



Sorting (Chapter 9)

Alexandre David

B2-206



Sorting

Problem

Arrange an unordered collection of elements into monotonically increasing (or decreasing) order.

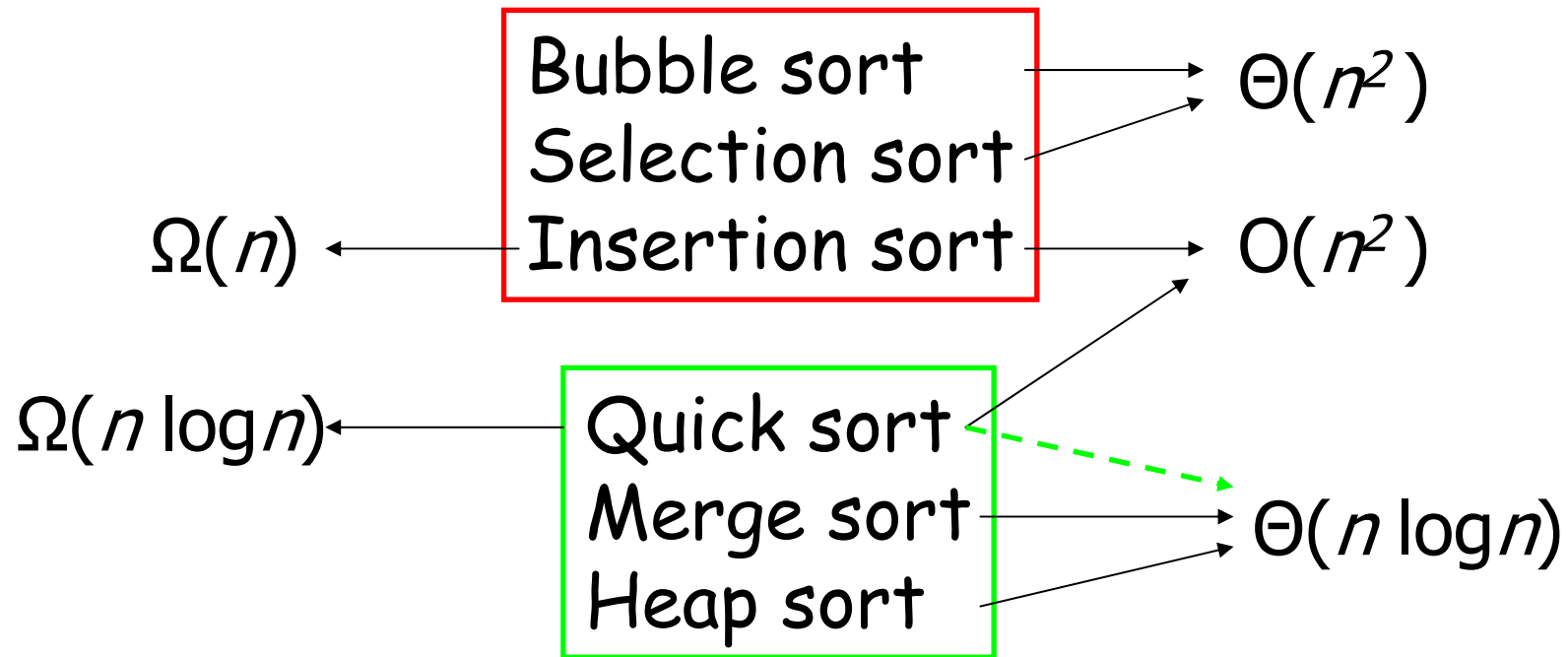
Let $S = \langle a_1, a_2, \dots, a_n \rangle$.

Sort S into $S' = \langle a'_1, a'_2, \dots, a'_n \rangle$ such that

$a'_i \leq a'_j$ for $1 \leq i \leq j \leq n$

and S' is a permutation of S .

Recall on Comparison Based Sorting Algorithms



Characteristics of Sorting

Algorithms



- **In-place** sorting: No need for additional memory (or only constant size).
- **Stable** sorting: Ordered elements keep their original relative position.
- **Internal** sorting: Elements fit in process memory.
- **External** sorting: Elements are on auxiliary storage.



Fundamental Distinction

- **Comparison based** sorting:
 - *Compare-exchange* of pairs of elements.
 - Lower bound is **$\Omega(n \log n)$** (proof based on decision trees).
 - Merge & heap-sort are optimal.
- **Non-comparison based** sorting:
 - Use information on the element to sort.
 - Lower bound is **$\Omega(n)$** .
 - Counting & radix-sort are optimal.



Issues in Parallel Sorting

- Where to store input & output?
 - One process or distributed?
 - Enumeration of processes used to distribute output.
- How to compare?
 - How many elements per process?
 - As many processes as element \Rightarrow poor performance because of inter-process communication.

Parallel Compare-Exchange

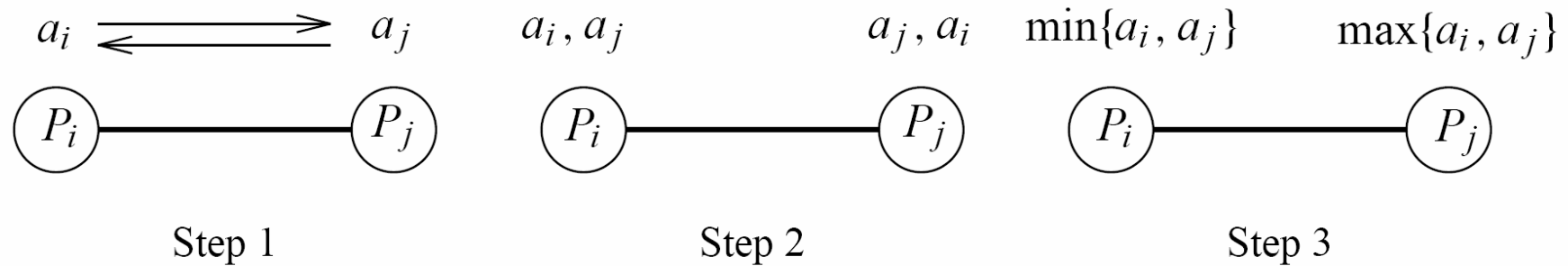
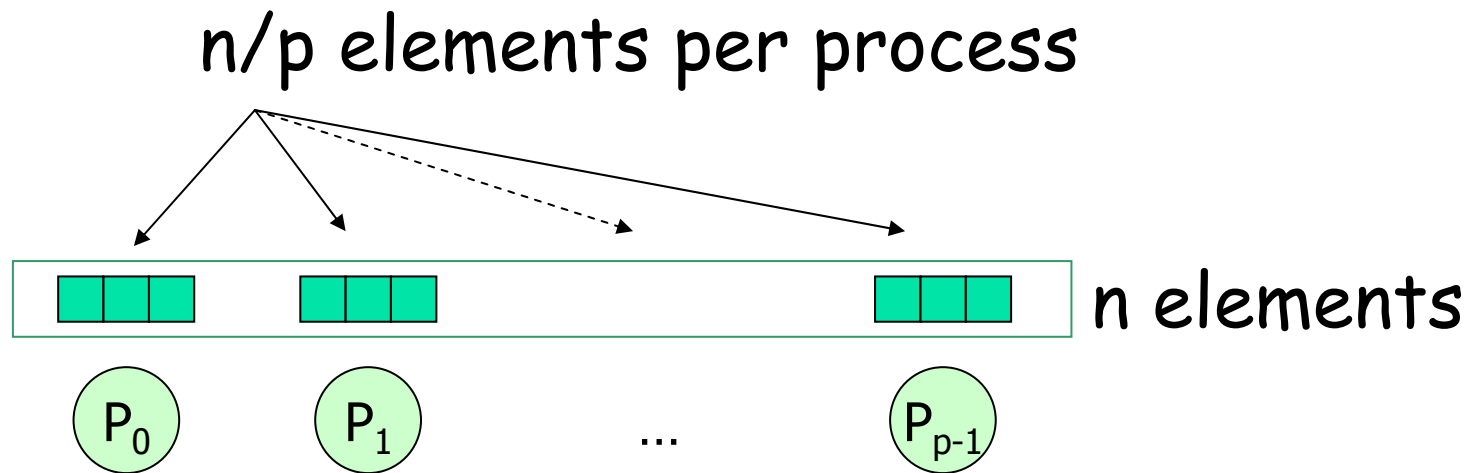


Figure 9.1 A parallel compare-exchange operation. Processes P_i and P_j send their elements to each other. Process P_i keeps $\min\{a_i, a_j\}$, and P_j keeps $\max\{a_i, a_j\}$.

Communication cost: $t_s + t_w$.

Comparison cost much cheaper \Rightarrow communication time dominates.

Blocks of Elements Per Process

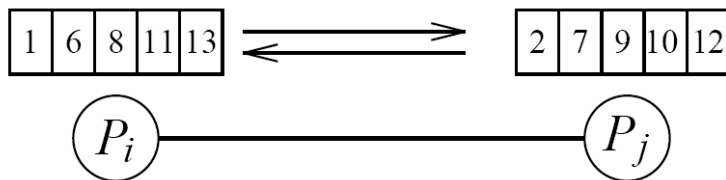


Blocks: $A_0 \leq A_1 \leq \dots \leq A_{p-1}$

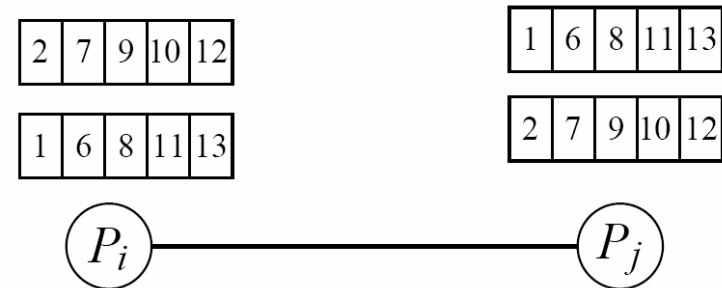
Compare-Split

For large blocks: $\Theta(n/p)$

Exchange: $\Theta(t_s + t_w n/p)$

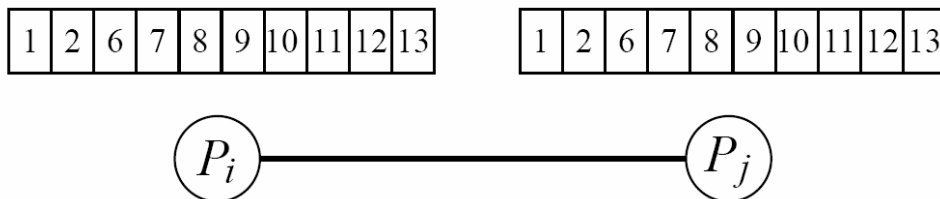


Step 1



Step 2

Merge: $\Theta(n/p)$



Step 3

Split: $O(n/p)$



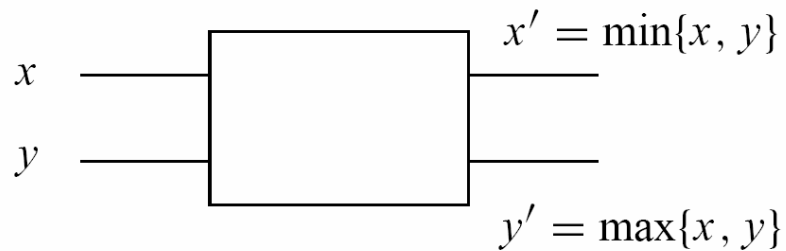
Step 4



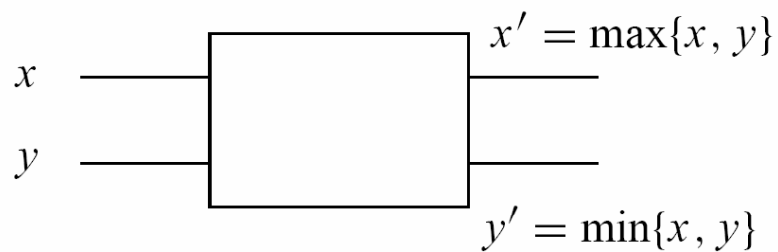
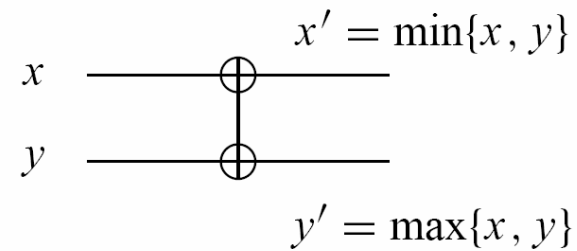
Sorting Networks

- Mostly of theoretical interest.
- Key idea: Perform many comparisons in parallel.
- Key elements:
 - Comparators: 2 inputs, 2 outputs.
 - Network architecture: Comparators arranged in columns, each performing a permutation.
 - Speed proportional to the depth.

Comparators



(a)



(b)

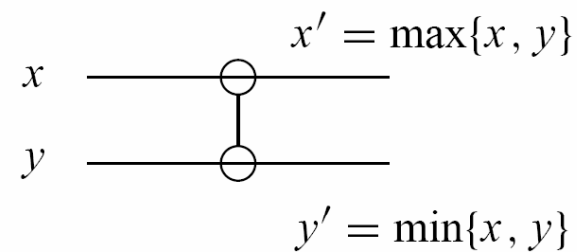


Figure 9.3 A schematic representation of comparators: (a) an increasing comparator, and (b) a decreasing comparator.

Sorting Networks

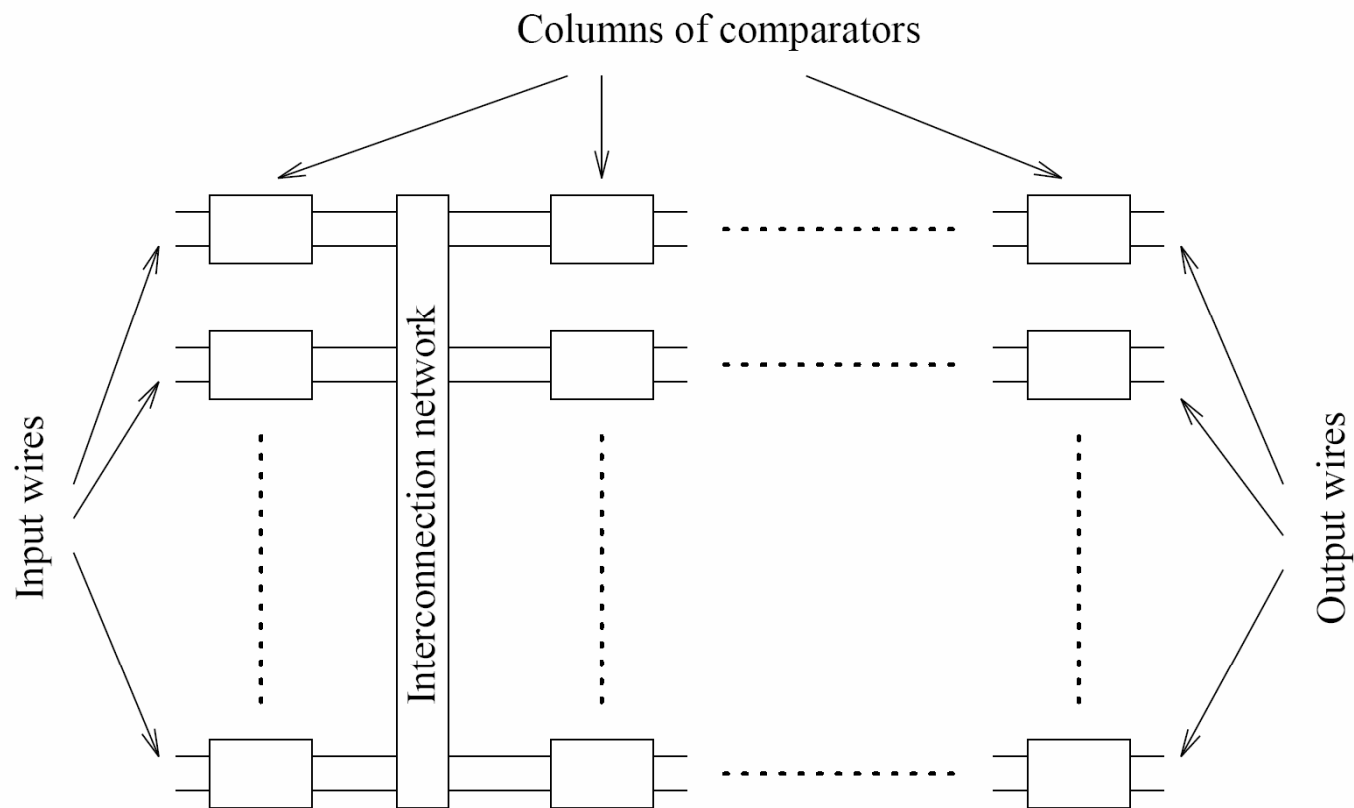


Figure 9.4 A typical sorting network. Every sorting network is made up of a series of columns, and each column contains a number of comparators connected in parallel.



Bitonic Sequence

Definition

A bitonic sequence is a sequence of elements $\langle a_0, a_1, \dots, a_n \rangle$ s.t.

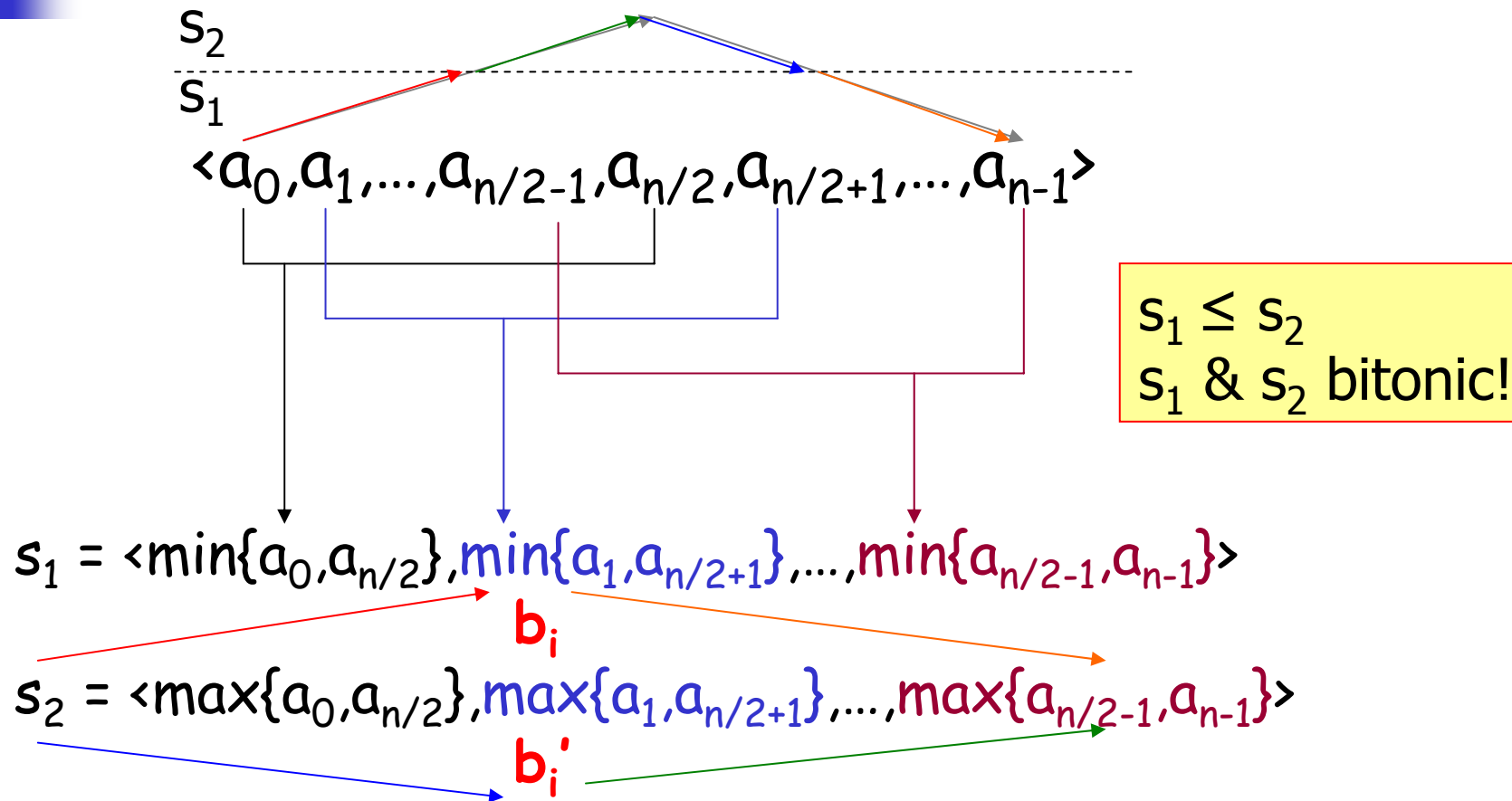
1. $\exists i, 0 \leq i \leq n-1$ s.t. $\langle a_0, \dots, a_i \rangle$ is monotonically increasing and $\langle a_{i+1}, \dots, a_{n-1} \rangle$ is monotonically decreasing,
2. or there is a cyclic shift of indices so that 1) is satisfied.



Bitonic Sort

- Rearrange a bitonic sequence to be sorted.
- Divide & conquer type of algorithm (similar to quicksort) using **bitonic splits**.
 - Sorting a bitonic sequence using bitonic splits = bitonic merge.
 - But we need a bitonic sequence...

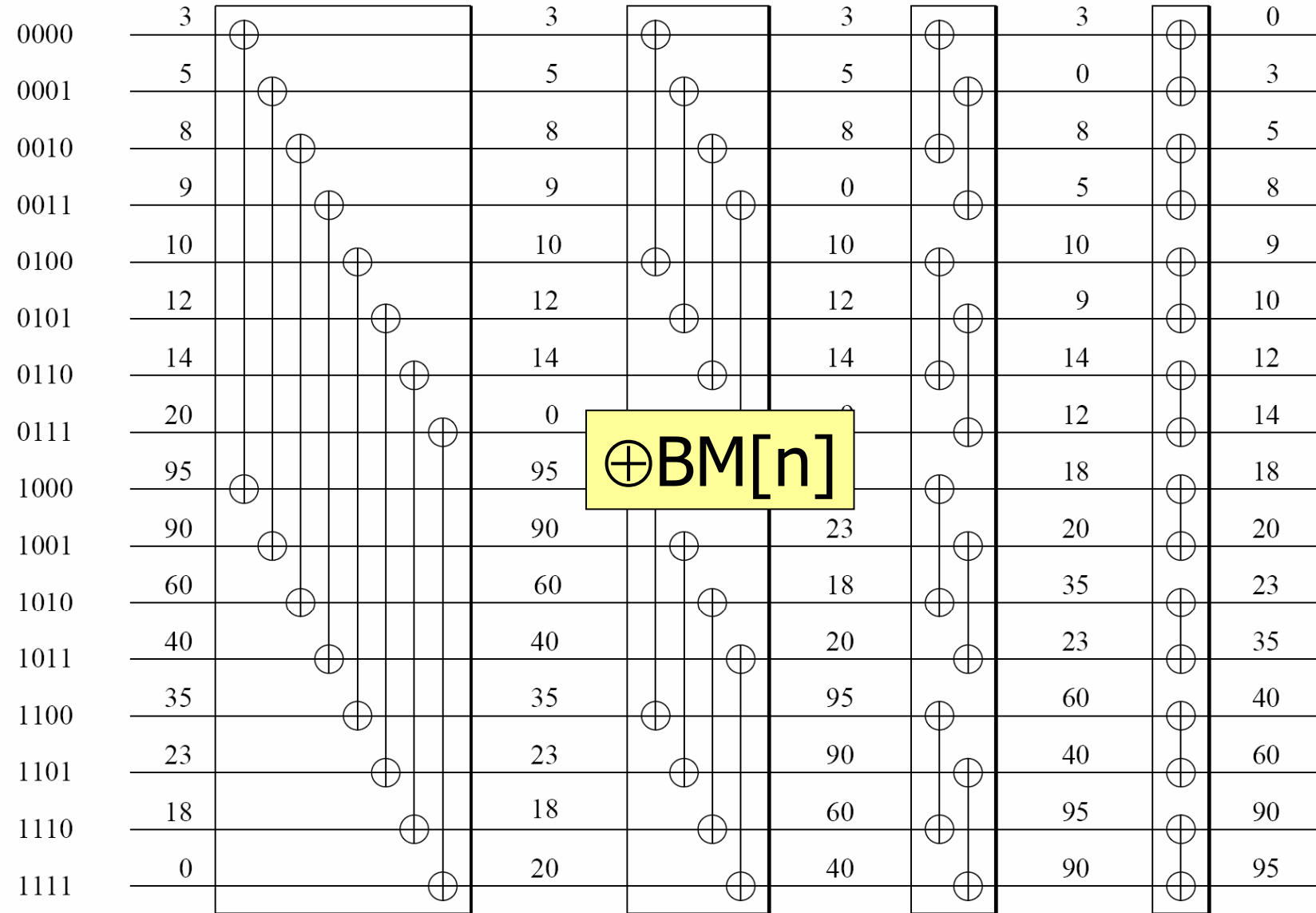
Bitonic Split



Bitonic Merging Network

$\log n$ stages

Wires



$n/2$ comparators



Bitonic Sort

- Use the bitonic network to merge bitonic sequences of increasing length... starting from 2, etc.
- Bitonic network is a component.

Bitonic Sort

$\log n$ stages

Wires

0000
0001
0010
0011
0100
0101

\oplus BM[2]

\ominus BM[2]

\oplus BM[2]

\oplus BM[4]

\oplus BM[8]

BM[16]

1001
1010
1011
1100
1101
1110
1111

\oplus BM[2]

\ominus BM[2]

\oplus BM[2]

\ominus BM[2]

\oplus BM[4]

\ominus BM[4]

\ominus BM[8]

Cost: $O(\log^2 n)$.
Simulated on a serial computer: $O(n \log^2 n)$.

Mapping to Hypercubes & Mesh

– Idea

- Communication intensive, so special care for the mapping.
- How are the input wires paired?
 - Pairs have their labels differing by only one bit
⇒ mapping to hypercube straightforward.
 - For a mesh lower connectivity, several solutions but worse than the hypercube because the sequential algorithm is suboptimal.
 $T_p = O(\log n) + O(1/n)$ for 1 element/process.
 - Block of elements: sort locally $(n/p \log n/p)$ & use bitonic merge ⇒ cost optimal.



Bubble Sort

```
procedure BUBBLE_SORT(n)
begin
  for i := n-1 downto 1 do
    for j := 1 to i do
      compare_exchange(aj, aj+1);
    end
  end
```

$\Theta(n^2)$

- Difficult to parallelize as it is because it is inherently sequential.

Odd-Even Transposition Sort

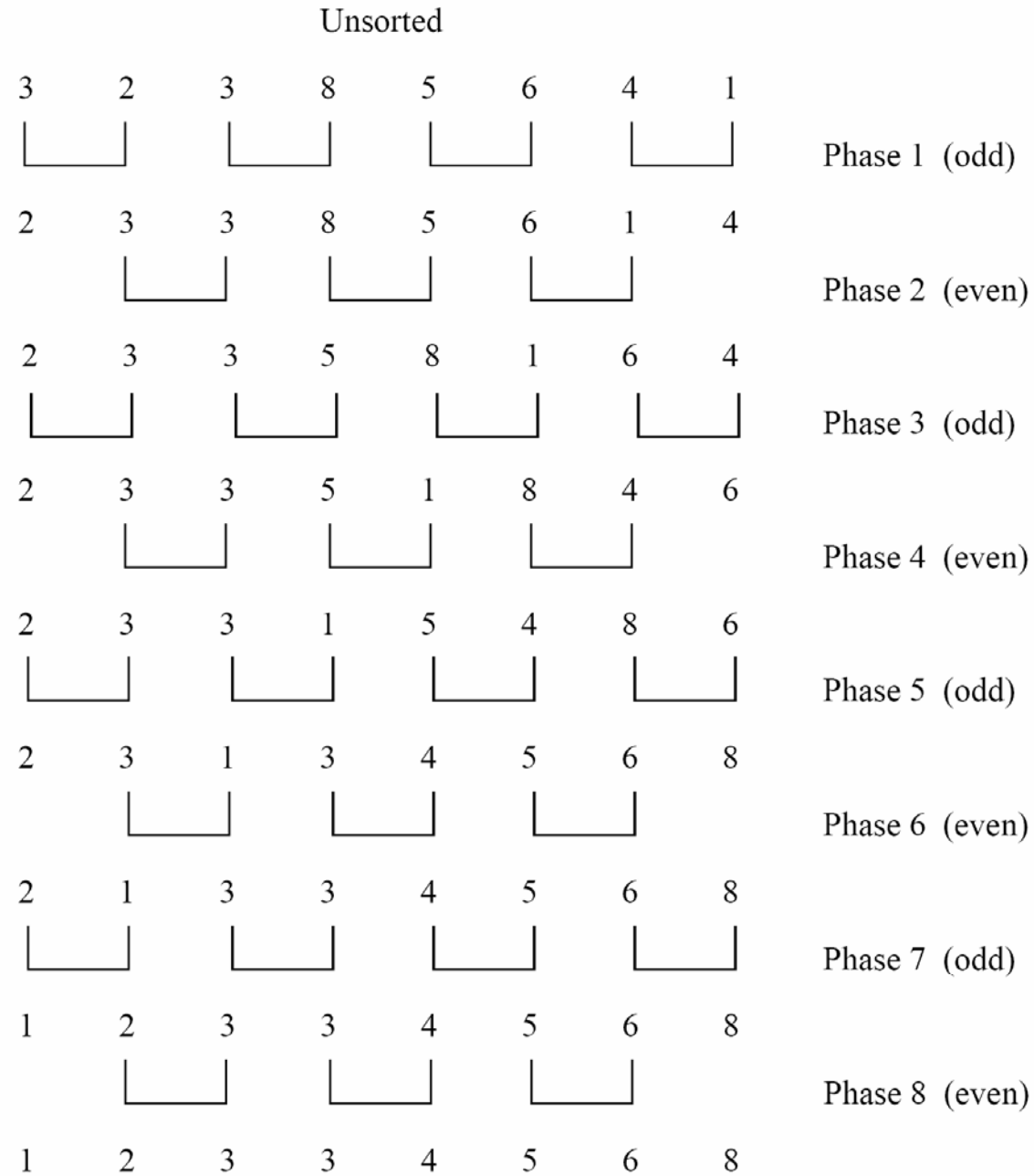
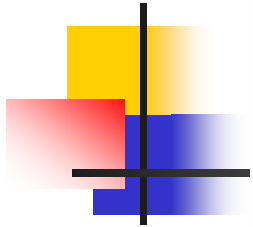
```
1.  procedure ODD-EVEN(n)
2.  begin
3.    for i := 1 to n do
4.    begin
5.      if i is odd then
6.        for j := 0 to n/2 - 1 do
7.          compare-exchange(a2j+1, a2j+2);
8.      if i is even then
9.        for j := 1 to n/2 - 1 do
10.         compare-exchange(a2j, a2j+1);
11.    end for
12. end ODD-EVEN
```

$\Theta(n^2)$

(*a*₁,*a*₂),(*a*₃,*a*₄)...

(*a*₂,*a*₃),(*a*₄,*a*₅)...

Algorithm 9.3 Sequential odd-even transposition sort algorithm.



Odd-Even Transposition Sort

- Easy to parallelize!
 - $\Theta(n)$ if 1 process/element.
 - Not cost optimal but use fewer processes, an optimal local sort, and compare-splits:

$$T_P = \Theta\left(\frac{n}{p} \log \frac{n}{p}\right) + \Theta(n) + \Theta(n)$$

local sort

Cost optimal for $p = O(\log n)$
but not scalable (few processes).

tion



Improvement: Shellsort

- 2 phases:
 - Move elements on longer distances.
 - Odd-even transposition but stop when no change.
- Idea: Put quickly elements near their final position to reduce the number of iterations of odd-even transposition.

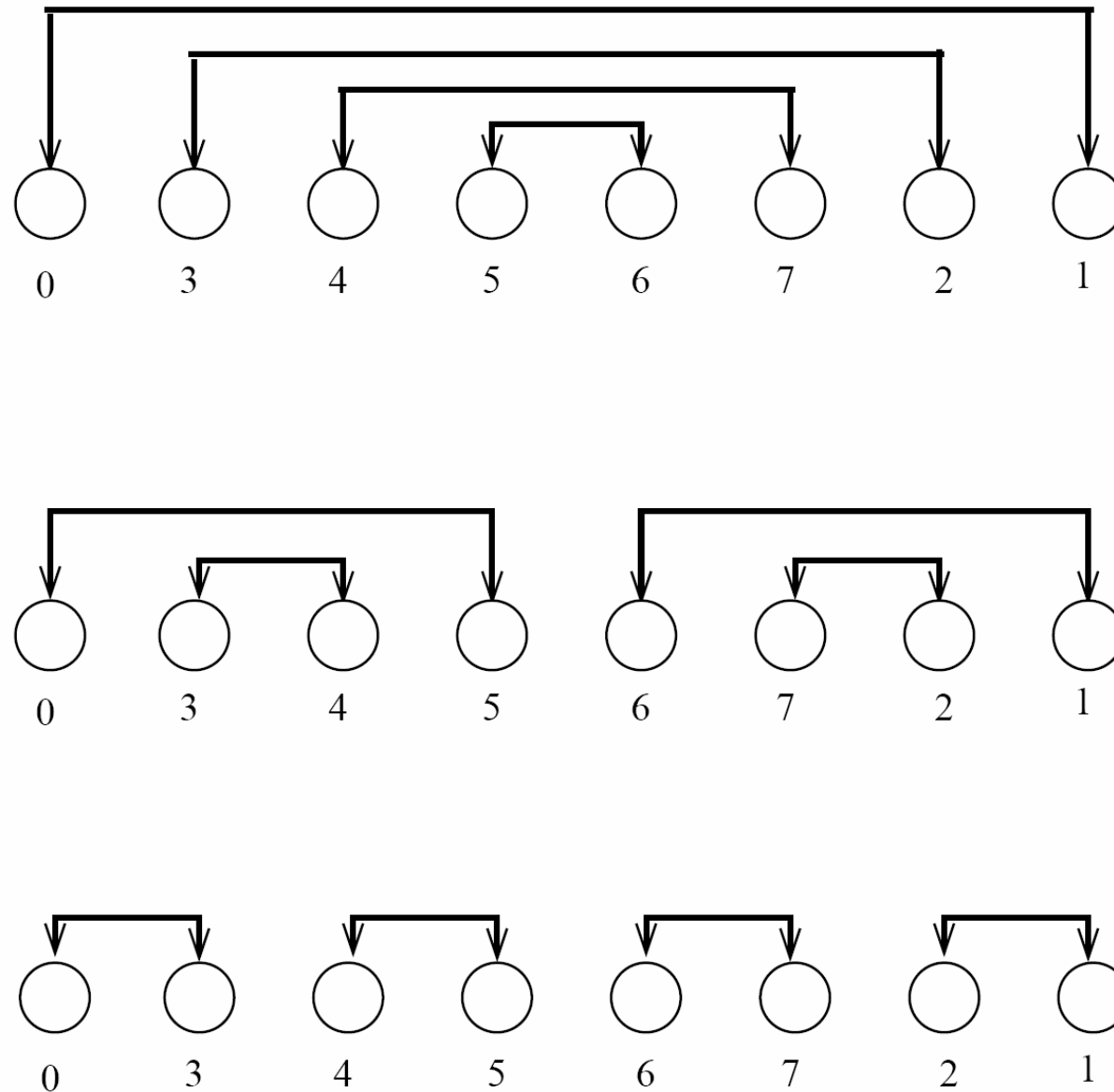


Figure 9.14 An example of the first phase of parallel shellsort on an eight-process array.



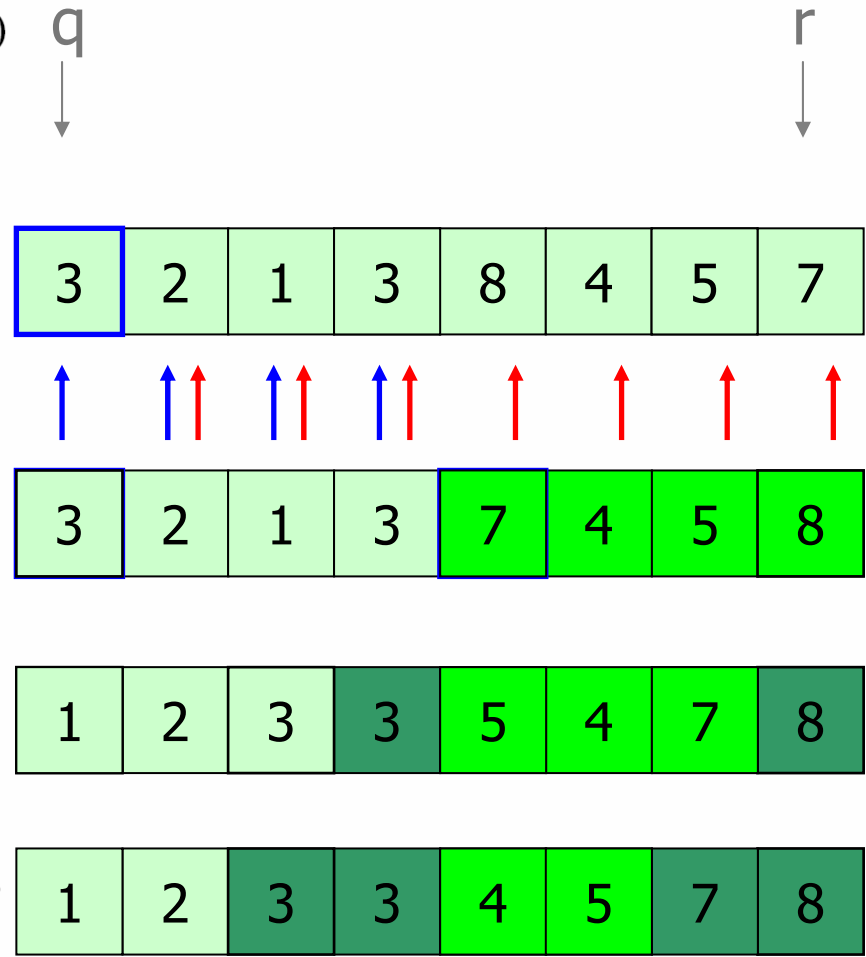
Quicksort

- Average complexity: $O(n \log n)$.
 - But very efficient in practice.
 - Average “robust”.
 - Low overhead and very simple.
- Divide & conquer algorithm:
 - Partition $A[q..r]$ into $A[q..s] \leq A[s+1..r]$.
 - Recursively sort sub-arrays.
 - Subtlety: How to partition?

```

1.  procedure QUICKSORT ( $A, q, r$ )
2.  begin
3.    if  $q < r$  then
4.      begin
5.         $x := A[q]$ ;
6.         $\rightarrow s := q$ ;
7.         $\rightarrow$  for  $i := q + 1$  to  $r$  do
8.          if  $A[i] \leq x$  then
9.            begin
10.              $s := s + 1$ ;
11.              $\text{swap}(A[s], A[i])$ ;
12.           end if
13.            $\text{swap}(A[q], A[s])$ ;
14.           QUICKSORT ( $A, q, s$ );
15.           QUICKSORT ( $A, s + 1, r$ );
16.         end if
17.       end QUICKSORT

```



Algorithm 9.5 The sequential quicksort algorithm.

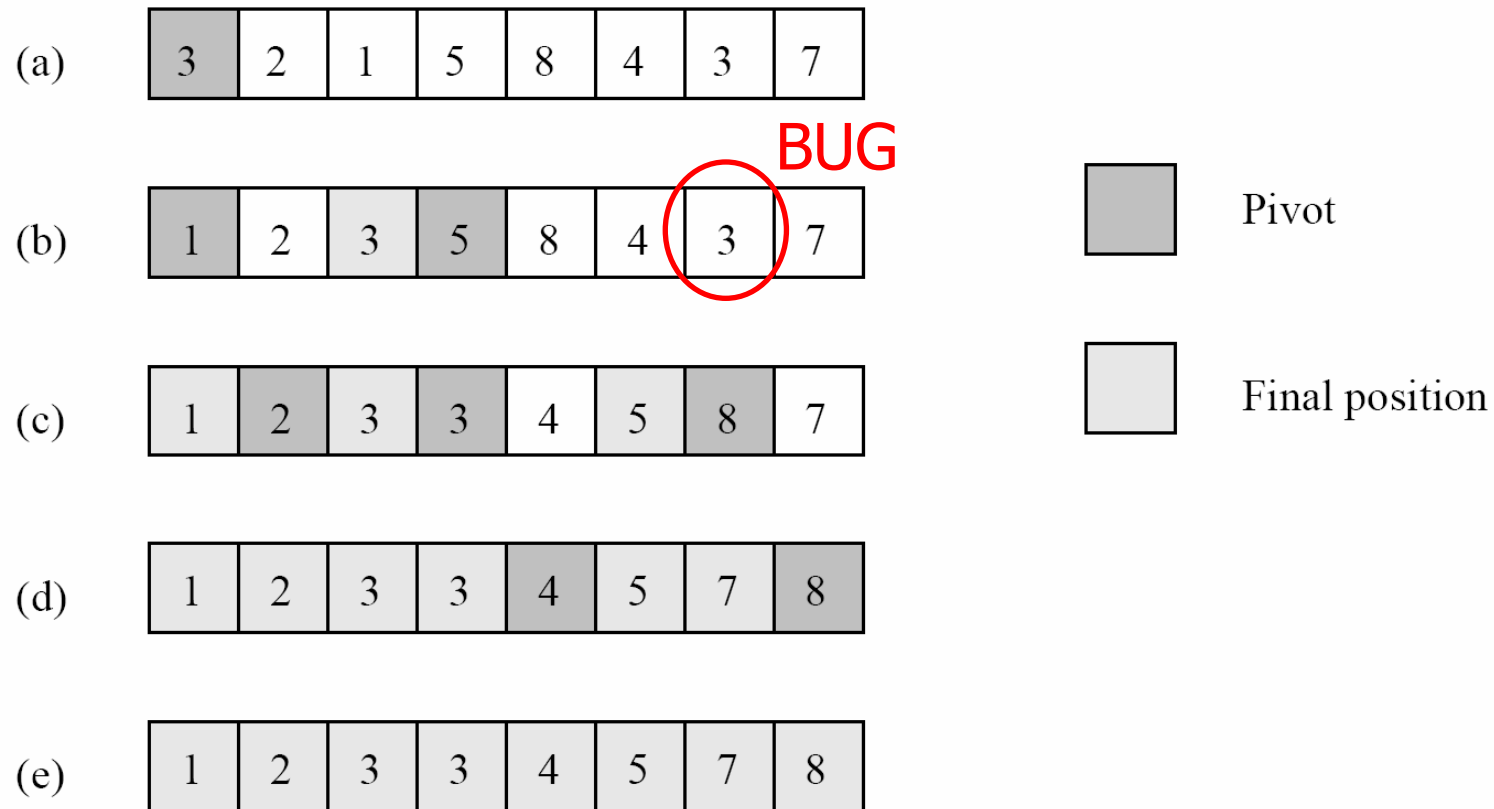


Figure 9.15 Example of the quicksort algorithm sorting a sequence of size $n = 8$.



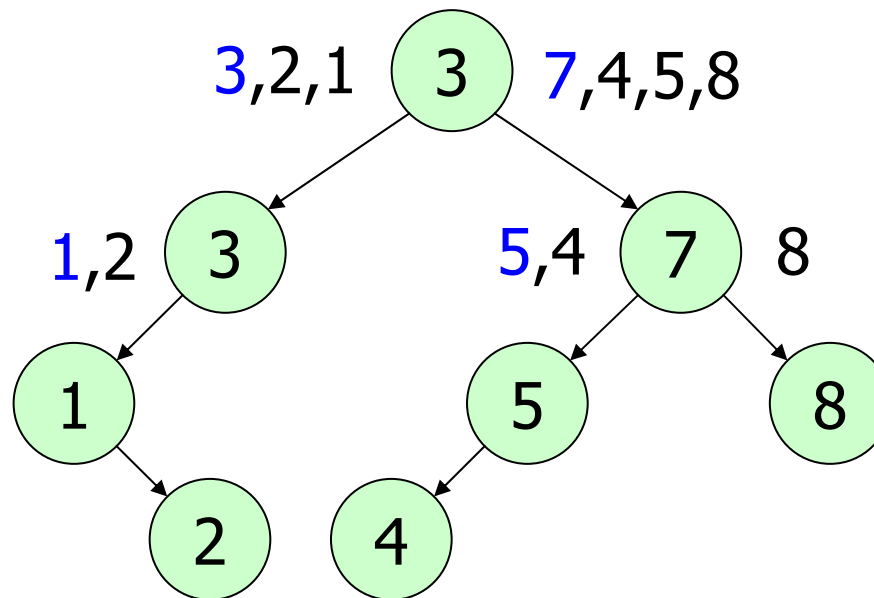
Parallel Quicksort

- Simple version:
 - Recursive decomposition with one process per recursive call.
 - Not cost optimal: Lower bound = n (initial partitioning).
 - Best we can do: Use $O(\log n)$ processes.
 - Need to parallelize the partitioning step.

Parallel Quicksort for CRCW

PRAM

- See execution of quicksort as constructing a binary tree.



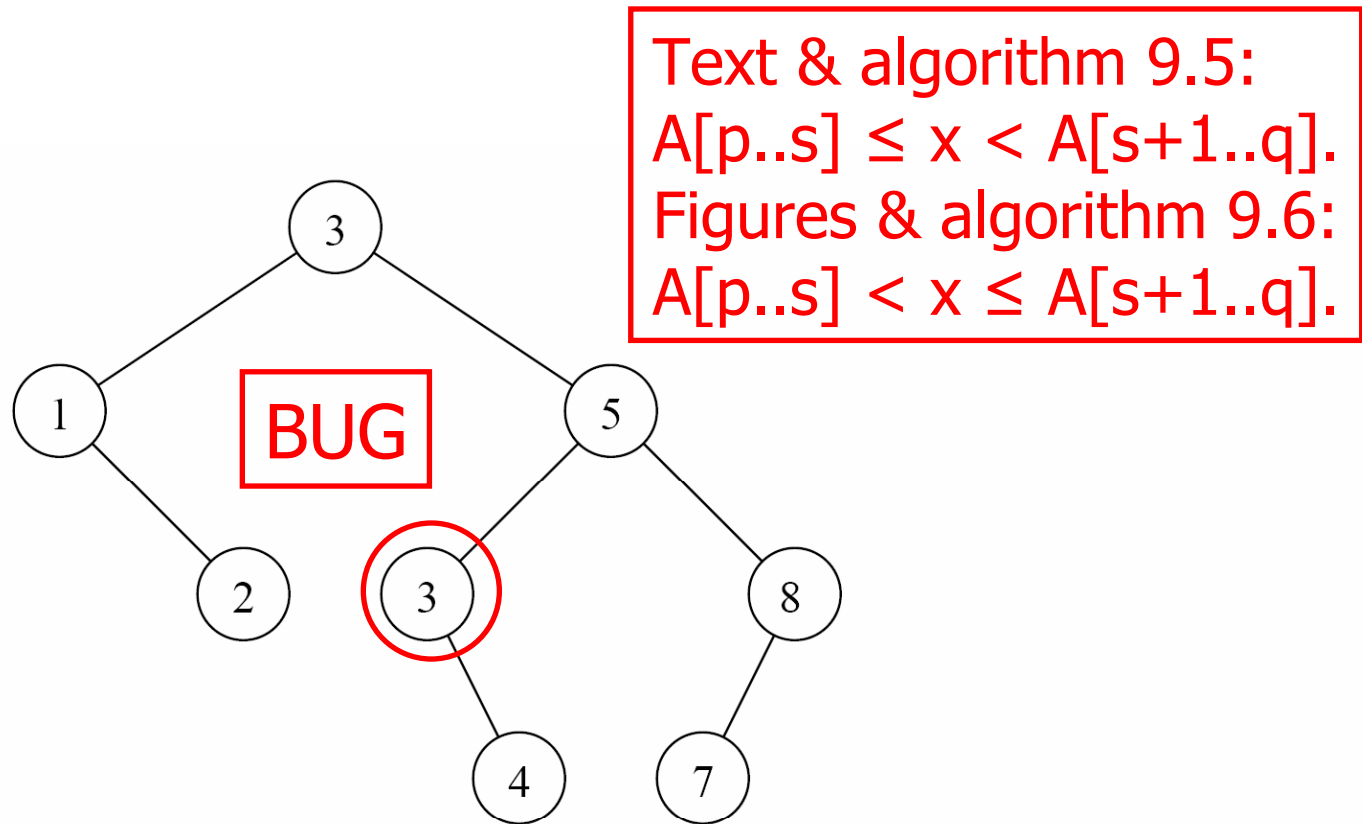
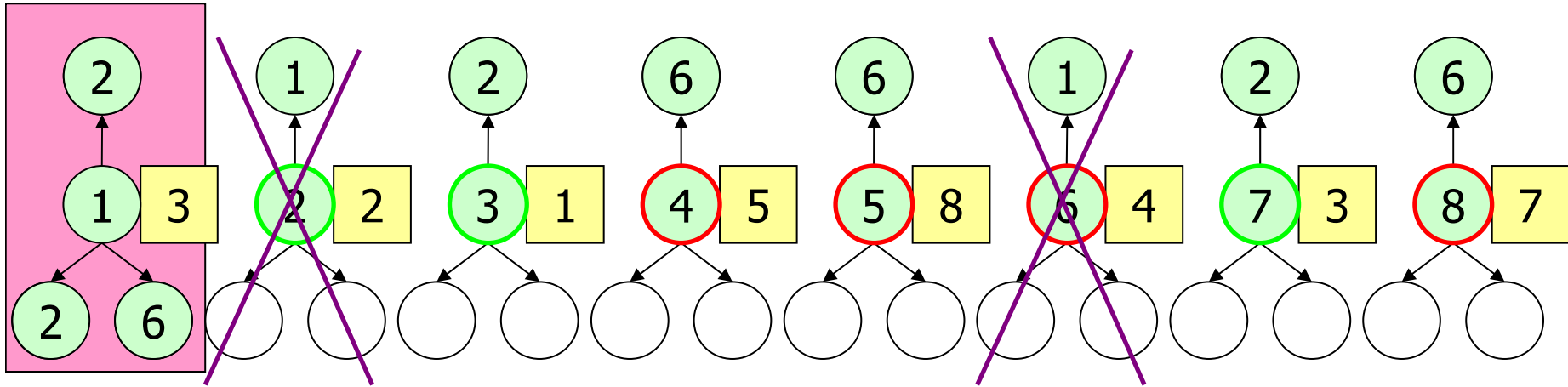


Figure 9.16 A binary tree generated by the execution of the quicksort algorithm. Each level of the tree represents a different array-partitioning iteration. If pivot selection is optimal, then the height of the tree is $\Theta(\log n)$, which is also the number of iterations.

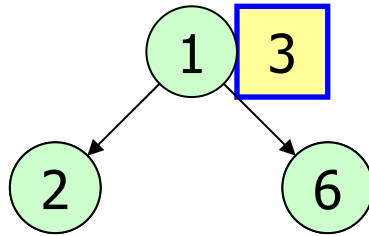
```

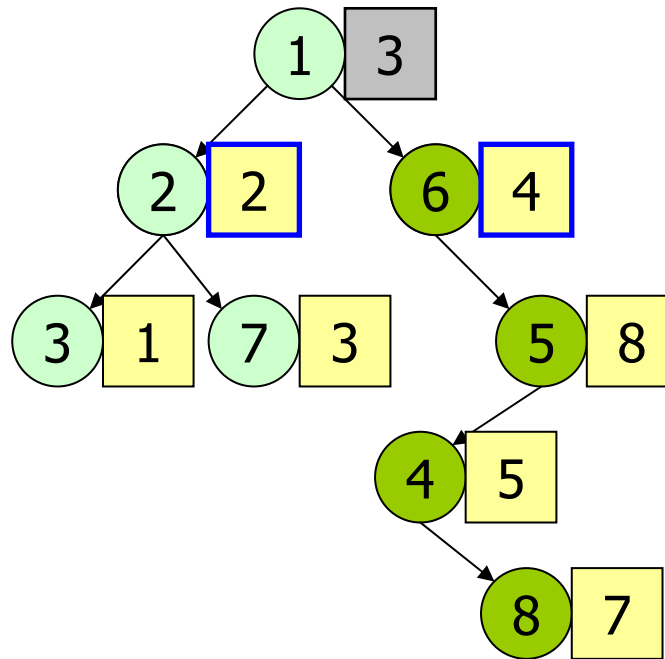
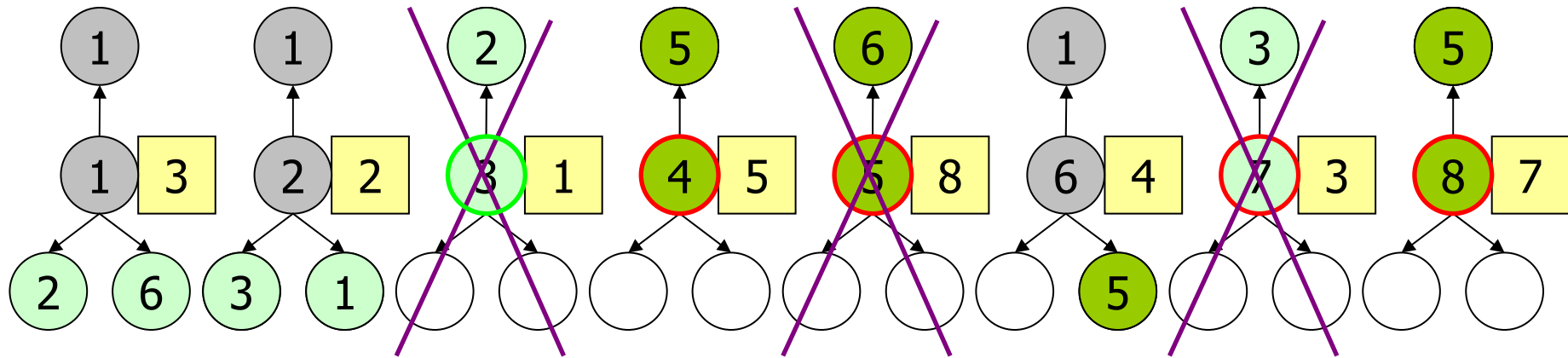
1.  procedure BUILD_TREE ( $A[1 \dots n]$ )
2.  begin
3.    for each process  $i$  do
4.    begin
5.       $root := i;$  only one succeeds
6.       $parent_i := root;$ 
7.       $leftchild[i] := rightchild[i] := n + 1;$ 
8.    end for
9.    repeat for each process  $i \neq root$  do
10.   begin
11.     if ( $A[i] \leq A[parent_i]$ ) then
12.       begin
13.          $leftchild[parent_i] := i;$ 
14.         if  $i = leftchild[parent_i]$  then exit
15.         else  $parent_i := leftchild[parent_i];$ 
16.       end for
17.     else
18.       begin
19.          $rightchild[parent_i] := i;$ 
20.         if  $i = rightchild[parent_i]$  then exit
21.         else  $parent_i := rightchild[parent_i];$ 
22.       end else
23.     end repeat
24.   end BUILD_TREE

```

root=1



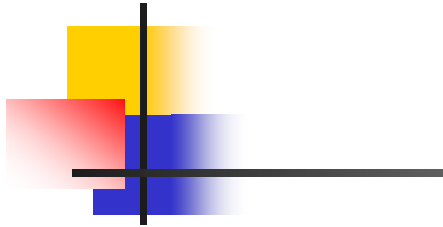


Each step: $\Theta(1)$. Average height: $\Theta(\log n)$.
 This is cost-optimal – but it is only a model.

Parallel Quicksort – Shared Address (Realistic)



- Same idea but remove contention:
 - Choose the pivot & broadcast it.
 - Each process rearranges its block of elements *locally*.
 - *Global* rearrangement of the blocks.
 - When the blocks reach a certain size, local sort is used.



First Step

P_0	P_1	P_2	P_3	P_4
7	13	18	2	17
1	14	20	6	10
15	9	3	16	19
4	11	12	5	8

pivot=7

P_0	P_1	P_2	P_3	P_4
7	2	18	13	1
17	14	20	6	10
15	9	3	4	19
16	5	12	11	8

after local rearrangement

7	2	1	6	3	4	5	18	13	17	14	20	10	15	9	19	16	12	11	8
---	---	---	---	---	---	---	----	----	----	----	----	----	----	---	----	----	----	----	---

after global rearrangement

Second Step

P_0	P_1	P_2	P_3	P_4
7	2	1	6	3
4	5	18	13	17
14	20	10	15	9
19	16	12	11	8

pivot=5 pivot=17

P_0	P_1	P_2	P_3	P_4
1	2	7	6	3
4	5	14	13	17
18	20	10	15	9
19	16	12	11	8

after local rearrangement

1	2	3	4	5	7	6	14	13	17	10	15	9	16	12	11	8	18	20	19
---	---	---	---	---	---	---	----	----	----	----	----	---	----	----	----	---	----	----	----

after global rearrangement

Third Step

P_0	P_1	P_2	P_3	P_4
1	2	3	4	5
7	6	14	13	17
10	15	9	16	12
11	8	18	20	19

pivot=11

P_0	P_1	P_2	P_3	P_4
1	2	3	4	5
6	7	10	13	17
14	15	9	8	12
11	16	18	19	20

after local rearrangement

10	9	8	12	11	13	17	14	15	16
----	---	---	----	----	----	----	----	----	----

after global rearrangement

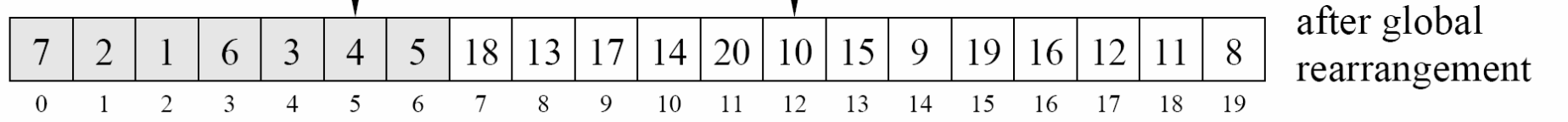
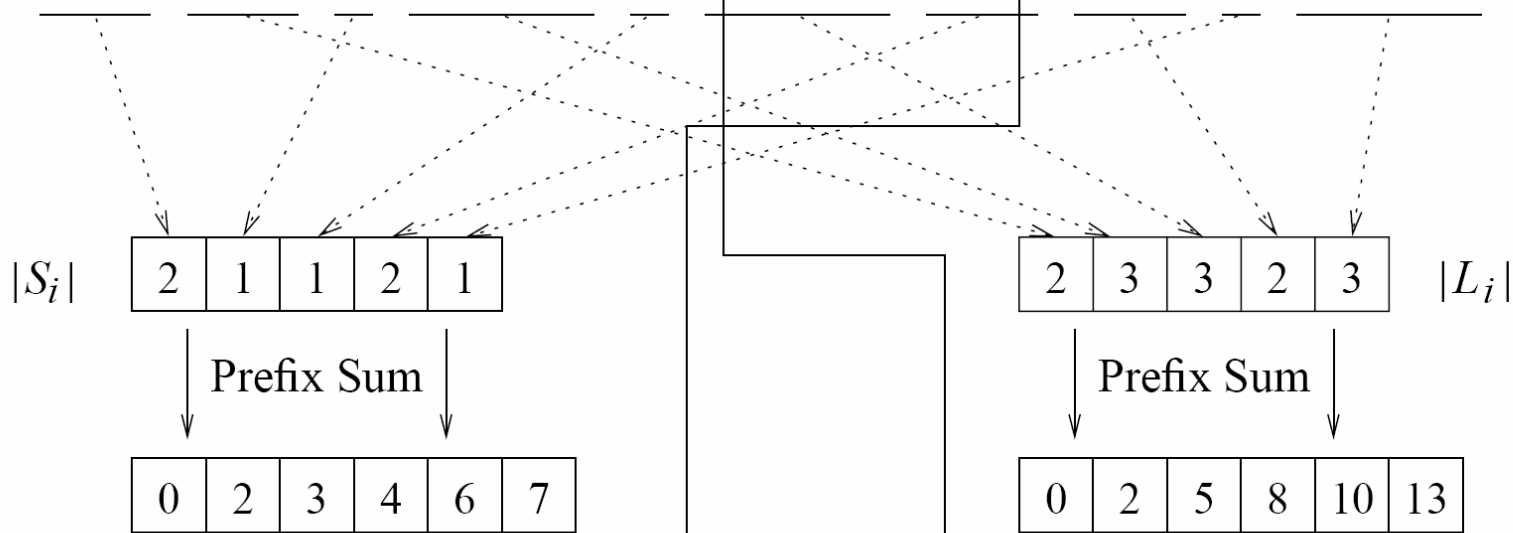
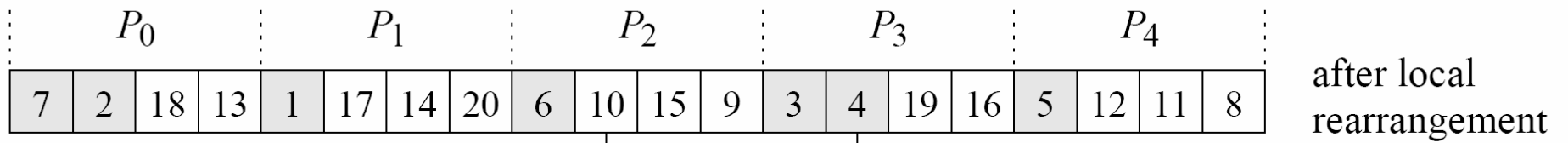
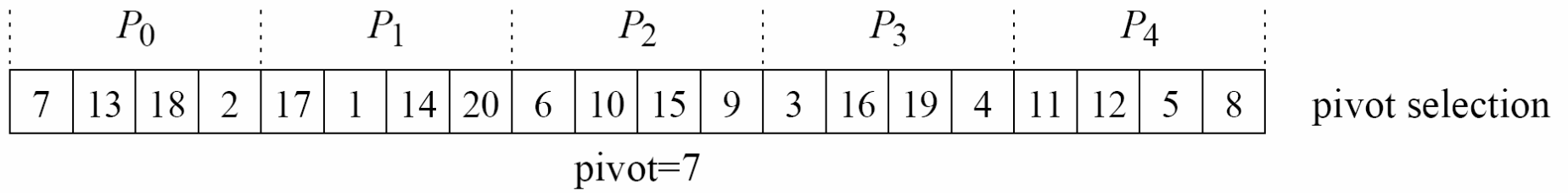
Fourth Step

10	9	8	12	11	13	17	14	15	16
----	---	---	----	----	----	----	----	----	----

after local rearrangement

P_0	P_1	P_2	P_3	P_4
1	2	3	4	5
6	7	8	9	10
11	12	13	14	15
16	17	18	19	20

Solution





Cost

- Scalability determined by time to broadcast the pivot & compute the prefix-sums.
- Cost optimal.

$$T_P = \overbrace{\Theta\left(\frac{n}{p} \log \frac{n}{p}\right)}^{\text{local sort}} + \overbrace{\Theta\left(\frac{n}{p} \log p\right) + \Theta(\log^2 p)}^{\text{array splits}}. \quad (4)$$



MPI Formulation of Quicksort

- Arrays must be explicitly distributed.
- Two phases:
 - Local partition smaller/larger than pivot.
 - Determine who will sort the sub-arrays.
 - And send the sub-arrays to the right process.



Final Word

- Pivot selection is very important.
- Affects performance.
- Bad pivot means idle processes.