1

# Parallel and Distributed Computing
## CS3006

Lecture 13

**Basic Communication Operations-III**

25th April 2022

Dr. Rana Asif Rehman

# All-Reduce

# Basic Communication Operations
## (All-Reduce)

- Precondition: Every process *i* has a single message $M_i$ of size *m words.*
- Post condition: All processes have a reduced message *M of size m words.*

**Strategies:**

1. Use **all-to-one reduction** followed by **one-to-all** broadcast $(2 * (t_s + mt_w) \log p)$
2. Use **modified All-to-All comm.** algorithm for hypercube $((t_s + mt_w) \log p)$
   - Replace Union with associative operator

4

# Prefix-Sum

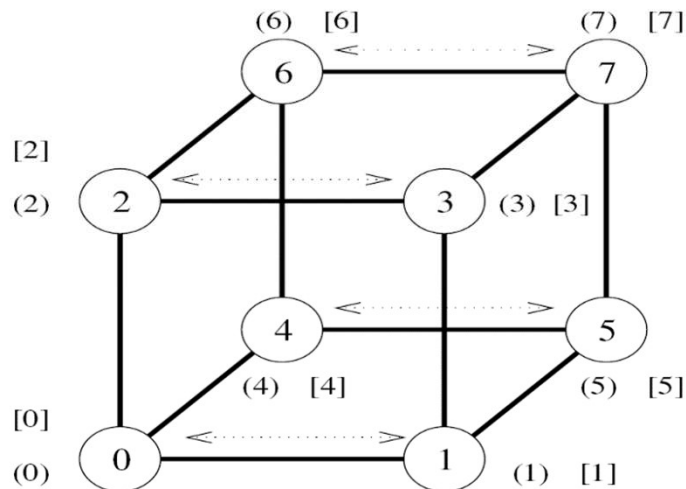# Basic Communication Operations
## (Prefix-Sums)

➡ Prefix-sums are also known as scan operations

➡ Given **p** numbers **n₀, n₁, ..., nₚ₋₁** (one on each node), the problem is to compute the sums such that: -

  ➡ $S_k = \sum_{i=0}^{K}(n_i)$

   ➡ Here $S_k$ is the prefix-sum computed at kth node after the operation.
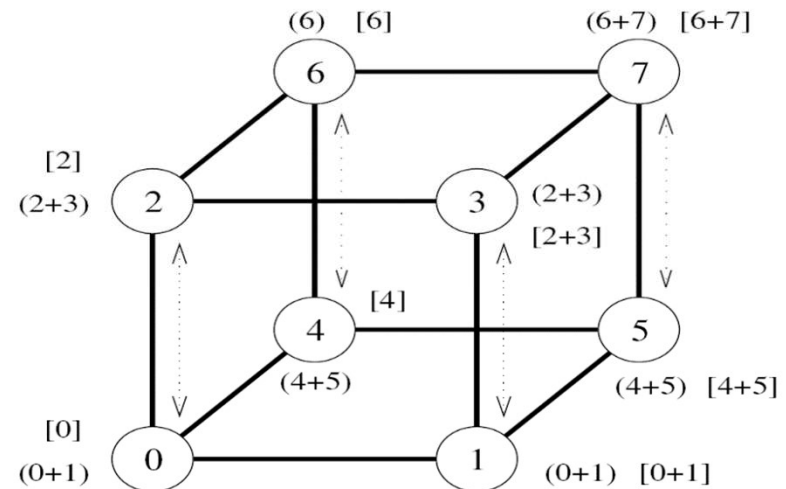
**Example:**

➡ Original sequence: <3, 1, 4, 0, 2>

➡ Sequence of prefix sums: <3, 4, 8, 8, 10>

# Basic Communication Operations
## (Prefix-Sums)

(6)  [6]                    (7)  [7]
     6 ──────────────── 7

[2]
(2)   2 ──────── 3  (3)  [3]

          4 ──────── 5

(4)  [4]                    (5)  [5]

[0]
(0)   0 ──────── 1  (1)  [1]

(a)  Initial distribution of values

(6)  [6]                    (6+7)  [6+7]
     6 ──────────────── 7

[2]
(2+3)   2 ──────── 3  (2+3)
                          [2+3]

               [4]
          4 ──────── 5

(4+5)                    (4+5)  [4+5]

[0]
(0+1)   0 ──────── 1  (0+1)  [0+1]

(b)  Distribution of sums before second step

(4+5+6)  [4+5+6]        (4+5+6+7)  [4+5+6+7]
        6 ──────────────── 7

[0+1+2]
(0+1+   2 ──────── 3
2+3)        [0+1+2+3]
                    (0+1+2+3)

               [4]
          4 ──────── 5  (4+5)
(4+5)                    [4+5]

[0]
(0+1+   0 ──────── 1  (0+1+
2+3)                    2+3)  [0+1]

(c)  Distribution of sums before third step

[0+ .. +6]              [0+ .. +7]
     6 ──────────────── 7

[0+1+2]
        2 ──────── 3
                    [0+1+2+3]

          4 ──────── 5
                          [0+ .. +5]
[0+1+2+3+4]
[0]
        0 ──────── 1  [0+1]

(d)  Final distribution of prefix sums

**Figure 4.12** Computing prefix sums on an eight-node hypercube. At each node, square brackets

# Basic Communication Operations
## (Prefix-Sums)

```
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 and Gather

8

# Basic Communication Operations
## (Scatter and Gather)

- Gather is different than reduction as it doesn't reduce the results with associative operator
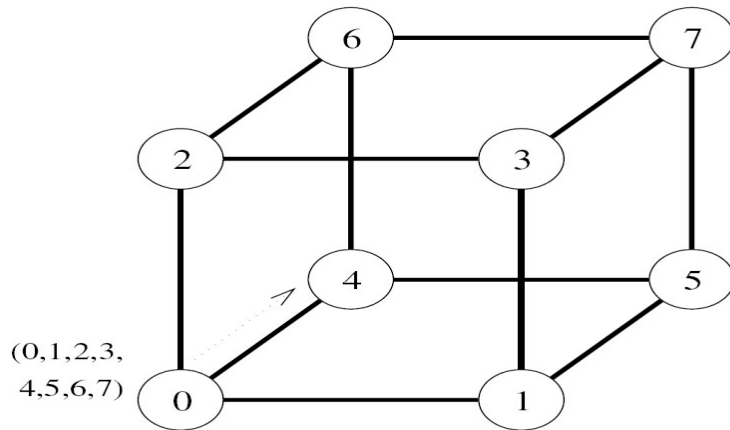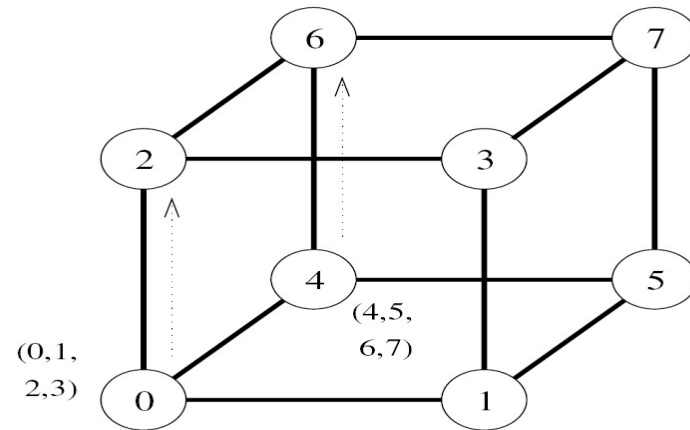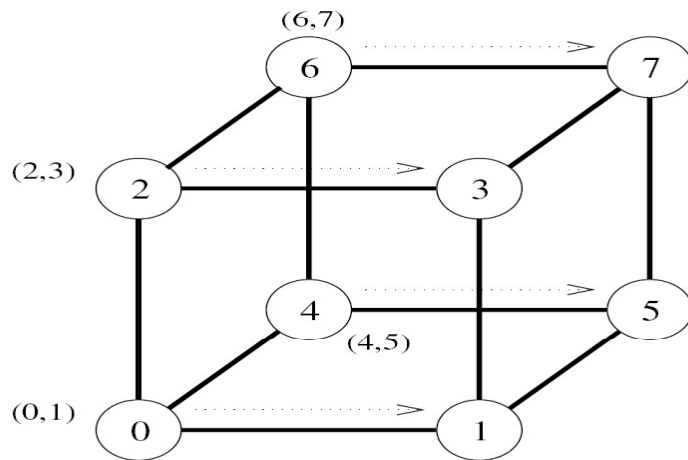


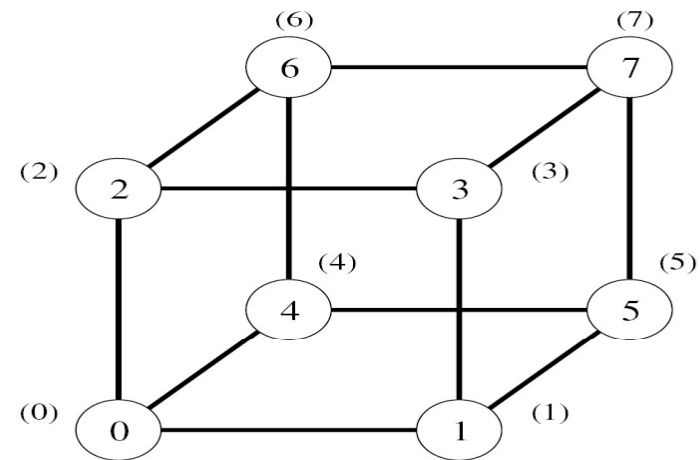**Figure 4.14** Scatter and gather operations.

(a) Initial distribution of messages

(b) Distribution before the second step

(c) Distribution before the third step

(d) Final distribution of messages

**Figure 4.15** The scatter operation on an eight-node hypercube.

# All-to-All personalized Communication

11

# Basic Communication Operations
## (All-to-All personalized)

- Each node sends a distinct message of size **m** to every other node.
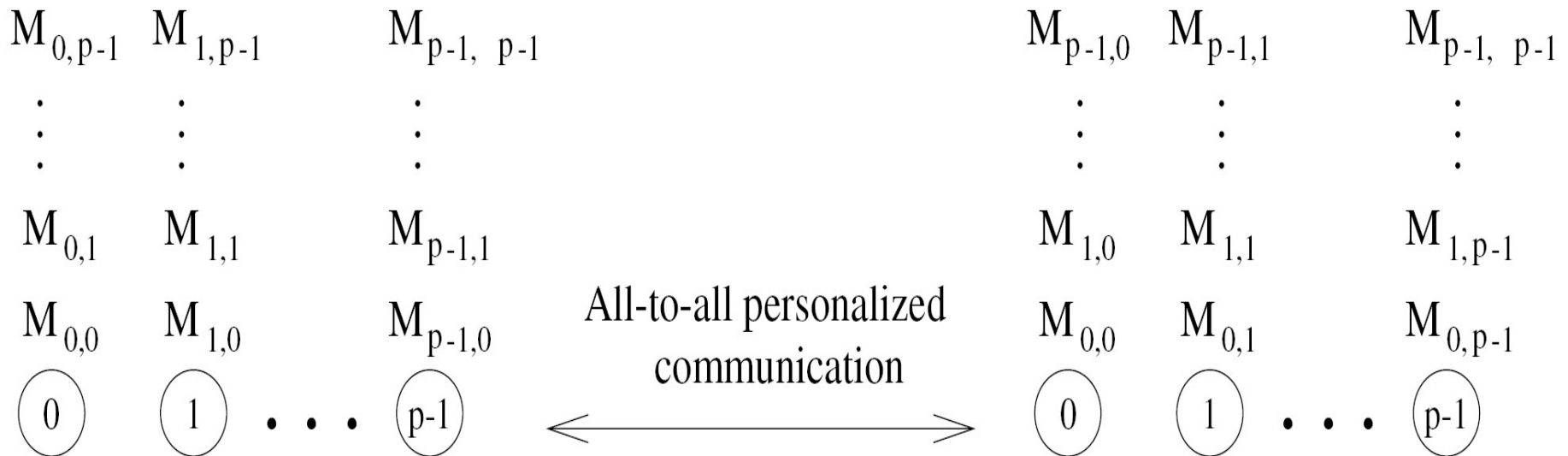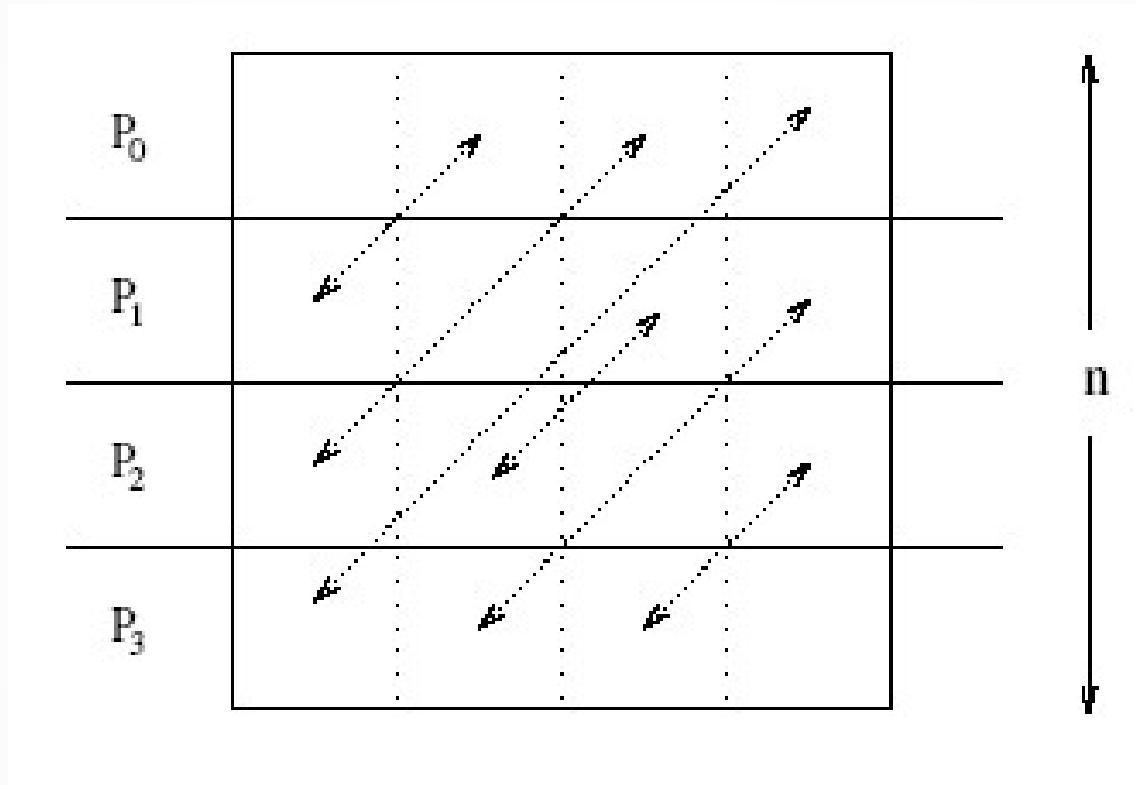- Also known **total exchange**

$$M_{0,p-1} \quad M_{1,p-1} \qquad M_{p-1,\ p-1} \qquad\qquad M_{p-1,0} \quad M_{p-1,1} \qquad M_{p-1,\ p-1}$$

$$\vdots \qquad\quad \vdots \qquad\qquad \vdots \qquad\qquad\qquad\quad \vdots \qquad\quad \vdots \qquad\qquad \vdots$$

$$M_{0,1} \qquad M_{1,1} \qquad M_{p-1,1} \qquad\qquad M_{1,0} \qquad M_{1,1} \qquad M_{1,p-1}$$

$$M_{0,0} \qquad M_{1,0} \qquad M_{p-1,0} \qquad\qquad M_{0,0} \qquad M_{0,1} \qquad M_{0,p-1}$$

All-to-all personalized communication

$$(0) \quad (1) \quad \cdots \quad (p-1) \qquad\qquad (0) \quad (1) \quad \cdots \quad (p-1)$$

**Figure 4.16**   All-to-all personalized communication.

# Basic Communication Operations
## (All-to-All personalized)



All-to-all personalized communication in transposing a *4 x 4 matrix using four processes.*

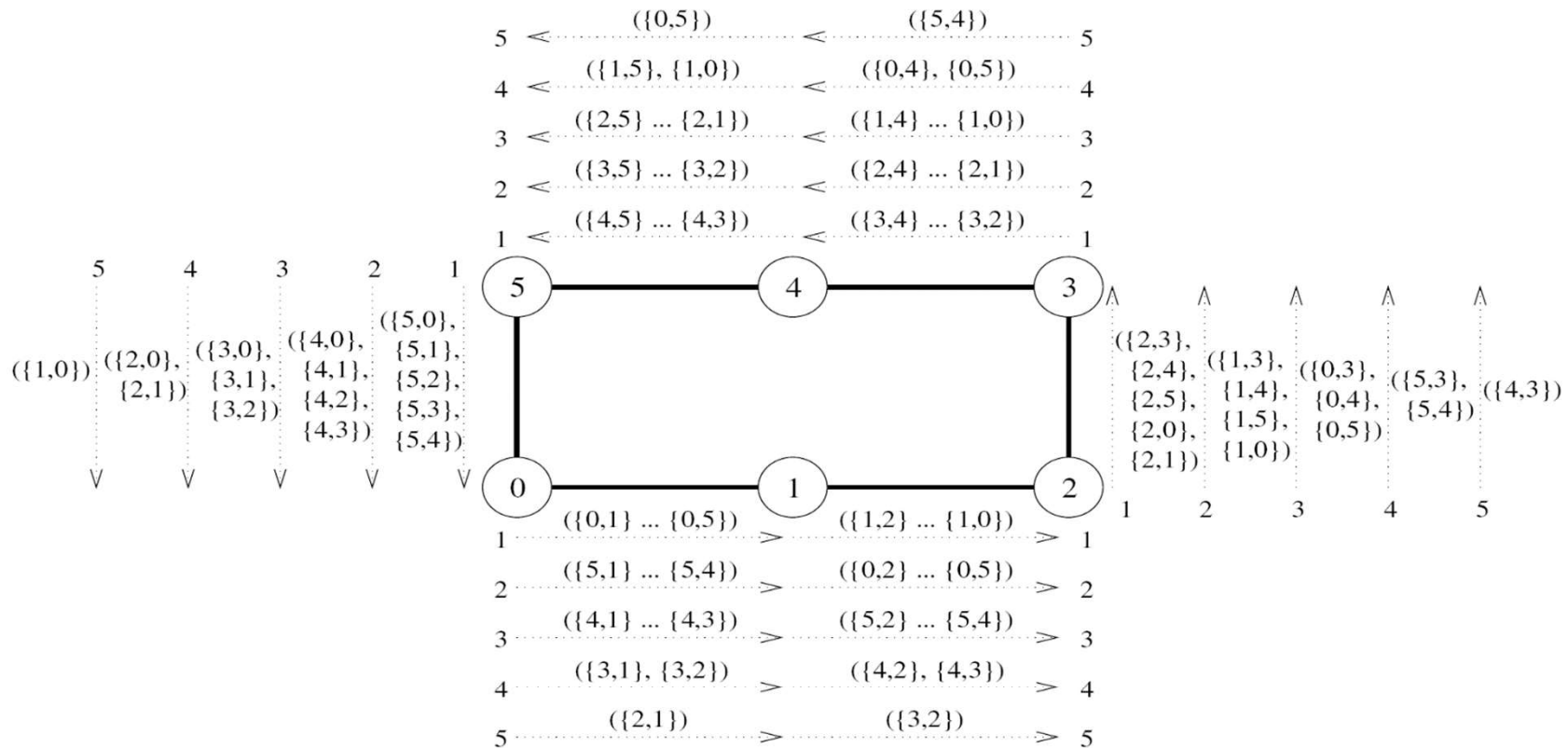## (All-to-All personalized [Ring])



**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 $(\{x_1, y_1\}, \{x_2, y_2\}, \ldots, \{x_n, y_n\})$ indicates a message that is formed by concatenating $n$ individual messages.
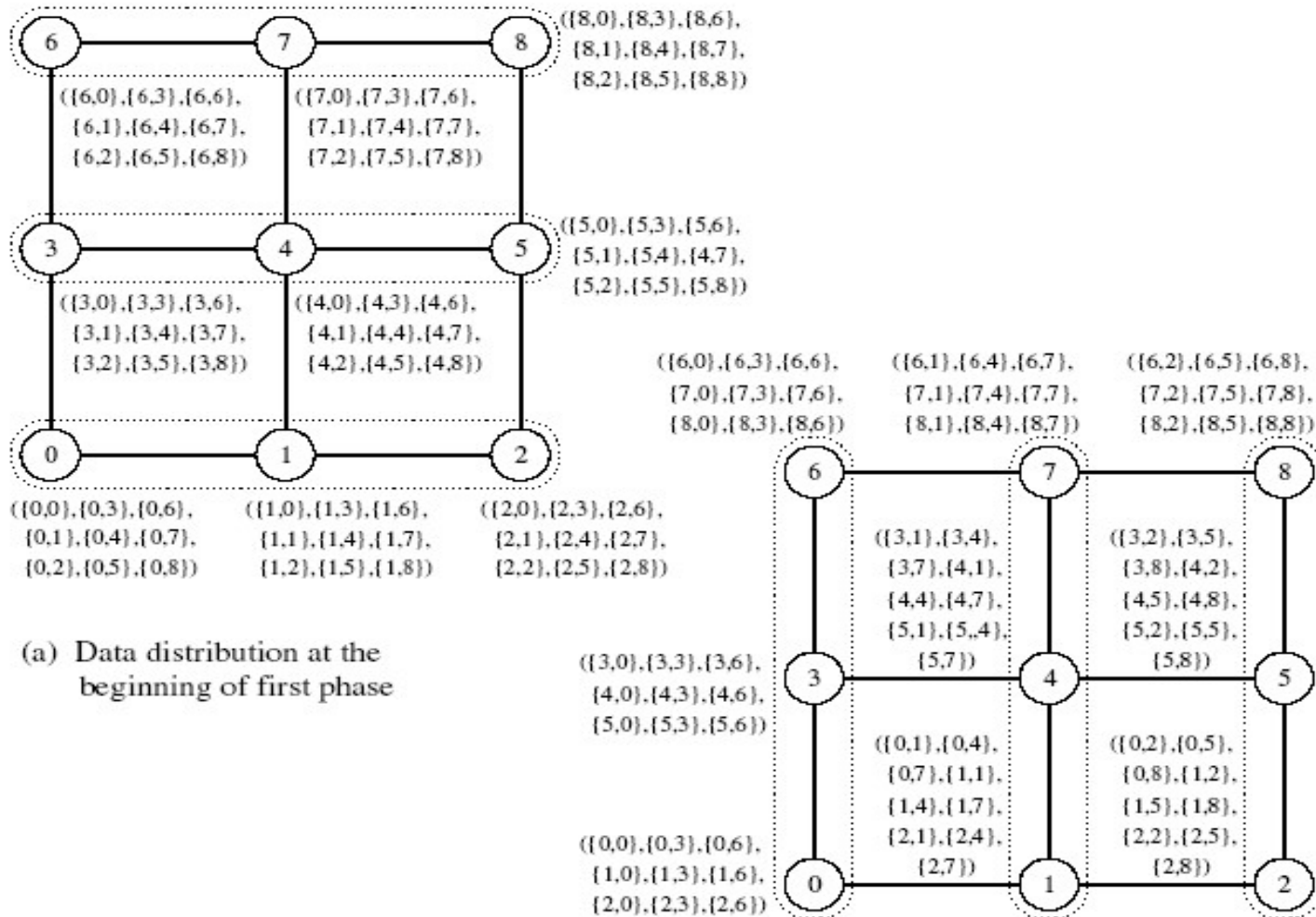
## (All-to-All personalized [Ring])

**Cost Analysis**

- $T = \sum_{i=1}^{(p-1)}(t_s + (p-i)mt_w)$
  - $= \sum_{i=1}^{(p-1)}(t_s) + mt_w\sum_{i=1}^{(p-1)}(p-i))$
    - $\rightarrow (\mathbf{p}-\mathbf{1})(t_s) + mt_w\sum_{i=1}^{(p-1)}(\boldsymbol{i})$
    - $\rightarrow \left((t_s + \left(\frac{1}{2}\right)\boldsymbol{pm}t_w\right)(p-1)$

## (All-to-All personalized [Mesh])



(a) Data distribution at the beginning of first phase

(b) Data distribution at the beginning of second phase

# Basic Communication Operations
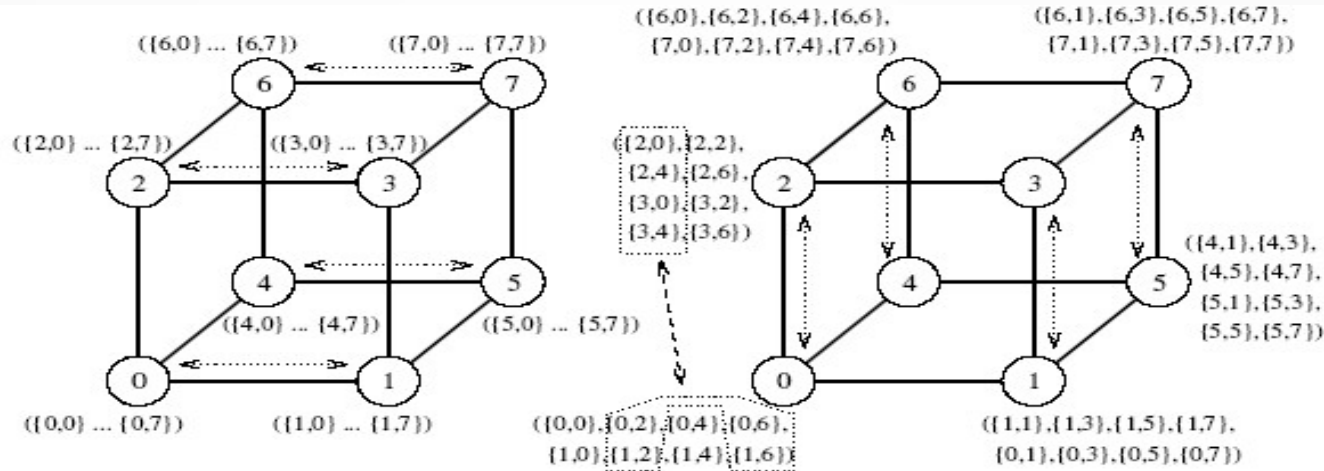## (All-to-All personalized [Mesh])

**Cost Analysis**

- ➧ Time for the first phase is identical to that in a ring with $\sqrt{p}$ processors, i.e., $(t_s + t_w mp/2)(\sqrt{p} - 1)$.
  - ➧ Here **$mt_w$** *becomes* $\sqrt{p}\,mt_w$ and **P** *becomes* $\sqrt{p}$
- ➧ Time in the second phase is identical to the first phase. Therefore, total time is twice of this time, i.e.,
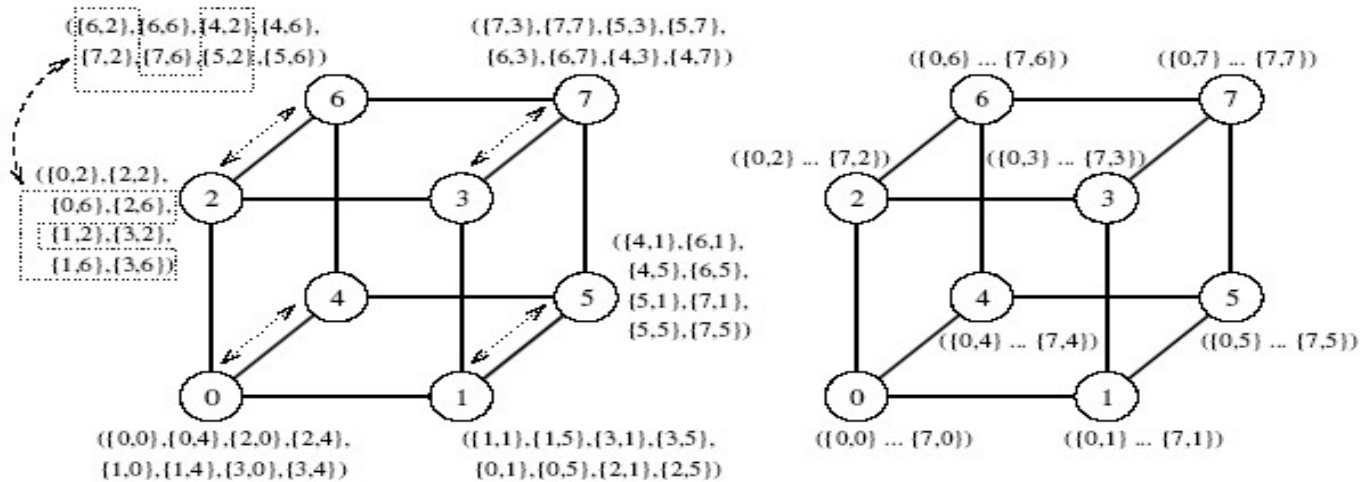
# Basic Communication Operations
## (All-to-All personalized [Hyper Cube])



(a) Initial distribution of messages

(b) Distribution before the second step

(c) Distribution before the third step

(d) Final distribution of messages

# Questions

# References

1.  Kumar, V., Grama, A., Gupta, A., & Karypis, G. (2017). *Introduction to parallel computing*. Redwood City, CA: Benjamin/Cummings.