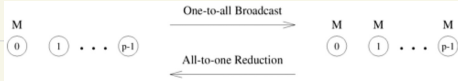


## Network is bidirectional

- send/recv는 따로 계산 안함

## One-All Communication



- 'one' 이 들어있는 single data 한개로 처리함
- broadcast ... MPI-Bcast  
: send 1 message (size M) to (P-1) nodes
- reduction ... MPI-Reduce  
: recv 1 message (size M) from each nodes  
combine with associative ops

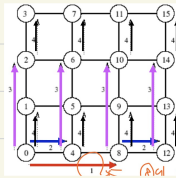
## Cost

$T$  = total time taken by procedure (whole)

$$= (t_s + t_w \cdot m) (\log P) \text{ bits}$$

start transfer msg size

step 증가할수록  
concurrency ↑↑



## Network Topologies

### - naive

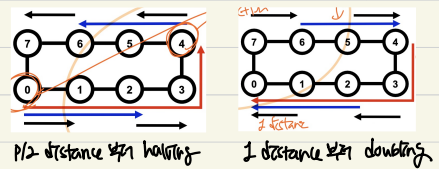
: source node가 다 처리하기 (hotspot/bottleneck)

→ (P-1) stages

### - Ring

→ recursive doubling이 좋음!

→  $(\log_2 P)$  stages



### - Mesh (2D)

: square mesh of P nodes  $\begin{matrix} \sqrt{P} \text{ rows} \\ \sqrt{P} \text{ columns} \end{matrix}$

→ steps

#### ① row에서의 ops

: recursive doubling을 broadcast/reduction

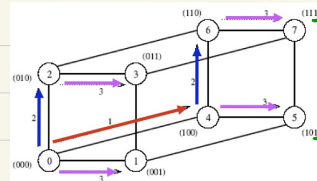
#### ② Column에서의 ops

: recursive doubling을 broadcast/reduction

### - Hypercube (3D)

:  $2^d$  nodes 가 각각 d-dimension으로 연결

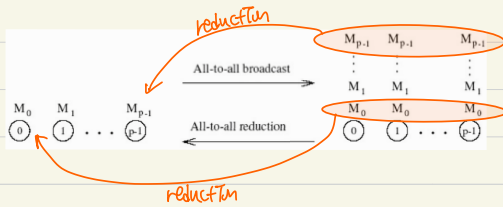
: each dimension에서 2nodes만 연결



→ communication은 MSB에서 LSB order로 진행됨

# All-All Collective Communication

: 모든 node가 one-All communication 하는 것



- Topology 따라 time cost가 다르다!

- step이 많

- msg size 다름

• Broadcast ... MPI-Allgather

• Reduction ... MPI-Reduce-scatter

• network topology

- naive

→ broadcast : 각 node가 1msg를 (P-1) time send

→ reduction : 각 node가 (P-1)개의 msg reduce

\* not recursive doubling 할 필요 x

- ring

→ broadcast : 각 node가 neighbor node로

prev step에 받은 msg를 전달한다.

→ Time :  $(t_s + t_w \cdot m) \cdot (P-1)$

→ reduction : 각 node가 neighbor node로

msg send.

recv node는 reduction하기.

→ Time :  $(t_s + t_w \cdot m) \cdot (P-1)$

- Mesh

→ broadcast

each msg size m

Step ① row-wise ring base broadcast

→ Time :  $(t_s + t_w \cdot m) \cdot (\sqrt{P}-1)$

② column-wise ring base broadcast

each msg size m/√P

→ Time :  $(t_s + t_w \cdot m/\sqrt{P}) \cdot (\sqrt{P}-1)$

→ reduction

: mesh상 모든 node의 방향으로만 한개씩만

- hypercube

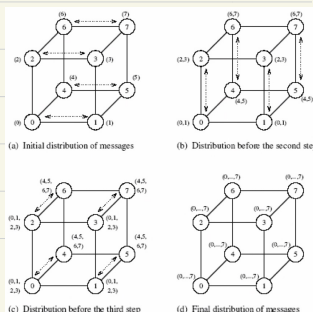
→ 1 node당 3 neighbor가 있는데,

3 steps을 거쳐면서 neighbor랑 모두 communicate 한다

broadcast → step 진행할수록 msg size doubling 되

→ Time

$$= \sum_{i=1}^{\log P} (t_s + 2^{i-1} \cdot m \cdot t_w)$$



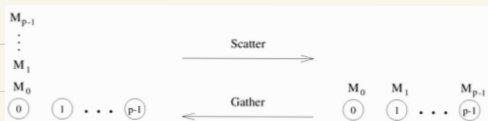
## Scatter/gather

• Scatter ... **MPI\_Scatter**

: 여러 개 MPI msg를 각각 P node에 나눠준다

• Gather ... **MPI\_Gather**

: (P-1) node는 1 msg만 append 해서



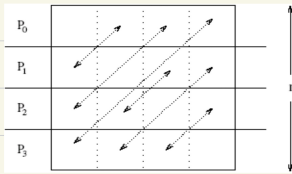
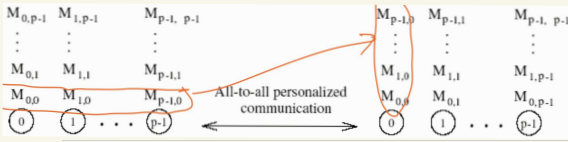
→ 모든 msg는 unique하다.



# All-All Personalized Communication

· matrix transpose 랑 똑같은

All이 Communicate 할여부같은데!



이때는 다들 알고리즘을 짜기

reshuffle

$$- \text{Time} = \sum_{T=1}^{P-1} (t_s + t_w \cdot m \cdot (P-T))$$

# One/All Collective Communication

## Assumption

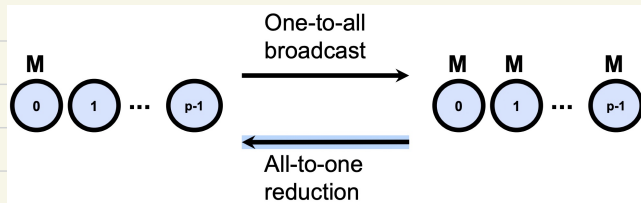
- Network is bidirectional
- Communication is single-pointed
  - ↳ step마다 1 message만 recv 할 수 있음
  - Multi-Message 불가능

## One-to-All

- broadcast
  - : Processor가  $M$  units of data를 갖고있고, 그 data를 다른 processor send

## All-to-one

- reduction
  - : each processor가  $M$  units of data 갖고있고,  
associative operator를 combine함  
target processor에 result 저장됨



# One/All Collective Communication (cont.)

## Ring Protocol

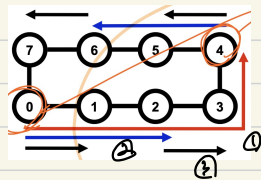
: network가 ring 구조라 가정할 때!

### Broadcast

- naive solution :  $(P-1)$  stages  
2번 1명이 나머지를 다 알려주어야 함!  $\left. \begin{array}{l} \end{array} \right\} 8-1 = 7 \text{ stage 필요}$

- recursive doubling :  $\log_2 P$  stages

send distance를 halving 하는 communication 구조



- initial sender No

$N_0 \rightarrow N_4$  (8/2)

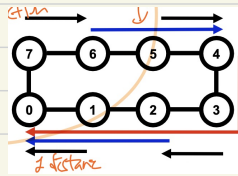
- data 가진 node는 바로 sender 역할을 수행함

$1 \rightarrow 2 \rightarrow 3$

$\left. \begin{array}{l} \end{array} \right\} \log_2 8 = 3 \text{ stage 필요}$

### Reduction

: Broadcast에서 방향 반대로 함



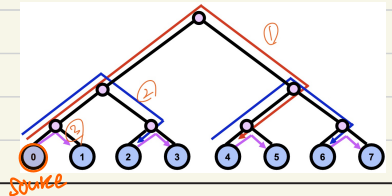
- 이번 반대로 distance 1/2로 줄여나가 수행

## Balanced Binary Tree Protocol

### Broadcast

- processor가 leaf에 있음
- left-most leaf = source of broadcast
- recursive doubling 사용함 (halving)

$\log_2 P$  stage 사용함



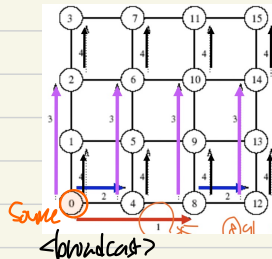
# One/All Collective Communication (cont.)

## 2D Mesh Protocol

- square mesh ( $\sqrt{p} \times \sqrt{p}$ ) of  $p$  nodes
  - row :  $\sqrt{p} \times \sqrt{p}$  linear array of node
  - column :  $\sqrt{p} \times \sqrt{p}$  linear array of node

### Steps

- ① Perform operations along a row
  - ② Perform operations each column concurrently
- the steps are concurrent! (concurrency  $\geq \sqrt{p}$ )



- $T_1$  : No send to  $N_8$  (row / recursive doubling)
- $T_2$  :  $N_0 \rightarrow N_4$  /  $N_8 \rightarrow N_{12}$  degree of concurrency = 2

Every row has data!

- $T_3$  :  $N_0 \rightarrow N_2$  /  $N_4 \rightarrow N_6$  /  $N_8 \rightarrow N_{10}$  /  $N_{12} \rightarrow N_{14}$
- $T_4$  : column 마다 data 있는 이들을 다 각동

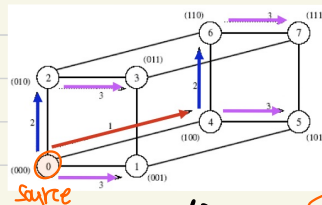
$T_2$ 와  $T_4$ 가 동시에 degree of concurrency  $\uparrow \uparrow$

\* reduction은 거꾸로 진행하는 방식! (degree of concurrency  $\geq \sqrt{p}$ 의 경우)

## Hypercube Protocol

: hypercube with  $2^d$  nodes

- d-dimensional mesh with 2 nodes in each dimension



: d steps 진행하는 ( $d = \log_2 2^d = \log_2 p$ )

$$\rightarrow \text{total cost} = \sum_{i=1}^d (b + 2^{i-1}wm)$$

broadcast 메시지의 size of message

# All / All Collective Communication

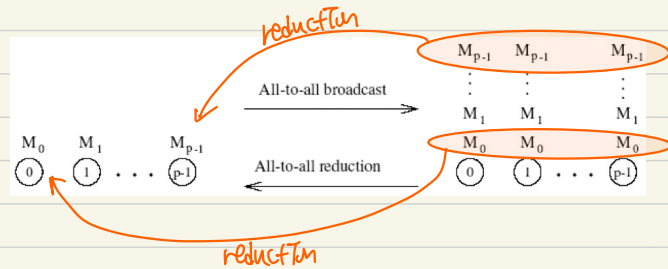
• each processor는 sender와 receiver

All-to-All Broadcast

: 각자가 자신의 data를 모두에게 copy

All-to-All Reduction

: each processor는 copy of result를 가짐 (one-to-many 1대n!)





## All / All Collective Communication (Cont.)

## Ring Protocol

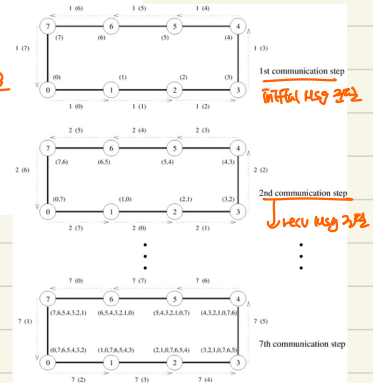
**Broadcast**: on finest step dist right 3 send to 2 left 3 then leave 3 msg to write 1 to 2

```

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  $\oplus$  msg;
11.     endfor;
12. end ALL_TO_ALL_BC_RING

```

시계 방향으로  
 $\oplus$  append



- Cost per mole

$$= t_s + t_w \cdot m$$

## Startup time

single-word transfer time

$T$  # of merges

**Reduction** : Combine each incoming message into local result

- for each step, recu msg  $\frac{2}{2}$  rest 2 send  $\frac{2}{2}$

```
procedure ALL_TO_ALL_RED_RING(my_id, my_msg, p, result)
begin
```

```

left := (my_id - 1) mod p;
right := (my_id + 1) mod p;
recv := 0;
for i := 1 to p - 1 do
    j := (my_id + i) mod p;
    temp := msg[j] + recv;
    send temp to left;
    receive recv from right;
endfor;
result := msg[my_id] + recv;
end ALL to ALL RED RING

```

→ node별 (p-1) communication 발생함

# All / All Collective Communication (Cont.)

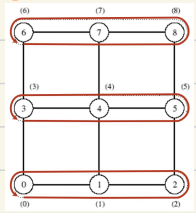
**Mesh**

**broadcast**

**reduction**

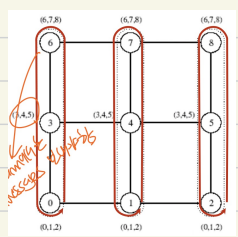
2 phases 3 이차원이다.

① **Row-wise All-to-All Communication**



→ each node collects  $\sqrt{p}$  messages  
 → 이걸 message size는 1번이 m/p 가 된다.  
 → **Cost Per node**  
 $= (t_{cs} + t_{wm})(\sqrt{p}-1)$

② **Column-wise all-to-all communication**



→ each node collects  $\sqrt{p}$  messages  
 → 이걸 message size는 m/p 가 된다.  
 → **Cost Per node**  
 $= (t_{cs} + t_{wm}/p)(\sqrt{p}-1)$

total =  $2t_{cs}(\sqrt{p}-1) + t_{wm}(p-1)$

: **reduction**은 broadcast랑 비슷하지만, reduction은 msg size는 1b

```

my_result = local_value
for each round
    send my_result to partner
    receive msg
    my_result = my_result  $\oplus$  msg
post condition: each my_result now contains global result
    
```

3 steps + send + receive  
 1. if it's 0!  
 associative op.

# Prefix Sum

데이터를 누적하여 합하는 것

## Pre-condition

- $N_k$  indicates that  $k^{\text{th}}$  node

## Problem statement

- Compute sums  $S_k = \sum_{i=0}^k A_i$  for  $\forall k$  between 0 and  $P-1$   
- (ex)  $\langle 3, 1, 4, 0, 2 \rangle$  sequence의 Prefix sum은  $\langle 3, 4, 8, 8, 10 \rangle$

## Post-condition

- $N_k$  node contains  $S_k$

- All-to-All reduction 사용가능함

```
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 2i;
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
```

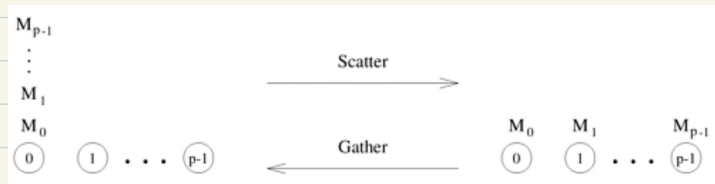
all-to-all로 다 받고 필요한 데이터만 사용하기

# Scatter · Gather

## Scatter

: node sends unique msg of size  $m$  to every other node  
(one-to-all personalized)

- broadcast와 다름: msg 크기가 정확히 같음.  
msg size constant recursive halving 이용 가능



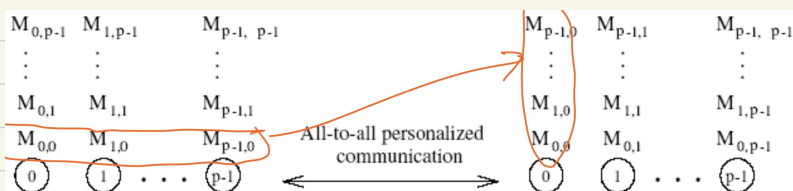
## Gather

: single node collects unique msg from each node

- reduction은 arithmetic 일지  
gather는 가 appending 'b

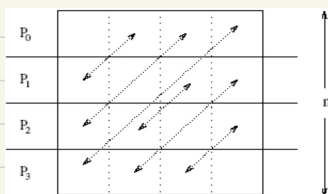
# All-to-All Personalized Communication

Total exchange



: matrix transpose 똑같은 것

matrix transpose



- each node holds dense row of matrix
- do transpose

• ring protocol

- each node<sup>s</sup> send all pieces of data of slice  $m(p-1)$  → neighbors만 보낸다
- recu node는 2중씩 needed 1 msg만 extract 하긴  $m(p-2)$ 를 보낸다.

:

(p-1) step 반복!

- cost

$$T = \sum_{i=1}^{p-1} (b + tw \cdot m(p-i))$$

step별 msg slice 감소