



Basic Communication Operations

Alexandre David
B2-206



Today

- One-to-all broadcast & all-to-one reduction (4.1).
- All-to-all broadcast and reduction (4.2).
- All-reduce and prefix-sum operations (4.3).

Collective Communication Operations

- Represent regular communication patterns.
- Used extensively in most data-parallel algorithms.
- Critical for efficiency.
- Available in most parallel libraries.
- Very useful to “get started” in parallel processing.

28-02-2006

Alexandre David, MVP'06

3


Collective: involve group of processors.

The efficiency of data-parallel algorithms depends on the efficient implementation of these operations.

Recall: $t_s + mt_w$ time for exchanging a m -word message with cut-through routing.

All processes participate in a single **global** interaction operation or subsets of processes in **local** interactions.

Goal of this chapter: good algorithms to implement commonly used communication patterns.



Reminder

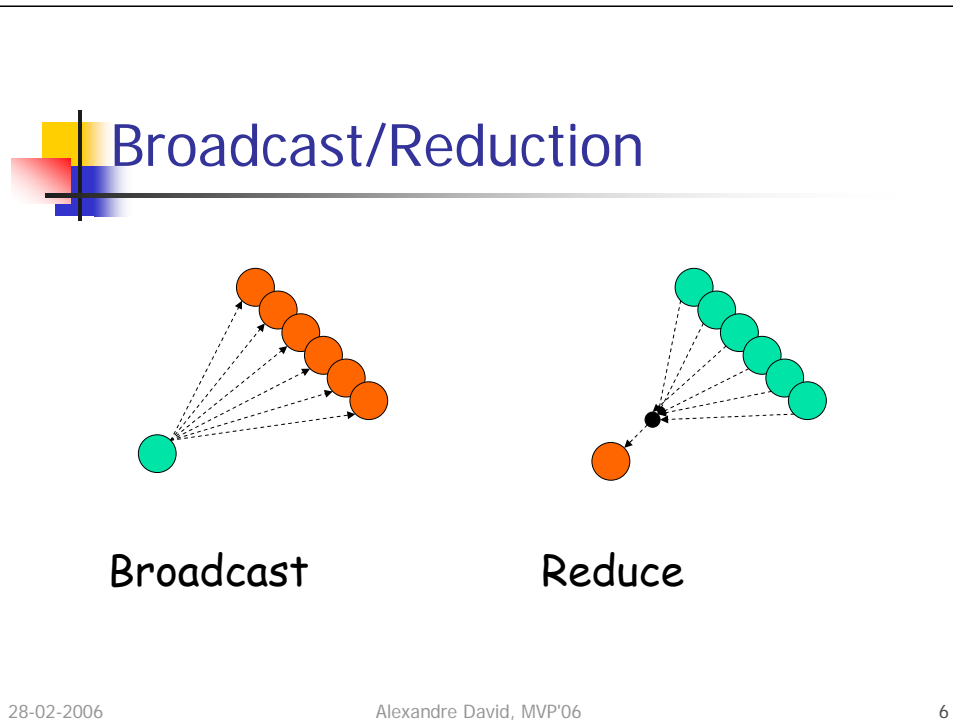
- Result from previous analysis:
 - Data transfer time is roughly the same between *all pairs* of nodes.
 - Homogeneity true on modern hardware (randomized routing, cut-through routing...).
 - $t_s + mt_w$
 - Adjust t_w for congestion: effective t_w .
- Model: bidirectional links, single port.
- Communication with point-to-point primitives.



Broadcast/Reduction

- One-to-all broadcast:
 - Single process sends identical data to all (or subset of) processes.
- All-to-one reduction:
 - Dual operation.
 - P processes have m words to send to one destination.
 - Parts of the message need to be *combined*.

Reduction can be used to find the sum, product, maximum, or minimum of sets of numbers.



This is the logical view, what happens from the programmer's perspective.



One-to-All Broadcast – Ring/Linear Array

- Naïve approach: send sequentially.
 - Bottleneck.
 - Poor utilization of the network.
- Recursive doubling:
 - Broadcast in $\log p$ steps (instead of p).
 - Divide-and-conquer type of algorithm.
 - Reduction is similar.

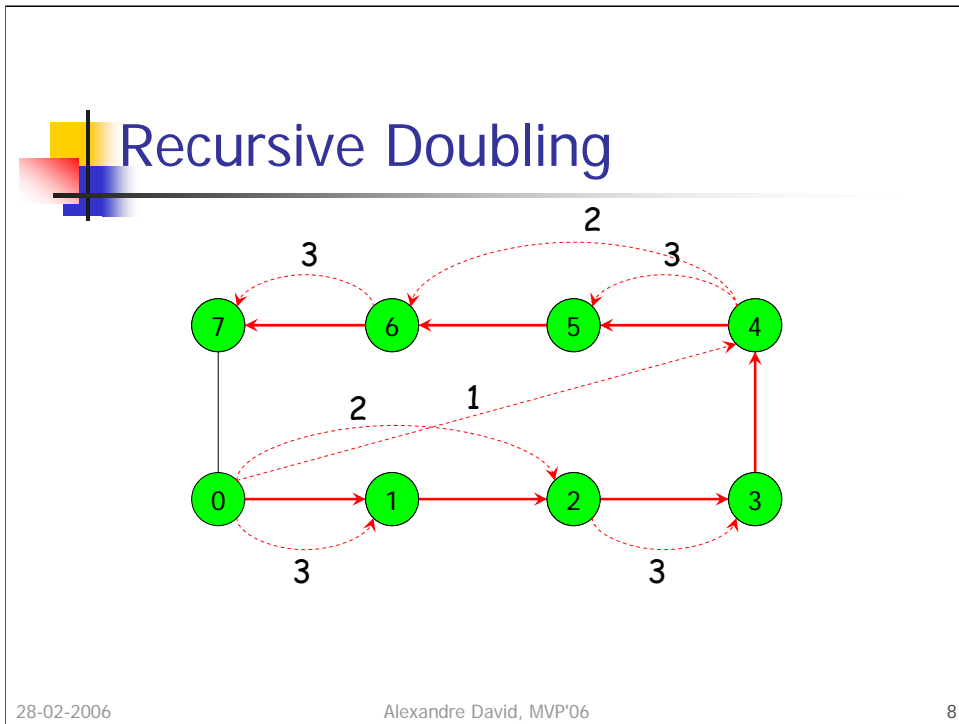
28-02-2006

Alexandre David, MVP'06

7

Source *process* is the bottleneck. Poor utilization: Only connections between single pairs of nodes are used at a time.

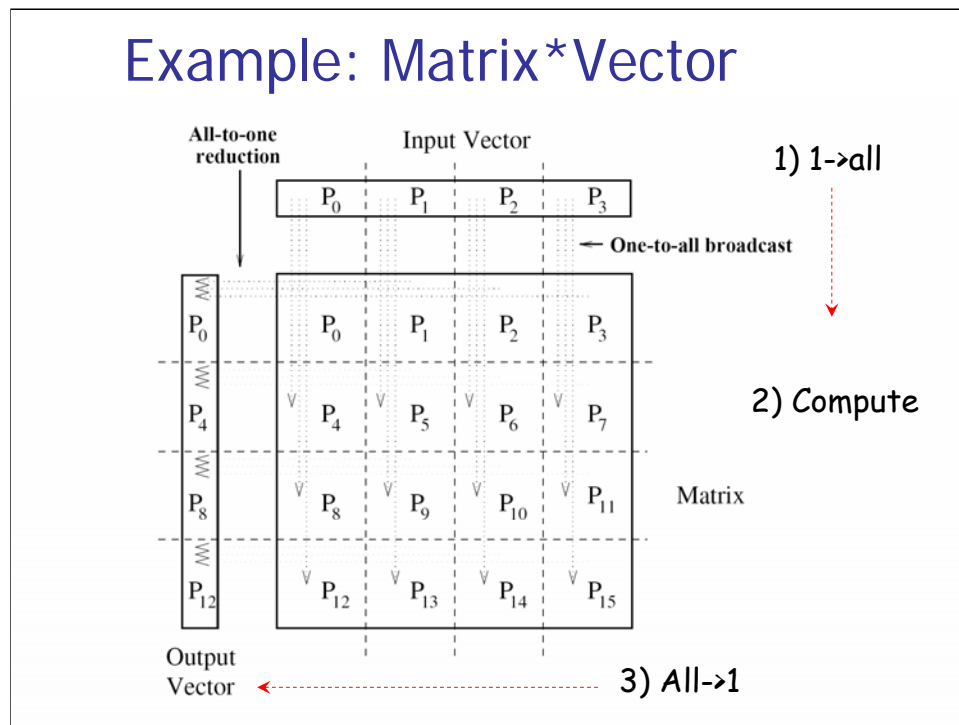
Recursive doubling: All processes that have the data can send it again.



Note:

- The nodes do not snoop the messages going “through” them. Messages are forwarded but the processes are not notified of this because they are not destined to them.
- Choose carefully destinations: furthest.
- Reduction symmetric: Accumulate results and send with the same pattern.

Example: Matrix*Vector



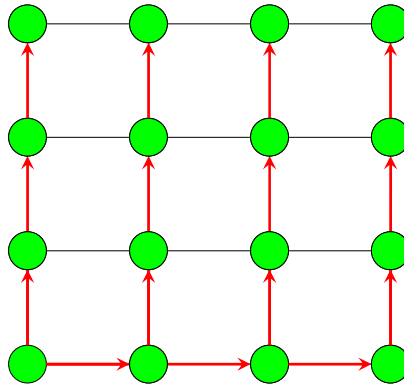
Although we have a matrix & a vector the broadcast are done on arrays.



One-to-All Broadcast – Mesh

- Extensions of the linear array algorithm.
 - Rows & columns = arrays.
 - Broadcast on a row, broadcast on columns.
 - Similar for reductions.
 - Generalize for higher dimensions (cubes...).

Broadcast on a Mesh



28-02-2006

Alexandre David, MVP'06

11

1. Broadcast like linear array.
2. Every node on the linear array has the data and broadcast on the columns with the linear array algorithm, *in parallel*.

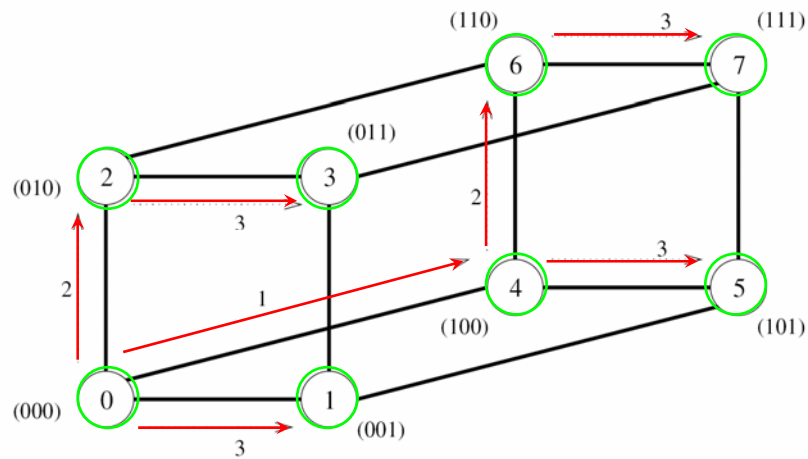
One-to-All Broadcast –



Hypercube

- Hypercube with 2^d nodes = d -dimensional mesh with 2 nodes in each direction.
- Similar algorithm in d steps.
- Also in $\log p$ steps.
- Reduction follows the same pattern.

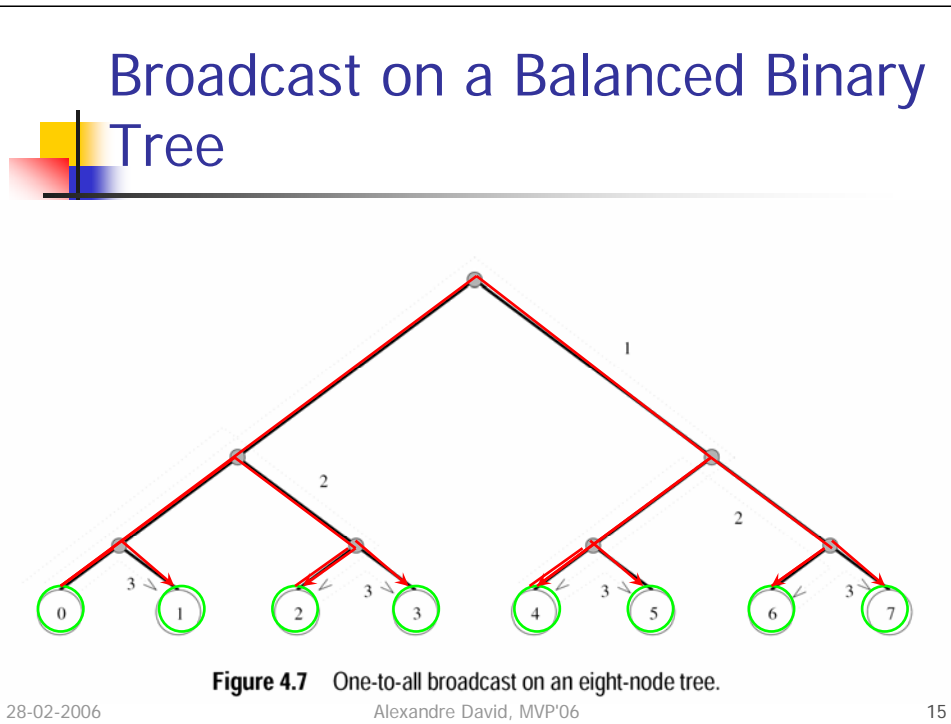
Broadcast on a Hypercube



Better for congestion: Use different links every time. Forwarding in parallel again.

All-to-One Broadcast – Balanced Binary Tree

- Processing nodes = leaves.
- Hypercube algorithm maps well.
- Similarly good w.r.t. congestion.



Divide-and-conquer type of algorithm again.

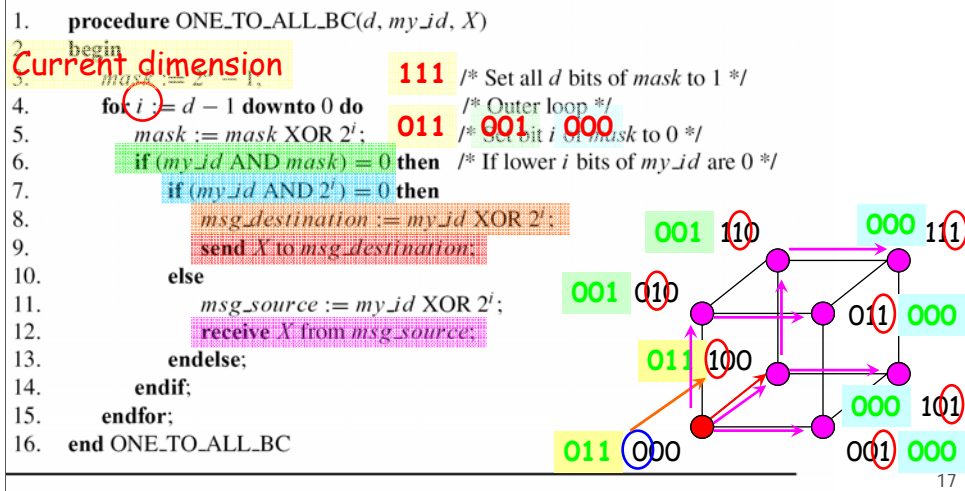


Algorithms

- So far we saw pictures.
- Not enough to implement.
- Precise description
 - to implement.
 - to analyze.
- Description for hypercube.
- Execute the following procedure on all the nodes.

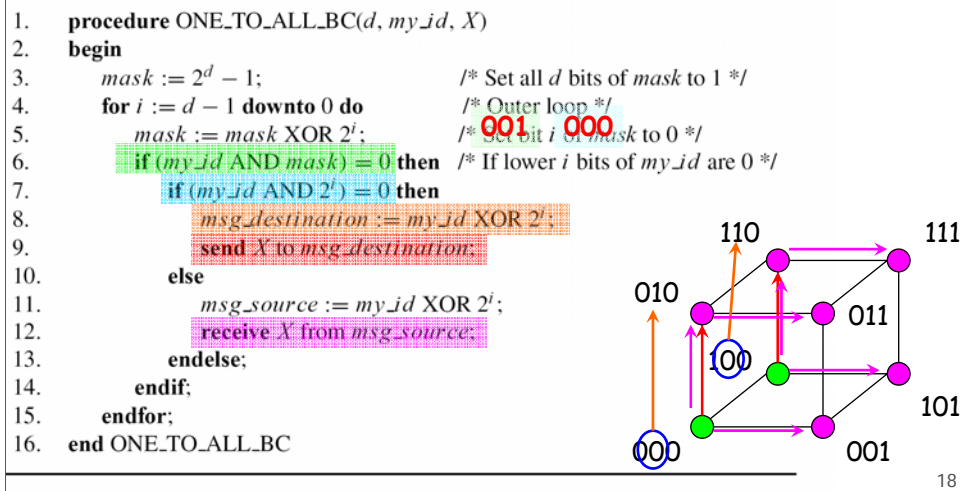
For sake of simplicity, the number of nodes is a power of 2.

Broadcast Algorithm



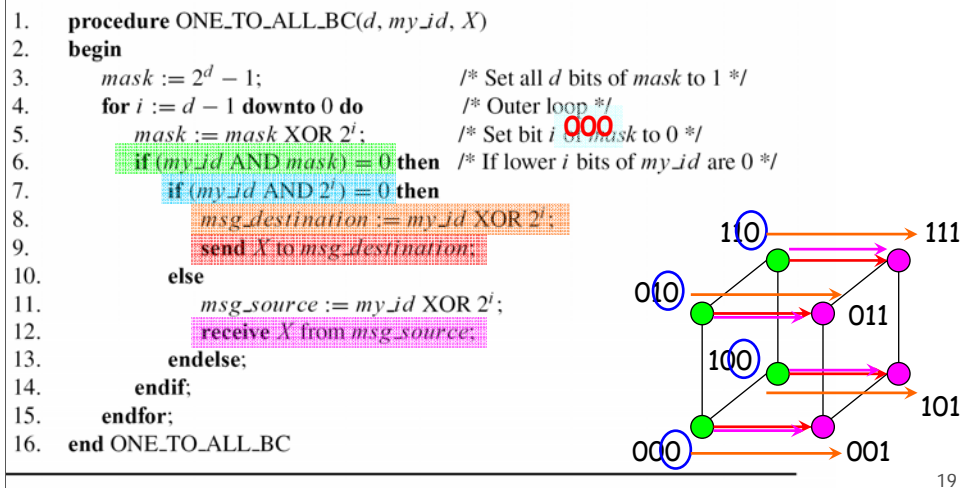
my_id is the label of the node the procedure is executed on. The procedure performs d communication steps, one along each dimension of the hypercube. Nodes with zero in i least significant bits (of their labels) participate in the communication.

Broadcast Algorithm



my_id is the label of the node the procedure is executed on. The procedure performs d communication steps, one along each dimension of the hypercube. Nodes with zero in i least significant bits (of their labels) participate in the communication.

Broadcast Algorithm



my_id is the label of the node the procedure is executed on. The procedure performs d communication steps, one along each dimension of the hypercube. Nodes with zero in i least significant bits (of their labels) participate in the communication.

Notes:

- Every node has to know when to communicate, i.e., call the procedure.
- The procedure is distributed and requires only point-to-point synchronization.
- Only from node 0.



Algorithm For Any Source

```
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 ( $my\_virtual\_id \text{ AND } mask$ ) = 0 then
8.       if ( $my\_virtual\_id \text{ AND } 2^i$ ) = 0 then
9.          $virtual\_dest := my\_virtual\_id \text{ XOR } 2^i$ ;
10.        send  $X$  to ( $virtual\_dest \text{ XOR } source$ );
11.        /* Convert  $virtual\_dest$  to the label of the physical destination */
12.      else
13.         $virtual\_source := my\_virtual\_id \text{ XOR } 2^i$ ;
14.        receive  $X$  from ( $virtual\_source \text{ XOR } source$ );
15.        /* Convert  $virtual\_source$  to the label of the physical source */
16.      endelse;
17.    endfor;
18. end GENERAL_ONE_TO_ALL_BC
```

20

XOR the source = renaming relative to the source. Still works because of the sub-cube property: changing 1 bit = navigate on one dimension, keep a set of equal bits = sub-cube.

Reduce Algorithm

```
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
6.     /* Select nodes whose lower  $i$  bits are 0 */
7.     if  $(my\_id \text{ AND } mask) = 0$  then
8.       if  $(my\_id \text{ AND } 2^i) \neq 0$  then
9.          $msg\_destination := my\_id \text{ XOR } 2^i$ ;
10.        In a nutshell:
11.        reverse the previous one.
12.        receive  $X$  from  $msg\_source$ ;
13.        for  $j := 0$  to  $m - 1$  do
14.           $sum[j] := sum[j] + X[j]$ ;
15.        endelse;
16.         $mask := mask \text{ XOR } 2^i$ ; /* Set bit  $i$  of  $mask$  to 1 */
17.      endfor;
18.   end ALL_TO_ONE_REDUCE
```



Cost Analysis

p processes $\rightarrow \log p$ steps (point-to-point transfers in parallel).
Each transfer has a time cost of $t_s + t_w m$.
Total time: $T = (t_s + t_w m) \log p$.

All-to-All Broadcast and Reduction

- Generalization of broadcast:
 - Each processor is a source and destination.
 - Several processes broadcast different messages.
- Used in matrix multiplication (and matrix-vector multiplication).
- Dual: all-to-all reduction.

28-02-2006

Alexandre David, MVP'06

23

How to do it?

If performed naively, it may take up to p times as long as a one-to-all broadcast (for p processors).

Possible to concatenate all messages that are going through the same path (reduce time because fewer t_s).

All-to-All Broadcast and Reduction

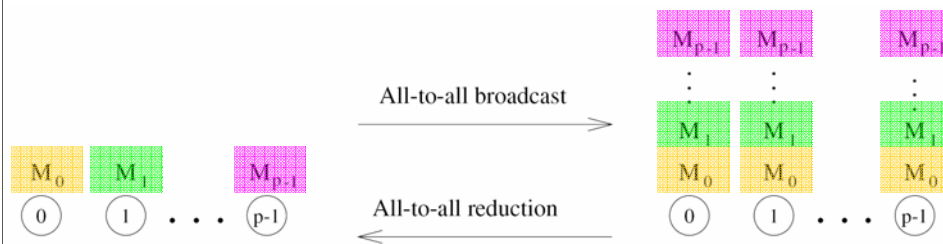
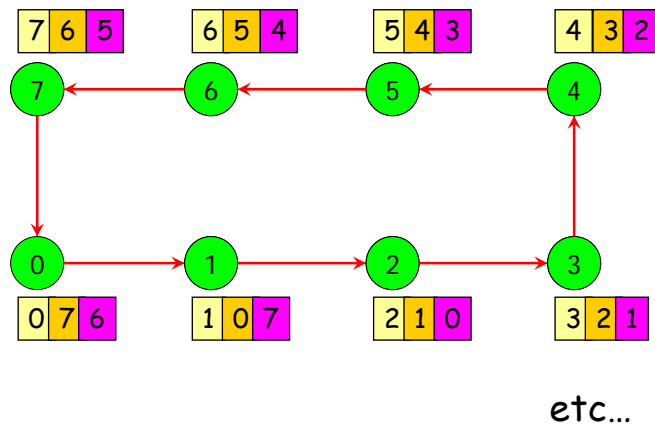


Figure 4.8 All-to-all broadcast and all-to-all reduction.

All-to-All Broadcast – Rings



28-02-2006

Alexandre David, MVP'06

25

All communication links can be kept busy until the operation is complete because each node has some information to pass. One-to-all in $\log p$ steps, all-to-all in $p-1$ steps instead of $p \log p$ (naïve).

How to do it for linear arrays? If we have bidirectional links (assumption from the beginning), we can use the same procedure.



All-to-All Broadcast Algorithm

```
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
```

Ring: mod p .
Receive & send - point-to-point.
Initialize the loop.
Forward msg.
Accumulate result.

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



All-to-All Reduce Algorithm

```
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
```

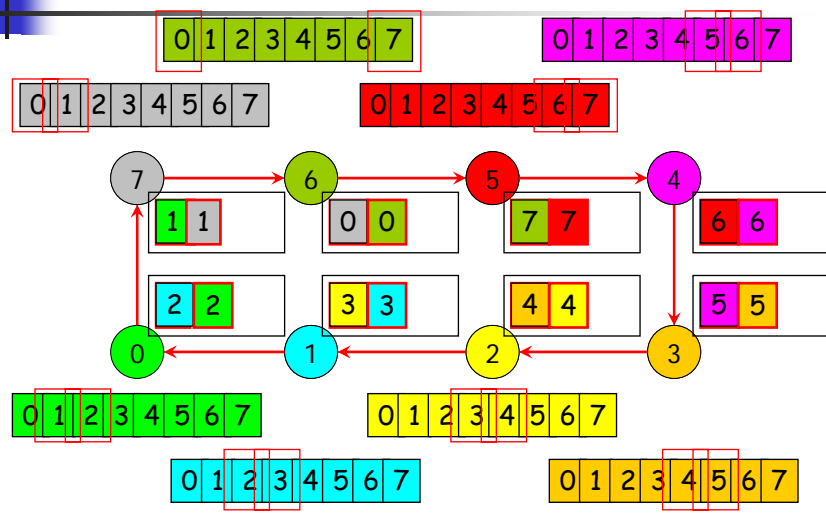
Accumulate and forward.

Last message for *my_id*.

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



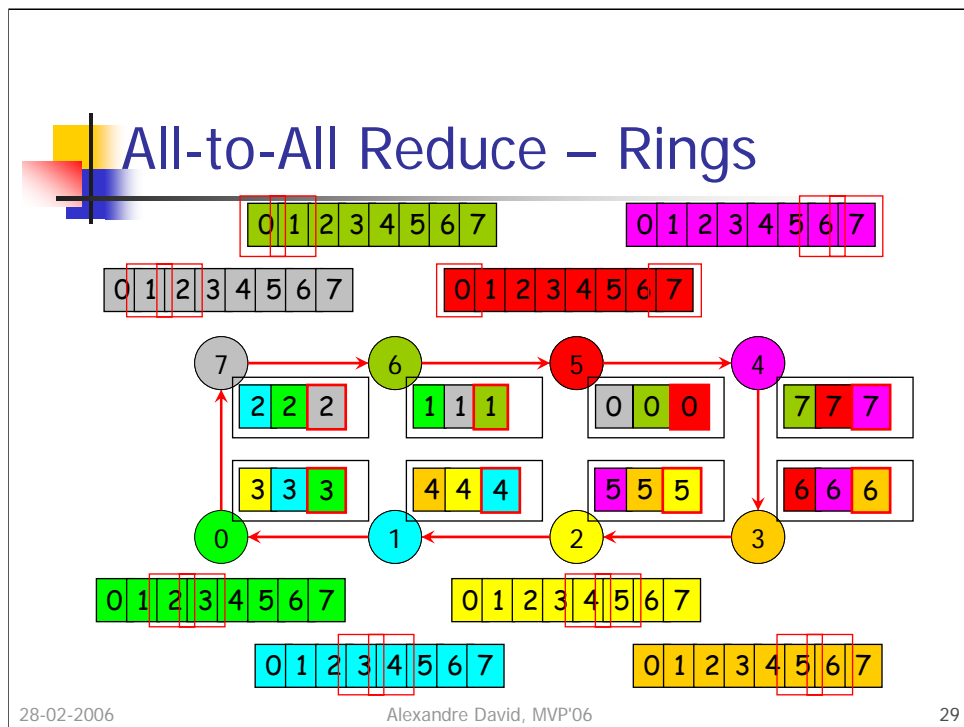
All-to-All Reduce – Rings



28-02-2006

Alexandre David, MVP'06

28



$p-1$ steps.

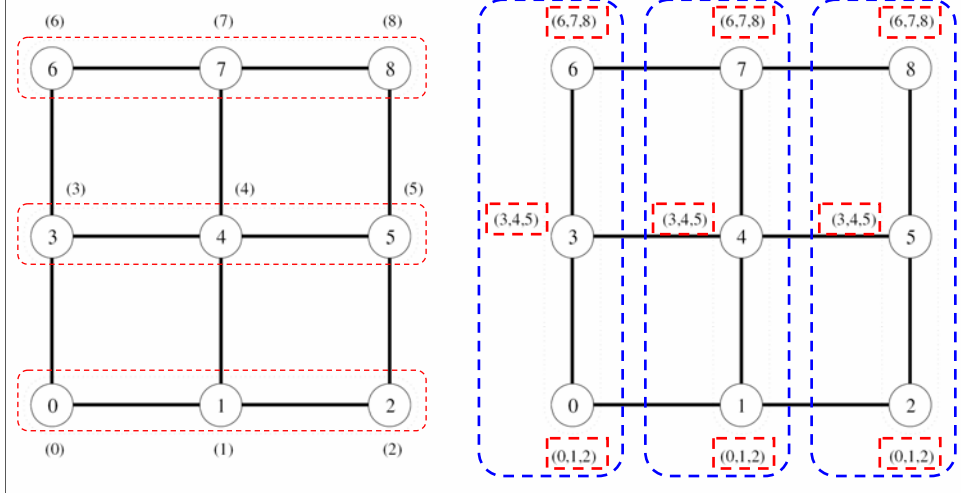


All-to-All Broadcast – Meshes

- Two phases:
 - All-to-all on rows – messages size m .
 - Collect \sqrt{p} messages.
 - All-to-all on columns – messages size $\sqrt{p} * m$.



All-to-All Broadcast – Meshes



Algorithm



```
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
```


All-to-All Broadcast - Hypercubes

- Generalization of the mesh algorithm to $\log p$ dimensions.
- Message size doubles at every step.
- Number of steps: $\log p$.

28-02-2006

Alexandre David, MVP'06

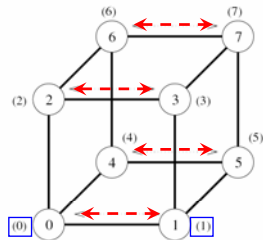
33

Remember the 2 extremes:

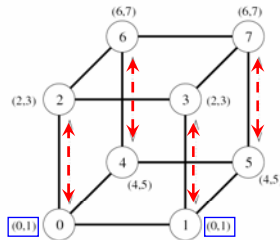
- Linear array: p nodes per (1) dimension – p^1 .
- Hypercubes: 2 nodes per $\log p$ dimensions – $2^{\log p}$.

And in between 2-D mesh \sqrt{p} nodes per (2) dimensions – \sqrt{p}^2 .

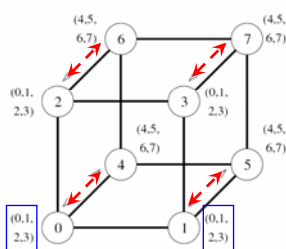
All-to-All Broadcast – Hypercubes



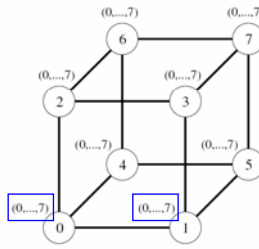
(a) Initial distribution of messages



(b) Distribution before the second step



(c) Distribution before the third step



(d) Final distribution of messages

28-02-2006

34



Algorithm

```
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  $\cup$  msg;
9.   endfor;
10. end ALL_TO_ALL_BC_HCUBE
```

Loop on the dimensions

Exchange messages

Forward (double size)

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

At every step we have a broadcast on sub-cubes. The size of the sub-cubes doubles at every step and all the nodes exchange their messages.

All-to-All Reduction – Hypercubes

```
1. procedure ALL_TO_ALL_RED_HCUBE(my_id, msg, d, result)
2. begin
3.   recloc := 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 := recloc + k;
9.     recloc := recloc + 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[recloc + j] := msg[recloc + j] + temp[j];
14.    endfor;
15.  endfor;
16.  result := msg[my_id];
17. end ALL_TO_ALL_RED_HCUBE
```

Similar pattern
in reverse order.

Combine results

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

Cost Analysis (Time)

- Ring:

- $T = (t_s + t_w m)(p-1).$

- Mesh:

- $T = (t_s + t_w m)(\sqrt{p}-1) + (t_s + t_w m \sqrt{p})(\sqrt{p}-1)$
 $= 2t_s(\sqrt{p}-1) + t_w m(p-1).$

- Hypercube:

- $T = \sum_{i=1}^{\log p} (t_s + 2^{i-1} t_w m)$
 $= t_s \log p + t_w m(p-1).$

$\log p$ steps
message of size $2^{i-1}m$.

28-02-2006

Alexandre David, MVP'06

37

Lower bound for the communication time of all-to-all broadcast for parallel computers on which a node can communicate on only one of its ports at a time $= t_w m(p-1)$. Each node receives at least $m(p-1)$ words of data. That's for **any** architecture.

The straight-forward algorithm for the simple ring architecture is interesting: It is a sequence of p one-to-all broadcasts with different sources every time. The broadcasts are pipelined. That's common in parallel algorithms.

We cannot use the hypercube algorithm on smaller dimension topologies because of congestion.

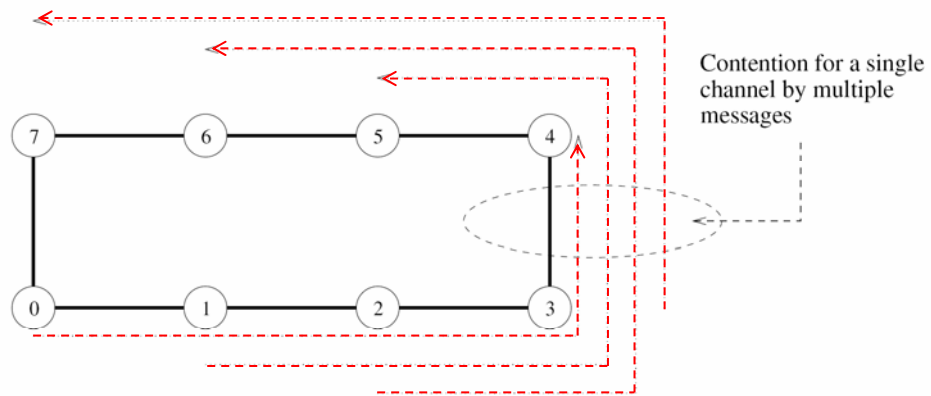


Figure 4.12 Contention for a channel when the communication step of Figure 4.11(c) for the hypercube is mapped onto a ring.

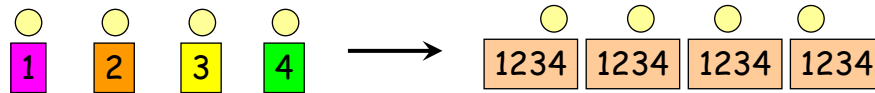
Alexandre David, MVP'06

38

Contention because communication is done on links with single ports. Contention is in the sense of the access to the link. The result is congestion on the traffic.

All-Reduce

- Each node starts with a buffer of size m .
- The final result is the same combination of all buffers on every node.
- Same as all-to-one reduce + one-to-all broadcast.
- Different from all-to-all reduce.




28-02-2006

Alexandre David, MVP'06

39

All-to-all reduce combines p different messages on p different nodes. All-reduce combines 1 message on p different nodes.



All-Reduce Algorithm

- Use all-to-all broadcast but
 - Combine messages instead of concatenating them.
 - The size of the messages does not grow.
 - Cost (in $\log p$ steps): $T = (t_s + t_w m) \log p$.



Prefix-Sum

- Given p numbers n_0, n_1, \dots, n_{p-1} (one on each node), the problem is to compute the sums $s_k = \sum_{i=0}^k n_i$ for all k between 0 and $p-1$.
- Initially, n_k is on the node labeled k , and at the end, the same node holds s_k .

This is a reminder.



Prefix-Sum Algorithm

```
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
```

All-reduce

Prefix-sum

Algorithm 4.9 Prefix sums on a d -dimensional hypercube.

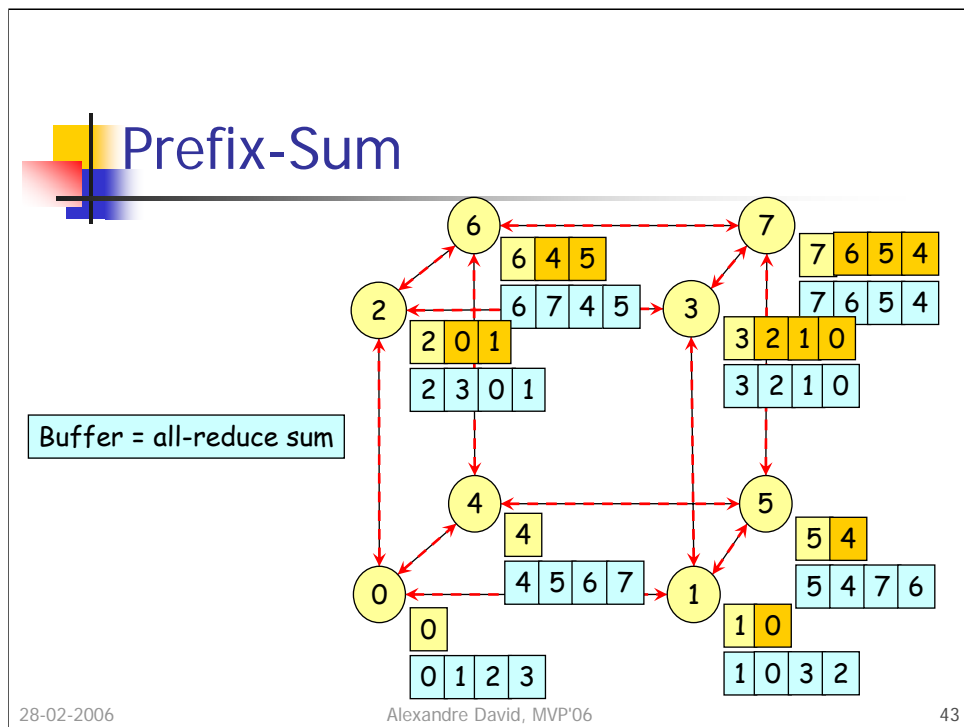


Figure in the book is messed up.