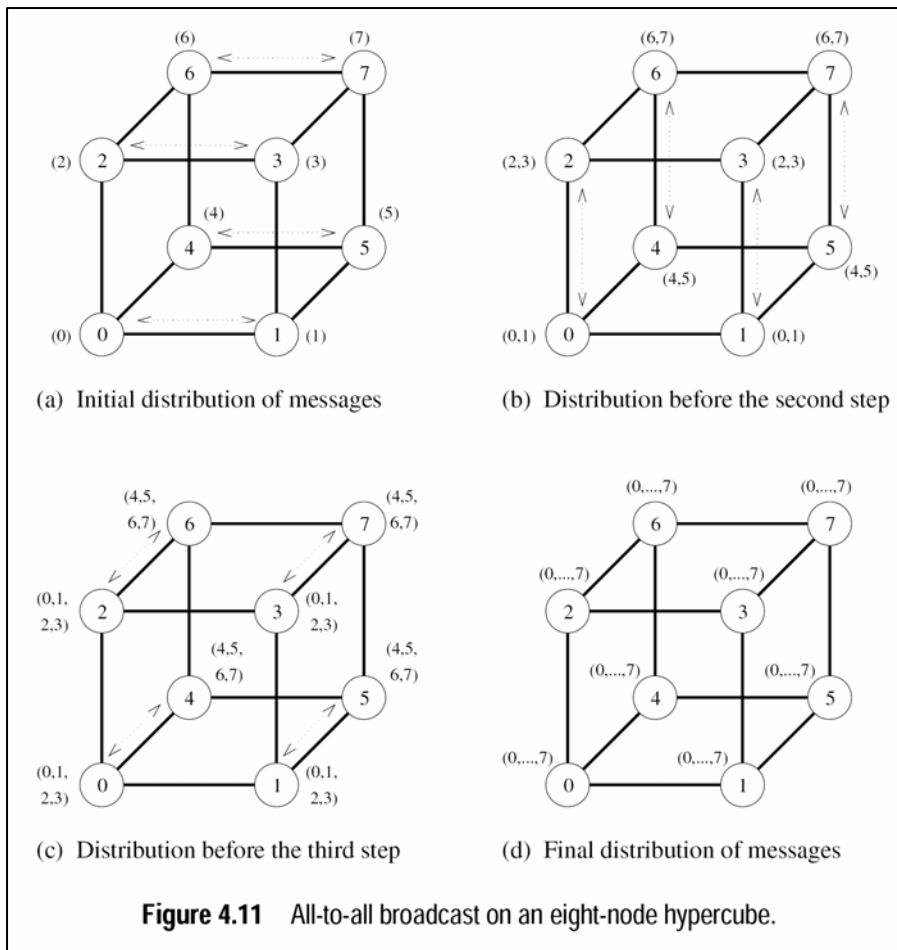
1. procedure ALL_TO_ALL_BC_MESH(my_id, my_msg, p, result)
2. begin

/* Communication along rows */
3.   left := my_id - (my_id mod  $\sqrt{p}$ ) + (my_id - 1) mod  $\sqrt{p}$ ;
4.   right := my_id - (my_id mod  $\sqrt{p}$ ) + (my_id + 1) mod  $\sqrt{p}$ ;
5.   result := my_msg;
6.   msg := result;
7.   for i := 1 to  $\sqrt{p} - 1$  do
8.     send msg to right;
9.     receive msg from left;
10.    result := result  $\cup$  msg;
11.  endfor;

/* Communication along columns */
12.  up := (my_id -  $\sqrt{p}$ ) mod p;
13.  down := (my_id +  $\sqrt{p}$ ) mod p;
14.  msg := result;
15.  for i := 1 to  $\sqrt{p} - 1$  do
16.    send msg to down;
17.    receive msg from up;
18.    result := result  $\cup$  msg;
19.  endfor;
20. end ALL_TO_ALL_BC_MESH
    
```

**Algorithm 4.6** All-to-all broadcast on a square mesh of  $p$  nodes.

# All-to-All Broadcast on a HCube



```

1. procedure ALL_TO_ALL_BC_HCUBE(my_id, my_msg, d, result)
2. begin
3.   result := my_msg;
4.   for i := 0 to d - 1 do
5.     partner := my_id XOR  $2^i$ ;
6.     send result to partner;
7.     receive msg from partner;
8.     result := result ∪ msg;
9.   endfor;
10. end ALL_TO_ALL_BC_HCUBE
  
```

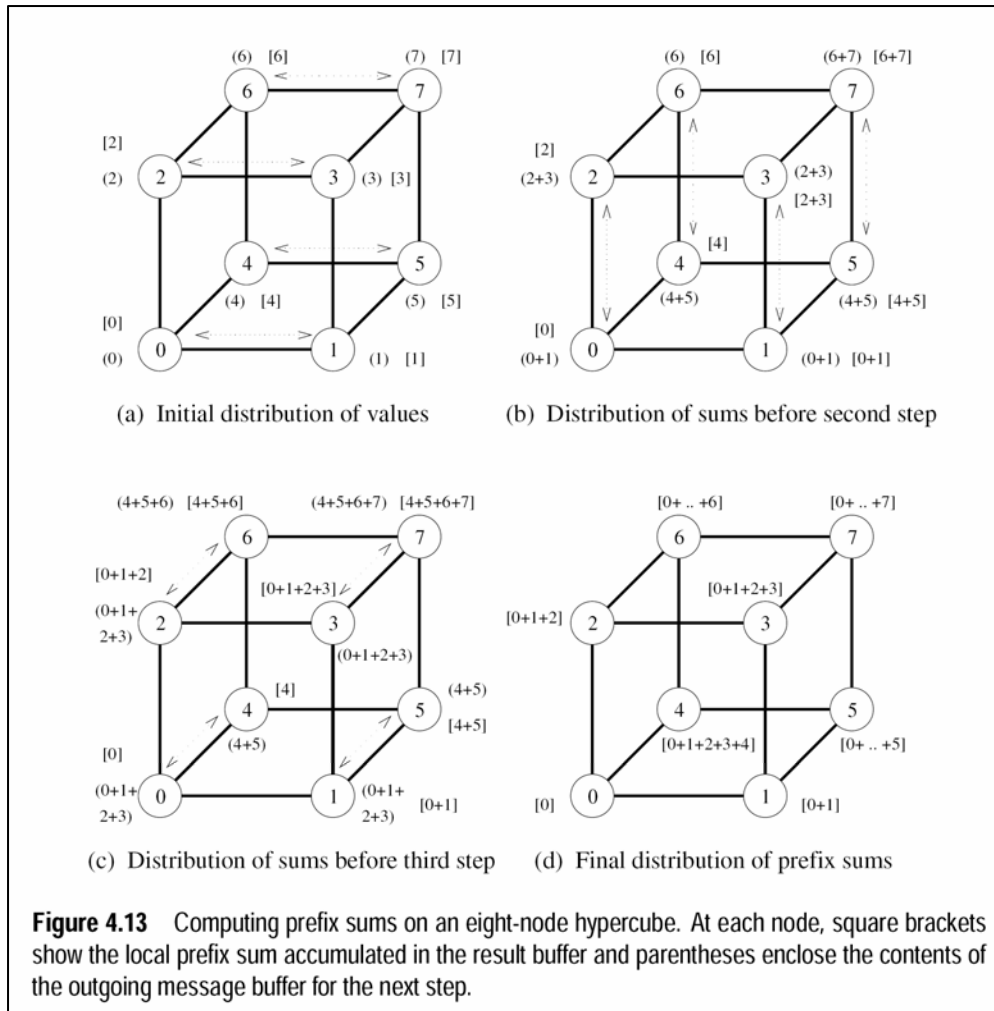
**Algorithm 4.7** All-to-all broadcast on a  $d$ -dimensional hypercube.

```

1. procedure ALL_TO_ALL_RED_HCUBE(my_id, msg, d, result)
2. begin
3.   reloc := 0;
4.   for i := d - 1 to 0 do
5.     partner := my_id XOR  $2^i$ ;
6.     j := my_id AND  $2^i$ ;
7.     k := (my_id XOR  $2^i$ ) AND  $2^i$ ;
8.     senloc := reloc + k;
9.     reloc := reloc + j;
10.    send msg[senloc .. senloc +  $2^i$  - 1] to partner;
11.    receive temp[0 ..  $2^i$  - 1] from partner;
12.    for j := 0 to  $2^i$  - 1 do
13.      msg[reloc + j] := msg[reloc + j] + temp[j];
14.    endfor;
15.  endfor;
16.  result := msg[my_id];
17. end ALL_TO_ALL_RED_HCUBE
  
```

**Algorithm 4.8** All-to-all broadcast on a  $d$ -dimensional hypercube. AND and XOR are bitwise logical-and and exclusive-or operations, respectively.

# All-Reduce & Prefix-Sum



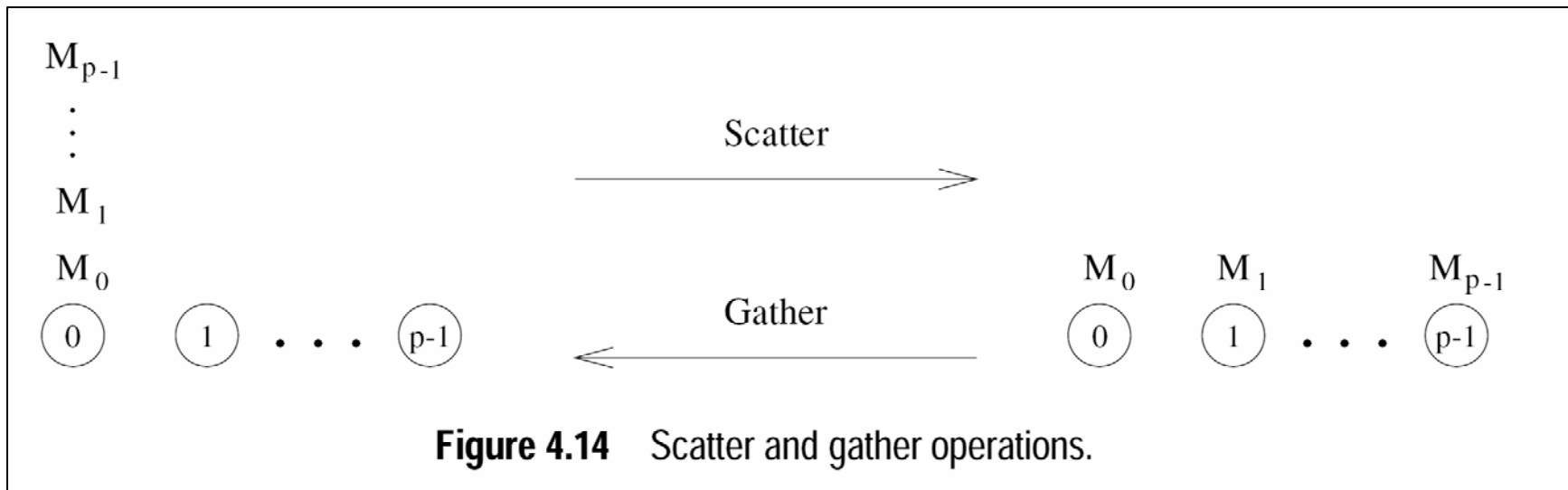
```

1. procedure PREFIX_SUMS_HCUBE(my_id, my_number, d, result)
2. begin
3.   result := my_number;
4.   msg := result;
5.   for i := 0 to d - 1 do
6.     partner := my_id XOR  $2^i$ ;
7.     send msg to partner;
8.     receive number from partner;
9.     msg := msg + number;
10.    if (partner < my_id) then result := result + number;
11.  endfor;
12. end PREFIX_SUMS_HCUBE
  
```

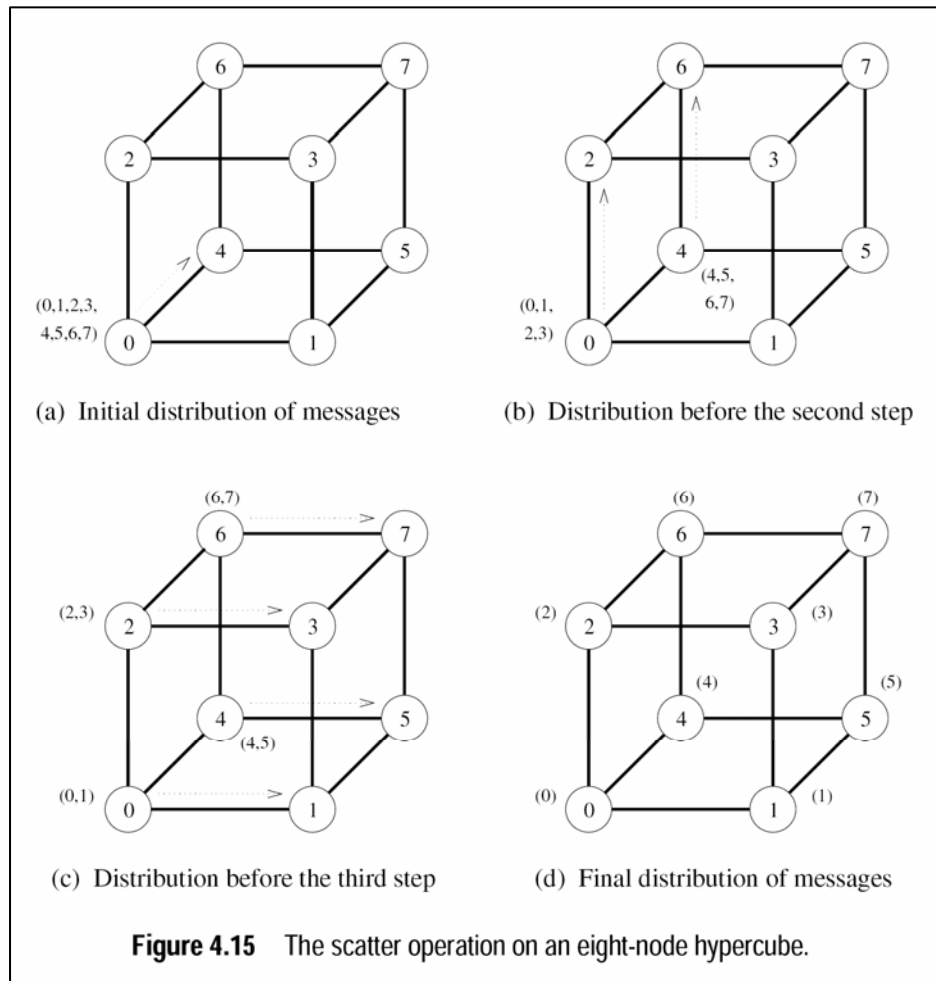
**Algorithm 4.9** Prefix sums on a  $d$ -dimensional hypercube.



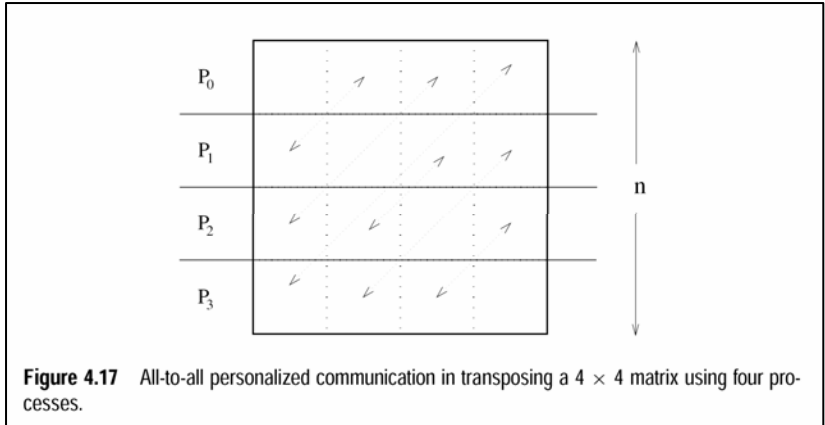
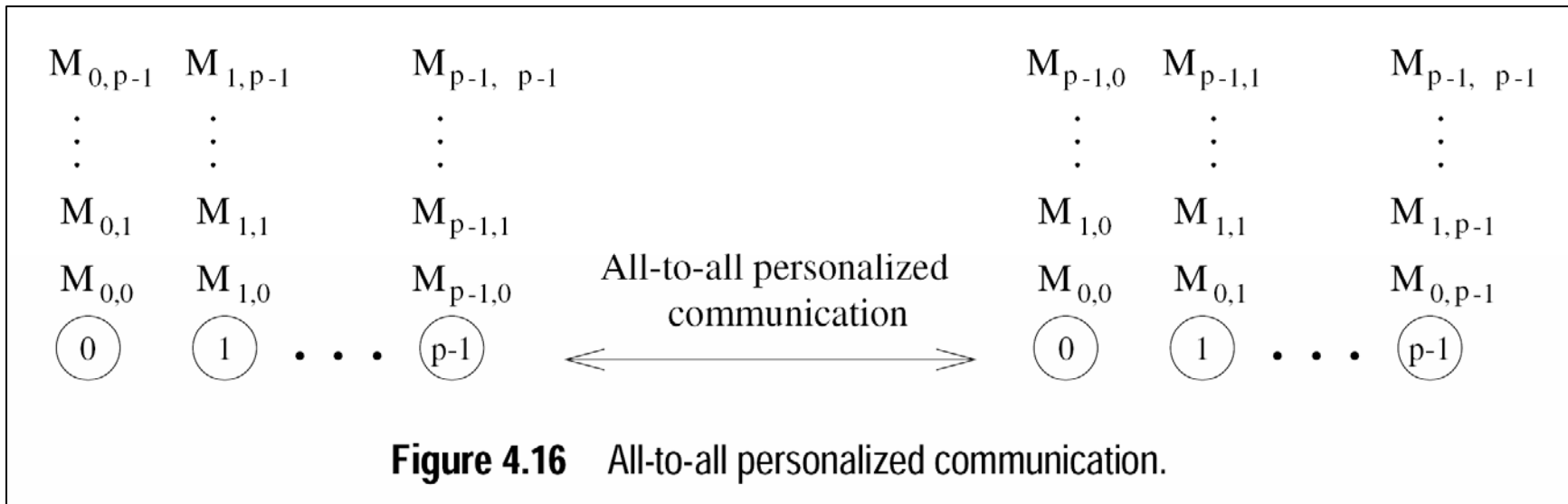
# Scatter & Gather



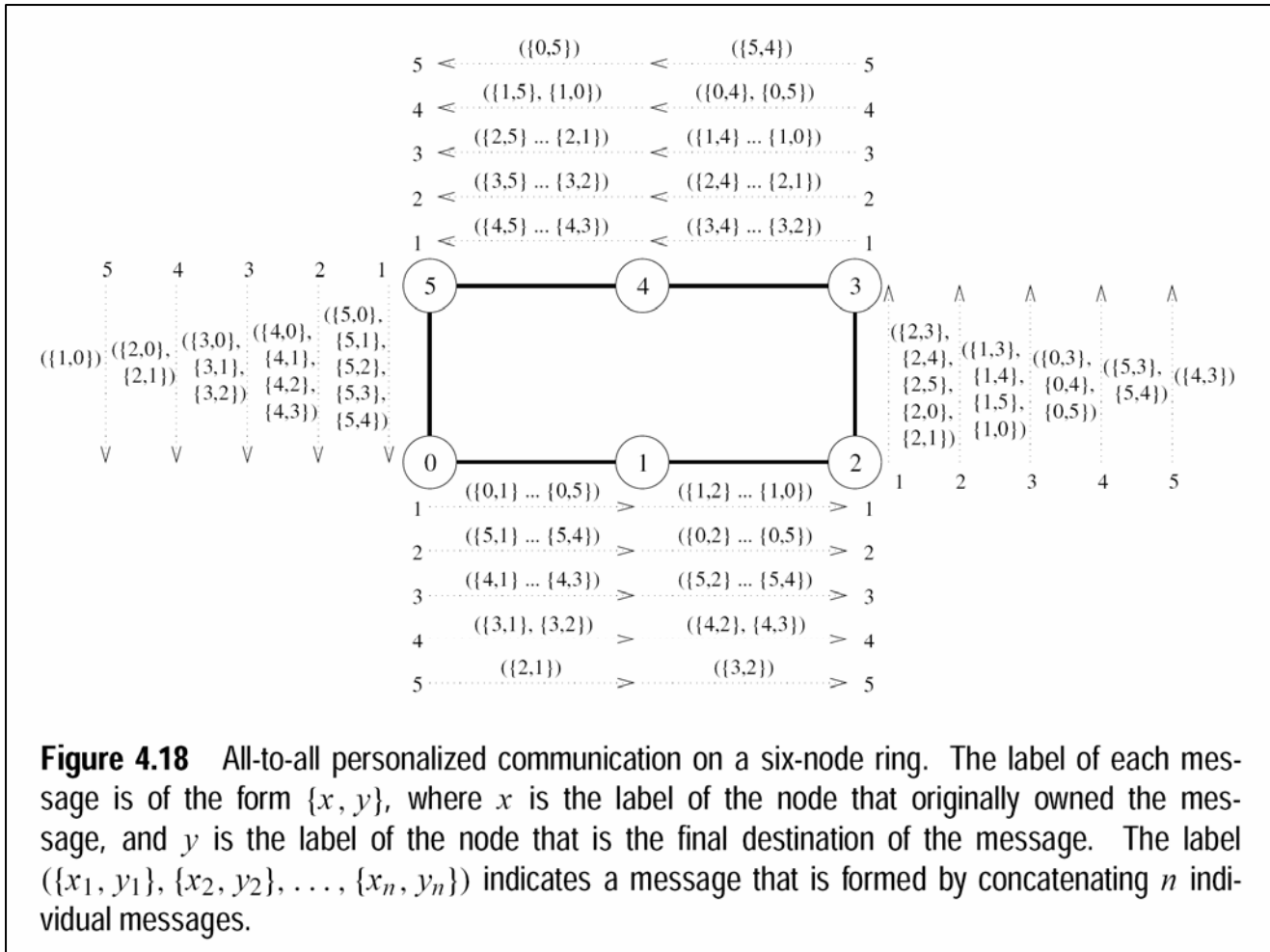
# Scatter Operation on HCube



# All-to-All Personalized (Transpose)

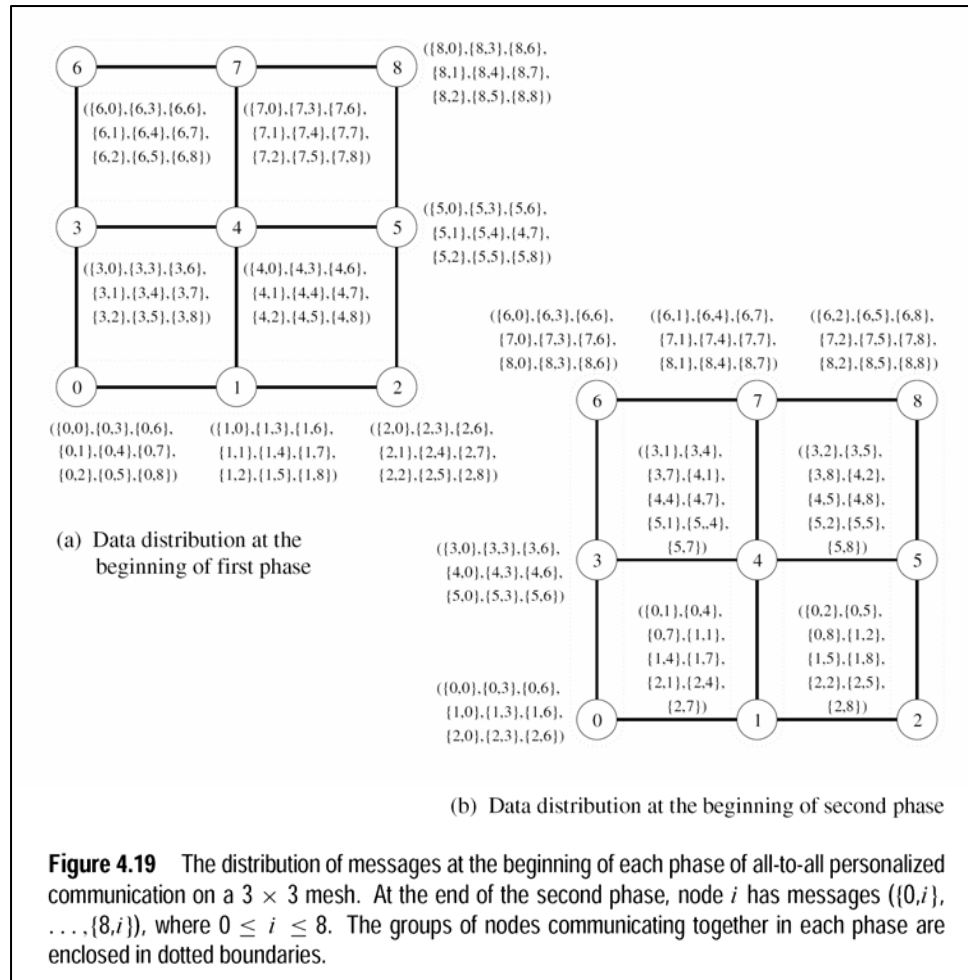


# All-to-all Personalized on a Ring

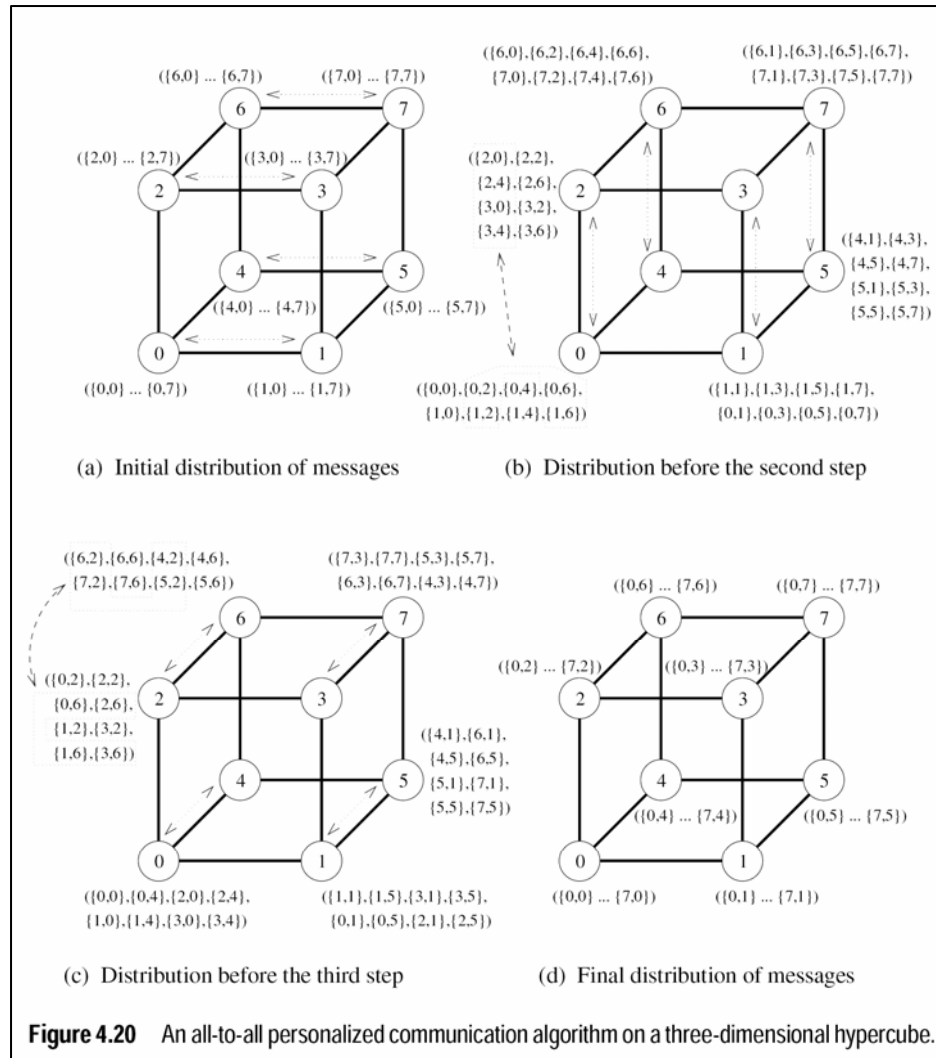


**Figure 4.18** All-to-all personalized communication on a six-node ring. The label of each message is of the form  $\{x, y\}$ , where  $x$  is the label of the node that originally owned the message, and  $y$  is the label of the node that is the final destination of the message. The label  $\{(x_1, y_1), \{x_2, y_2\}, \dots, \{x_n, y_n\}\}$  indicates a message that is formed by concatenating  $n$  individual messages.

# All-to-all Personalized on a Mesh



# All-to-all Personalized on a HCube



**Figure 4.20** An all-to-all personalized communication algorithm on a three-dimensional hypercube.

# All-to-all Personalized on a HCube Improved Algorithm

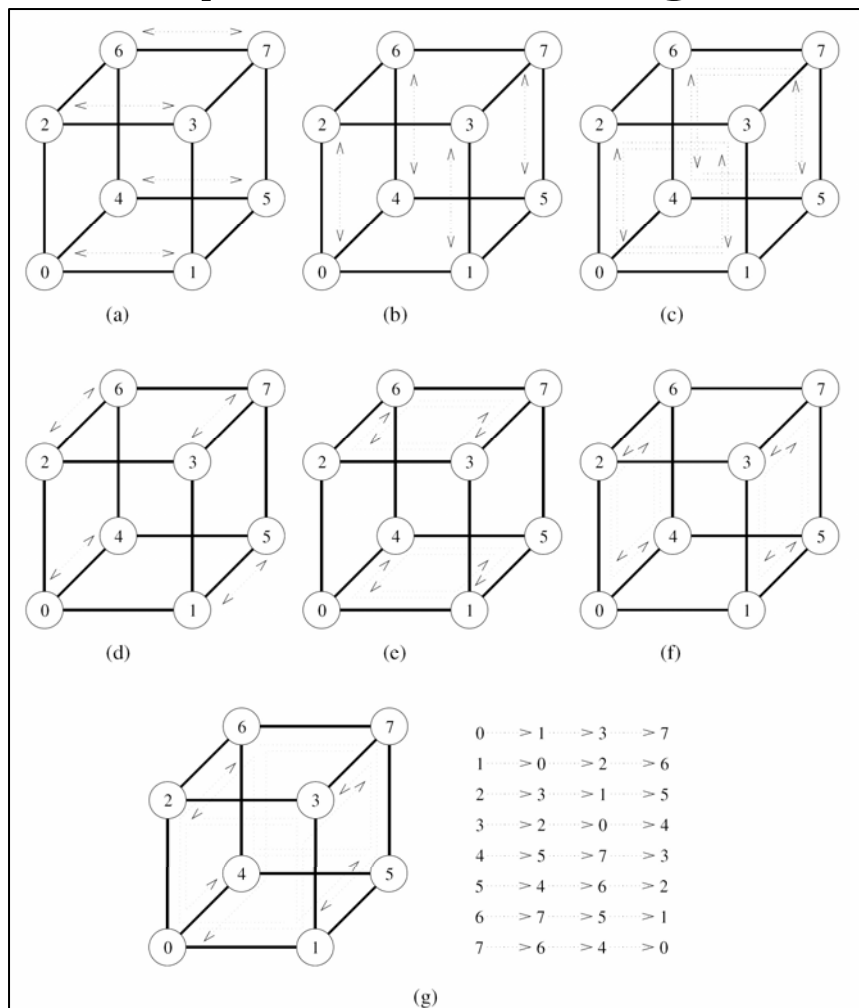


Figure 4.21 Seven steps in all-to-all personalized communication on an eight-node hypercube.

Perform  $\log(p)$  point-to-point communication steps

Processor  $i$  communicates with processor  $iXORj$  during the  $j$ th communication step.

```

1. procedure ALL_TO_ALL_PERSONAL( $d, my\_id$ )
2. begin
3.   for  $i := 1$  to  $2^d - 1$  do
4.     begin
5.        $partner := my\_id XOR i$ ;
6.       send  $M_{my\_id, partner}$  to  $partner$ ;
7.       receive  $M_{partner, my\_id}$  from  $partner$ ;
8.     endfor;
9.   end ALL_TO_ALL_PERSONAL
    
```

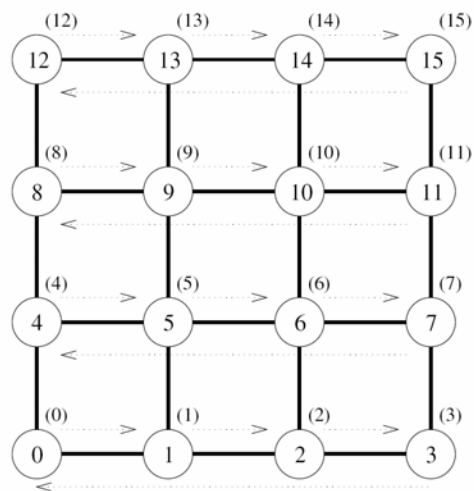
**Algorithm 4.10** A procedure to perform all-to-all personalized communication on a  $d$ -dimensional hypercube. The message  $M_{i,j}$  initially resides on node  $i$  and is destined for node  $j$ .

# Complexities

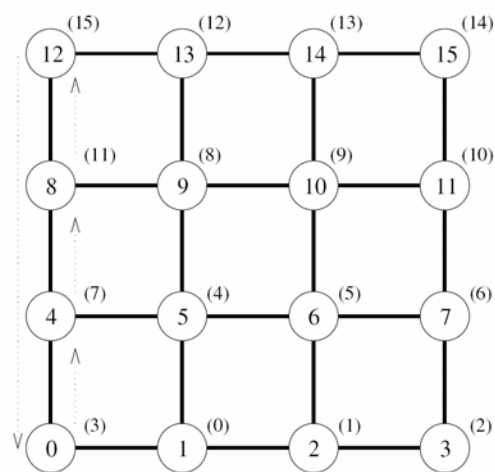
**Table 4.1** Summary of communication times of various operations discussed in Sections 4.1–4.7 on a hypercube interconnection network. The message size for each operation is  $m$  and the number of nodes is  $p$ .

Operation	Hypercube Time	B/W Requirement
One-to-all broadcast, All-to-one reduction	$\min((t_s + t_w m) \log p, 2(t_s \log p + t_w m))$	$\Theta(1)$
All-to-all broadcast, All-to-all reduction	$t_s \log p + t_w m(p - 1)$	$\Theta(1)$
All-reduce	$\min((t_s + t_w m) \log p, 2(t_s \log p + t_w m))$	$\Theta(1)$
Scatter, Gather	$t_s \log p + t_w m(p - 1)$	$\Theta(1)$
All-to-all personalized	$(t_s + t_w m)(p - 1)$	$\Theta(p)$
Circular shift	$t_s + t_w m$	$\Theta(p)$

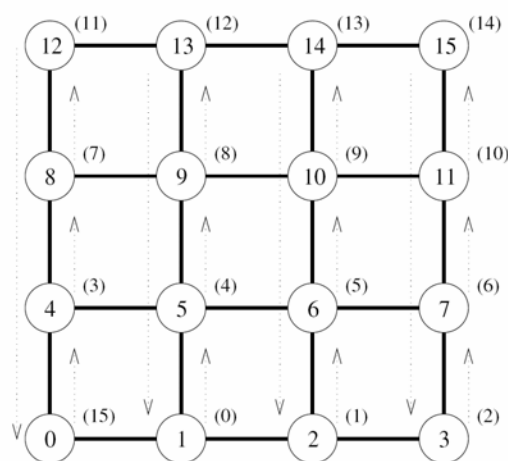




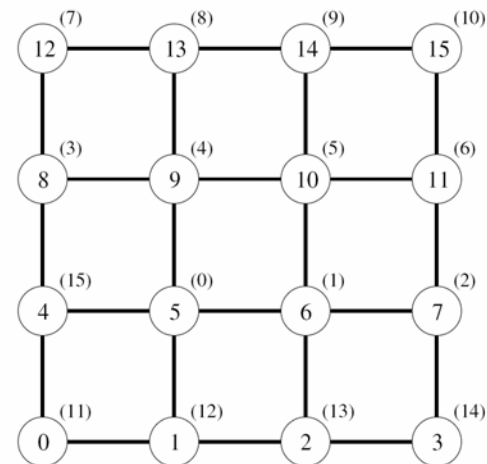
(a) Initial data distribution and the first communication step



(b) Step to compensate for backward row shifts

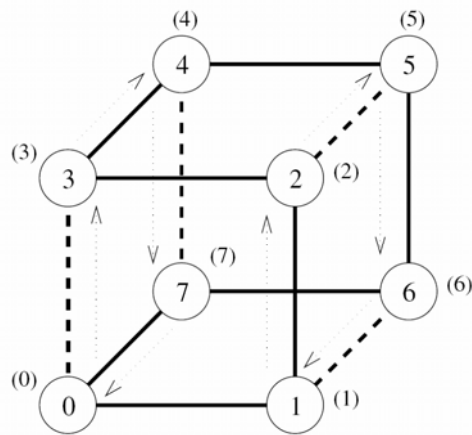


(c) Column shifts in the third communication step

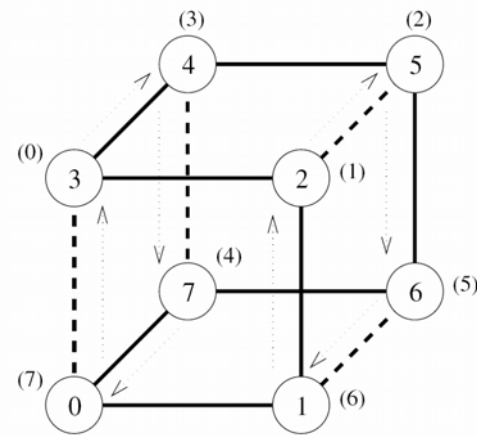


(d) Final distribution of the data

**Figure 4.22** The communication steps in a circular 5-shift on a  $4 \times 4$  mesh.

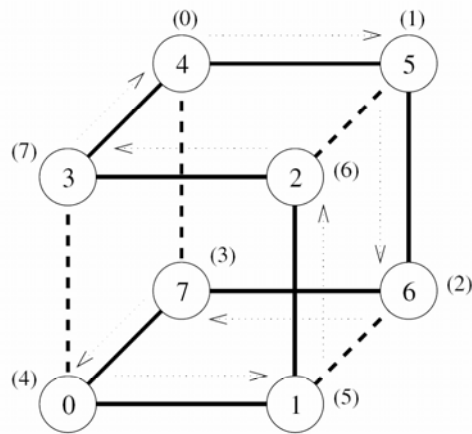


First communication step of the 4-shift

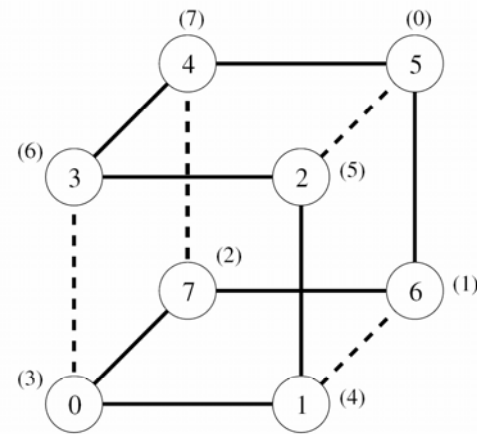


Second communication step of the 4-shift

(a) The first phase (a 4-shift)

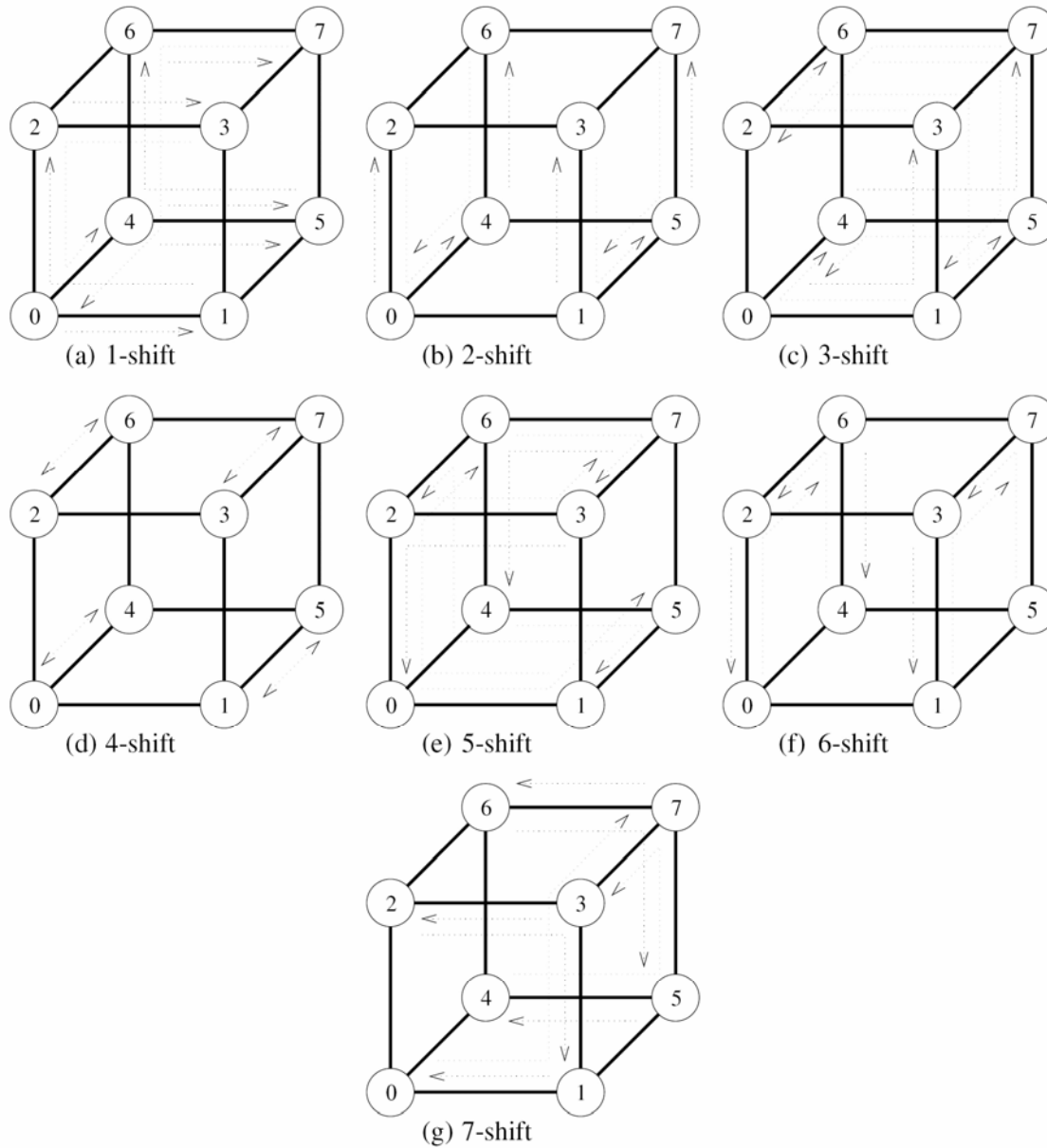


(b) The second phase (a 1-shift)



(c) Final data distribution after the 5-shift

**Figure 4.23** The mapping of an eight-node linear array onto a three-dimensional hypercube to perform a circular 5-shift as a combination of a 4-shift and a 1-shift.



**Figure 4.24** Circular  $q$ -shifts on an 8-node hypercube for  $1 \leq q < 8$ .