



Sorting (Chapter 9)

Alexandre David

B2-206



Sorting

Problem

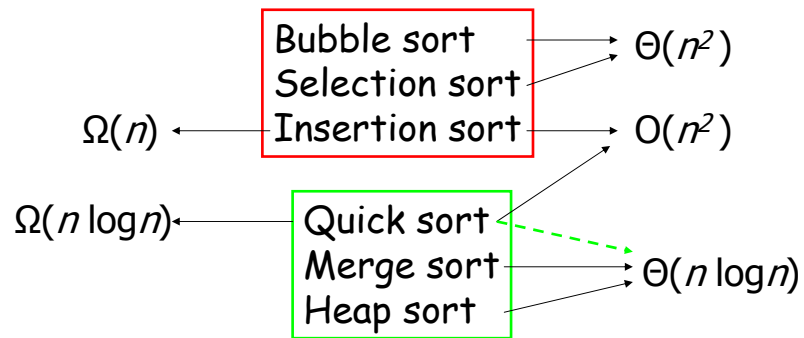
Arrange an unordered collection of elements into monotonically increasing (or decreasing) order.

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 .

The elements to sort (actually used for comparisons) are also called the keys.

Recall on Comparison Based Sorting Algorithms



21-04-2006

Alexandre David, MVP'06

3

You should know these complexities from a previous course on 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.

We assume internal sorting is possible.



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.

21-04-2006

Alexandre David, MVP'06

5

We assume comparison based sorting is used.



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

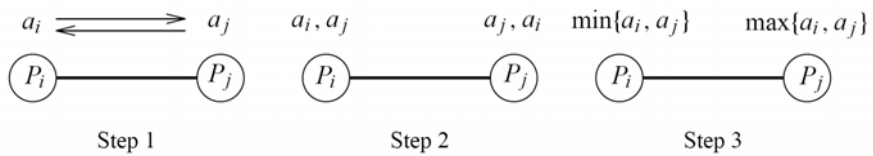
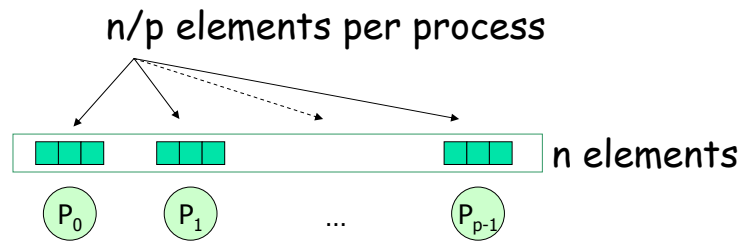


Figure 9.1 A parallel compare-exchange operation. Processes P_i and P_j send their elements to each other. Process P_i keeps $\min\{a_i, a_j\}$, and P_j keeps $\max\{a_i, a_j\}$.

Communication cost: $t_s + t_w$.

Comparison cost much cheaper \Rightarrow communication time dominates.

Blocks of Elements Per Process



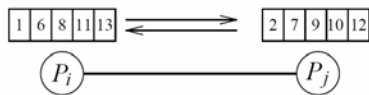
Blocks: $A_0 \leq A_1 \leq \dots \leq A_{p-1}$



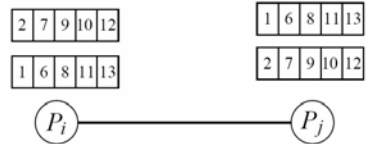
Compare-Split

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

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

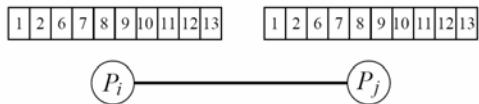


Step 1



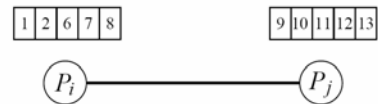
Step 2

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

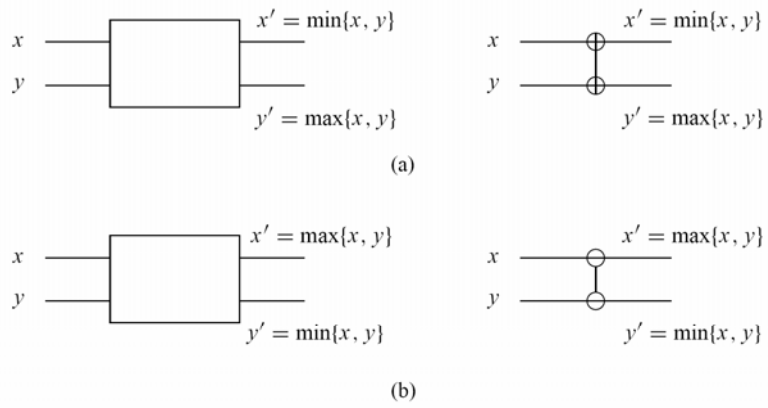


Figure 9.3 A schematic representation of comparators: (a) an increasing comparator, and (b) a decreasing comparator.

Sorting Networks

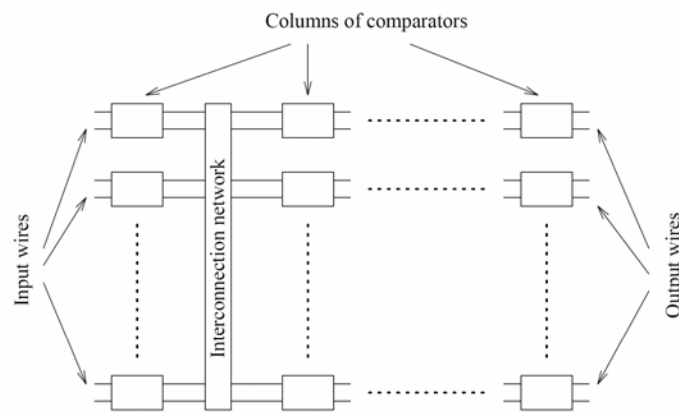


Figure 9.4 A typical sorting network. Every sorting network is made up of a series of columns, and each column contains a number of comparators connected in parallel.



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.

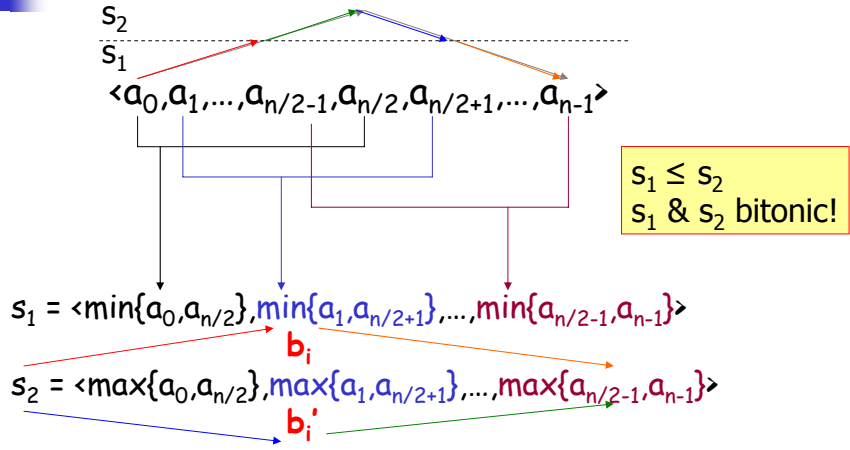
Example: $\langle 1, 2, 4, 7, 6, 0 \rangle$ & $\langle 8, 9, 2, 1, 0, 4 \rangle$ are bitonic sequences.



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

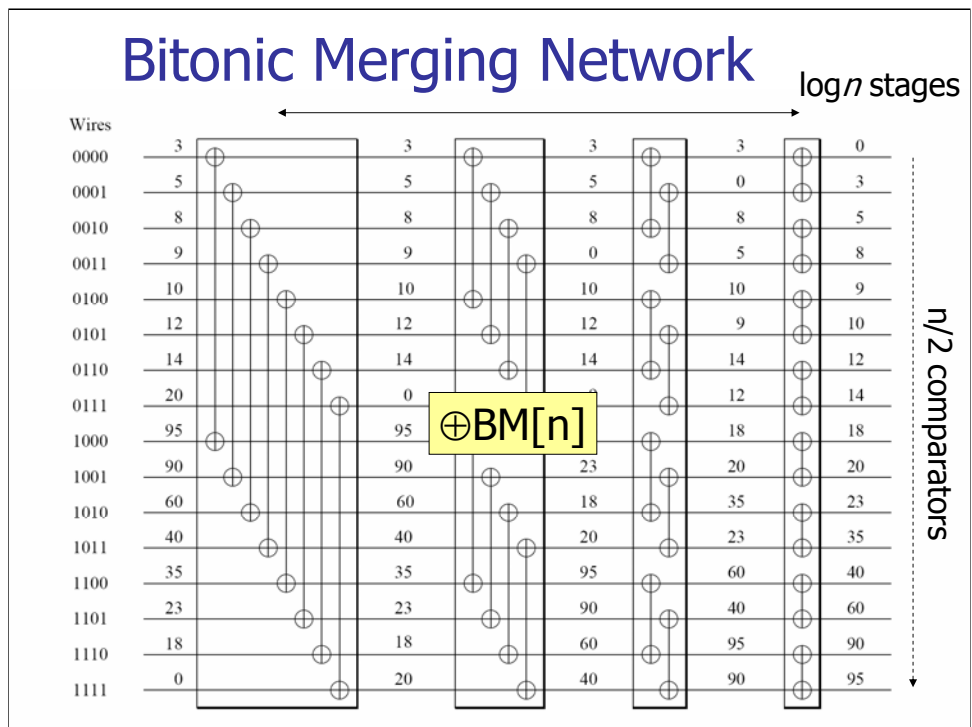


21-04-2006

Alexandre David, MVP'06

15

And in fact the procedure works even if the original sequence needs a cyclic shift to look like this particular case.



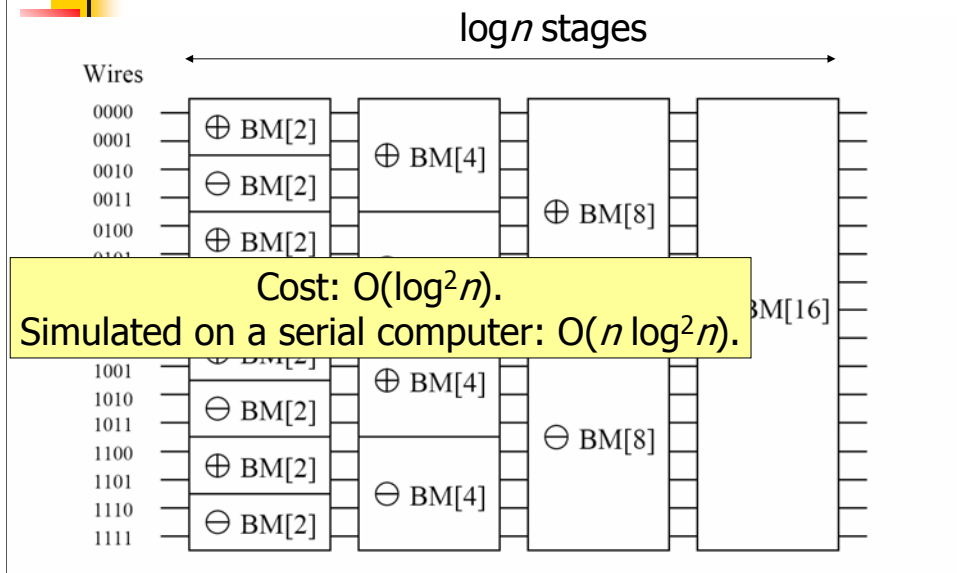
Cost: $\Theta(\log n)$ obviously.



Bitonic Sort

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

Bitonic Sort



Not cost optimal compared to the optimal serial algorithm.

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 with lower connectivity, several solutions but worse than the hypercube because the sequential algorithm is suboptimal.
 $T_p = O(\log n) + O(n/p)$ for 1 element/process.
 - Block of elements: sort locally ($n/p \log n/p$) & use bitonic merge ⇒ cost optimal.

21-04-2006

Alexandre David, MVP'06

19

Hypercube: Neighbors differ with each other by one bit.



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.

It is difficult to sort n elements in time $\log n$ using n processes (cost optimal w.r.t. the best serial algorithm in $n \log n$) but it is easy to parallelize other (less efficient) algorithms.



Odd-Even Transposition Sort

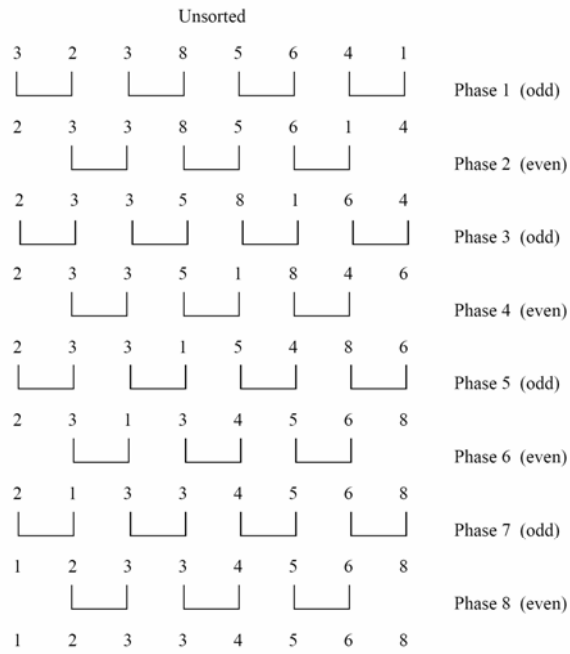
```
1. procedure ODD-EVEN( $n$ )
2. begin
3.   for  $i := 1$  to  $n$  do
4.     begin
5.       if  $i$  is odd then
6.         for  $j := 0$  to  $n/2 - 1$  do
7.           compare-exchange( $a_{2j+1}, a_{2j+2}$ );
8.       if  $i$  is even then
9.         for  $j := 1$  to  $n/2 - 1$  do
10.      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

$\Theta(n^2)$

$(a_1, a_2), (a_3, a_4) \dots$

$(a_2, a_3), (a_4, a_5) \dots$

Algorithm 9.3 Sequential odd-even transposition sort algorithm.



Phase 1 (odd)
Phase 2 (even)
Phase 3 (odd)
Phase 4 (even)
Phase 5 (odd)
Phase 6 (even)
Phase 7 (odd)
Phase 8 (even)

21-04-2006

Sorted

22

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

Write speedup & efficiency to find the bound on p but you can also see it with T_P .



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.

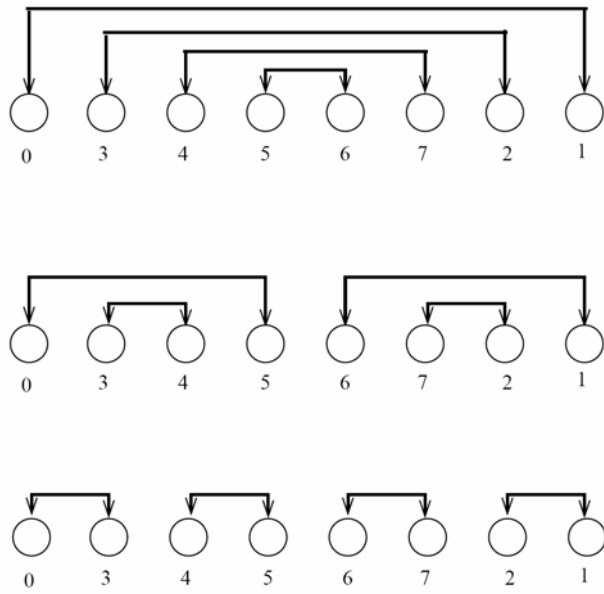


Figure 9.14 An example of the first phase of parallel shellsort on an eight-process array.

21-04-2006

Alexandre David, MVP'06

25



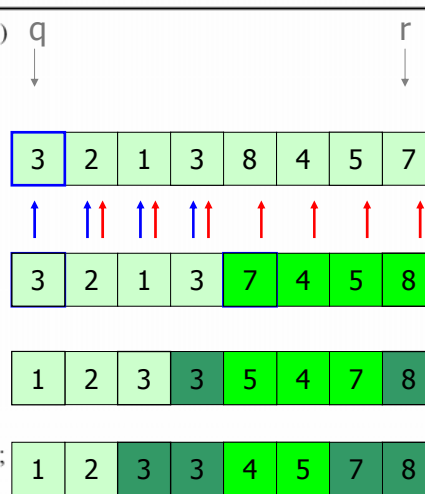
Quicksort

- 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.       s := q;
7.       for i := q + 1 to r do
8.         if A[i] ≤ x then
9.           begin
10.            s := s + 1;
11.            swap(A[s], A[i]);
12.          end if
13.        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.

Hoare partitioning is better. Check in your algorithm course.

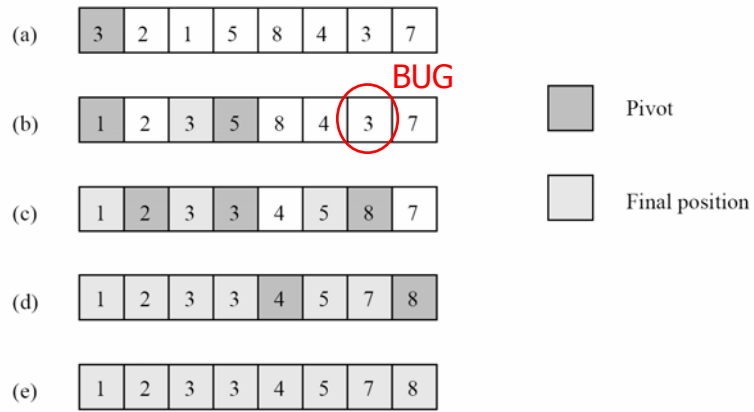


Figure 9.15 Example of the quicksort algorithm sorting a sequence of size $n = 8$.



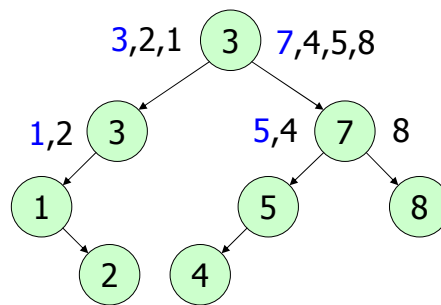
Parallel Quicksort

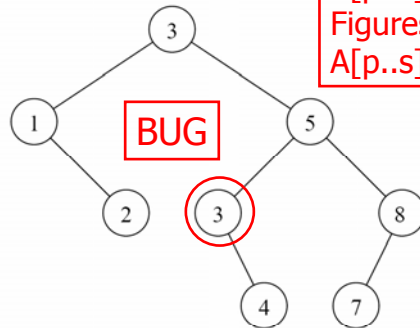
- Simple version:
 - Recursive decomposition with one process per recursive call.
 - Not cost optimal: Lower bound = n (initial partitioning).
 - Best we can do: Use $O(\log n)$ processes.
 - Need to parallelize the partitioning step.

Parallel Quicksort for CRCW

PRAM

- See execution of quicksort as constructing a binary tree.





Text & algorithm 9.5:
 $A[p..s] \leq x < A[s+1..q]$.
Figures & algorithm 9.6:
 $A[p..s] < x \leq A[s+1..q]$.

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.

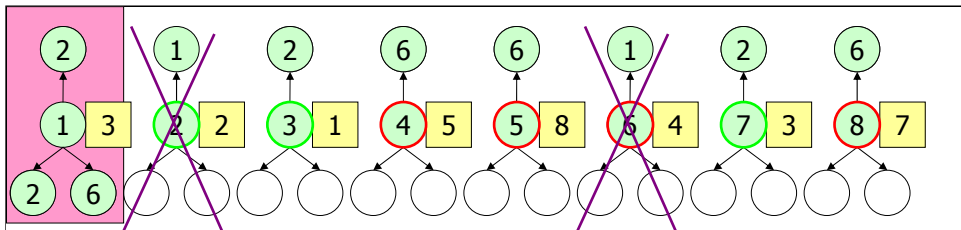
```

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

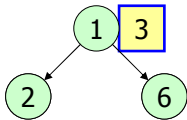
```

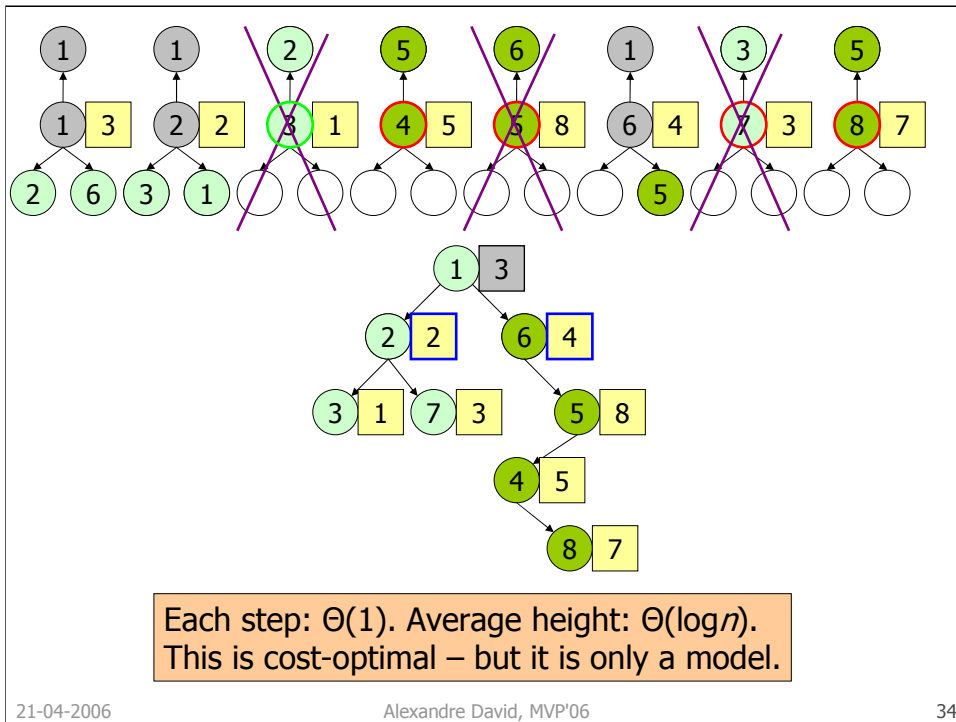
32

This algorithm does not correspond exactly to the serial version. Time for partitioning: $O(1)$.



root=1

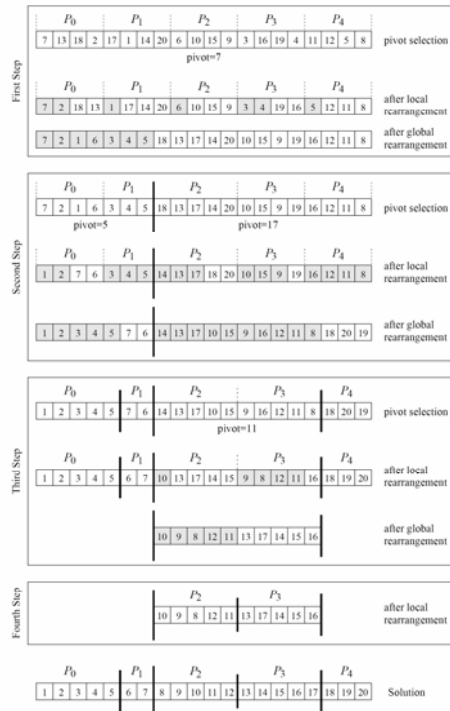






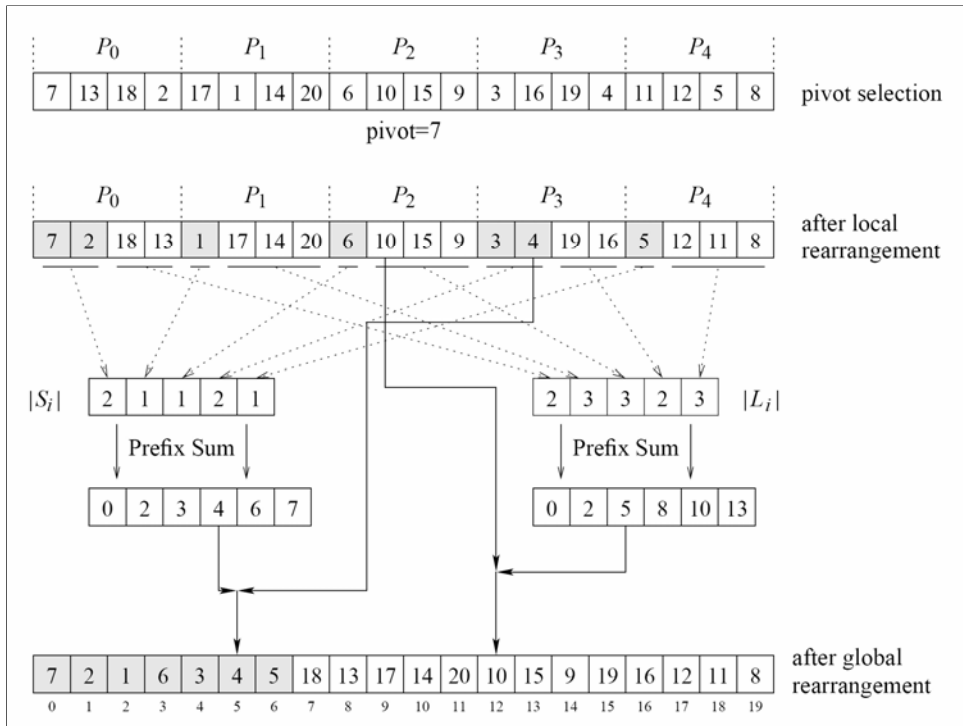
Parallel Quicksort – Shared Address (Realistic)

- 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.



21-04-2006

36





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.