
ALGORITHMES PARALLELES : tris

Daniel Etiemble
LRI, Université Paris Sud
de@lri.fr

Références

- Parallel sorting algorithms, <http://web.mst.edu/~ercal/387>
- Parallel Programming in C with MPI and OpenMP, Chapter 14: Sorting, Michael J. Quinn

Tris séquentiels

- Tris simples
 - Tri bulle
 - Simple
 - Pair impair
 - Tri insertion
 - Etc
- Tris plus efficaces
 - Tri rapide (quick sort)
 - Choix d'un pivot
 - Liste de tous les éléments inférieurs au pivot
 - Liste de tous les éléments supérieurs au pivot
 - Appel récursifs sur les listes créées
 - Algorithmes utiles
 - Fusion de deux listes triées

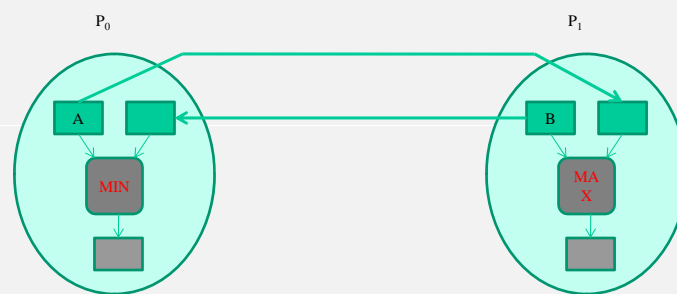
Algorithmes parallèles

- Hypothèses
 - Mémoire distribuée
 - La liste des N éléments est distribuée sur p processeurs avec N/p éléments par processeur
 - Mémoire partagée
 - Les N éléments sont dans la mémoire partagée
 - Accès par les p processeurs à la mémoire partagée
 - Caches

Compare-and-Exchange Sorting Algorithms

Form the basis of several, if not most, classical sequential sorting algorithms.

Two numbers, say A and B , are compared between P_0 and P_1 .



Polytech4 - Option
parallélisme - 2014

D. Etiemble

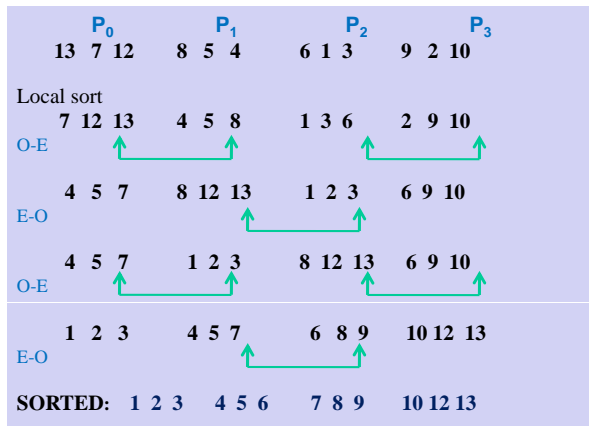
Odd-Even Transposition Sort - example

Step	P_0	P_1	P_2	P_3	P_4	P_5	P_6	P_7
0	4	2	7	8	5	1	3	6
1	2	4	7	8	1	5	3	6
2	2	4	7	1	8	3	5	6
3	2	4	1	7	3	8	5	6
4	2	1	4	3	7	5	8	6
5	1	2	3	4	5	7	6	8
6	1	2	3	4	5	6	7	8
7	1	2	3	4	5	6	7	8

Parallelism Option $T_{par} = O(n)$ (for $P=n$) D. Etiemble
parallélisme - 2014

Odd-Even Transposition Sort – Example (N >> P)

Each PE gets n/p numbers. First, PEs sort n/p locally, then they run odd-even trans. algorithm each time doing a merge-split for 2n/p numbers.

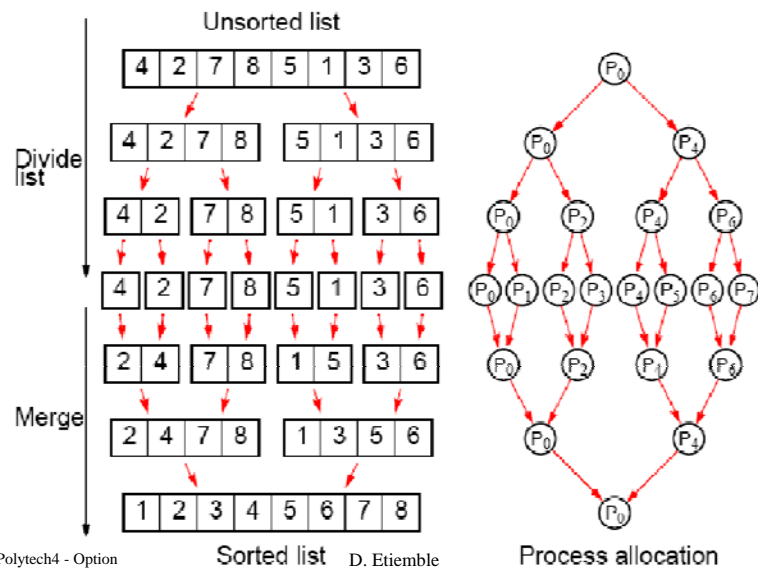


Time complexity: $T_{par} = (\text{Local Sort}) + (p \text{ merge-splits}) + (p \text{ exchanges})$

Polytech4 - Option
parallélisme - 2014

$$T_{par} = (n/p)\log(n/p) + p^*(n/p) + p^*(n/p) = (n/p)\log(n/p) + 2n$$

Parallelizing Mergesort



Polytech4 - Option
parallélisme - 2014

D. Etiemble

Mergesort - Time complexity

Sequential :

$$T_{seq} = 1 * n + 2 * \frac{n}{2} + 2^2 * \frac{n}{2^2} + \dots + 2^{\log n} * \frac{n}{2^{\log n}}$$

$$T_{seq} = O(n \log n)$$

Parallel :

$$T_{par} = 2 \left(\frac{n}{2^0} + \frac{n}{2^1} + \frac{n}{2^2} + \dots + \frac{n}{2^k} \dots + 2 + 1 \right)$$

$$= 2n(2^0 + 2^{-1} + 2^{-2} + \dots + 2^{-\log n})$$

$$T_{par} = O(4n)$$

Polytech4 - Option
parallélisme - 2014

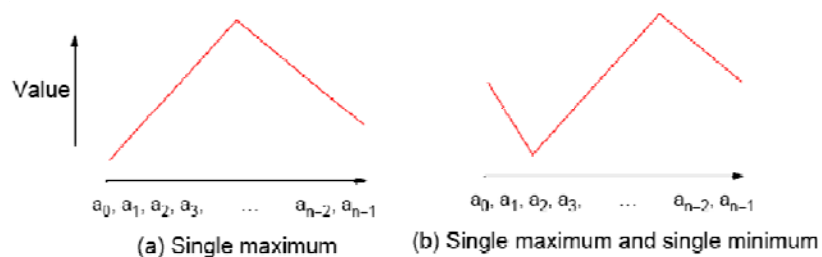
D. Etiemble

9

Bitonic Mergesort

Bitonic Sequence

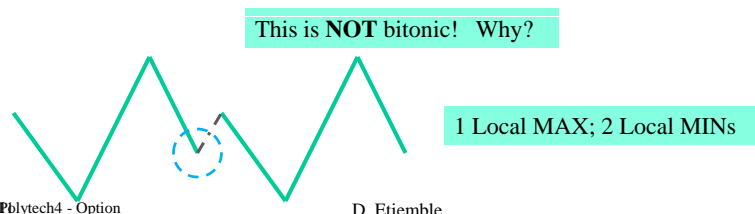
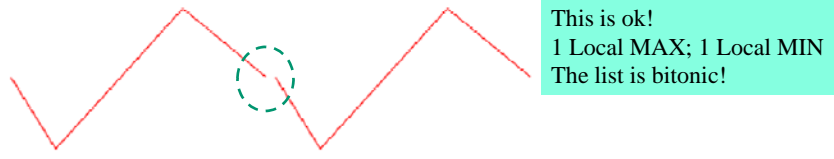
A *bitonic sequence* is defined as a list with no more than one **LOCAL MAXIMUM** and no more than one **LOCAL MINIMUM**.
(Endpoints must be considered - wraparound)



Polytech4 - Option
parallélisme - 2014

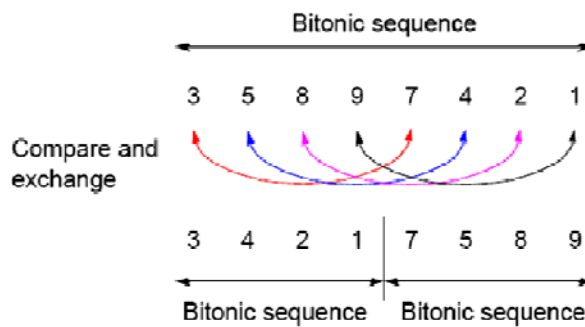
D. Etiemble

A *bitonic sequence* is a list with no more than one **LOCAL MAXIMUM** and no more than one **LOCAL MINIMUM**.
(Endpoints must be considered - wraparound)



Binary Split

1. Divide the **bitonic list** into two equal halves.
2. Compare-Exchange each item on the first half with the corresponding item in the second half.



Result:
Two bitonic sequences where the numbers in one sequence are all less than the numbers in the other sequence.

Repeated application of binary split

Bitonic list:

24 20 15 9 4 2 5 8 | 10 11 12 13 22 30 32 45

Result after Binary-split:

10 11 12 9 4 2 5 8 | 24 20 15 13 22 30 32 45

If you keep applying the **BINARY-SPLIT** to each half repeatedly, you will get a **SORTED LIST** !

10 11 12 9 4 2 5 8 | 24 20 15 13 22 30 32 45
 4 2 5 8 10 11 12 9 | 22 20 15 13 24 30 32 45
 4 2 5 8 10 9 12 11 | 15 13 22 20 24 30 32 45
 2 4 5 8 9 10 11 12 | 13 15 20 22 24 30 32 45

Q: How many parallel steps does it take to sort ?

Polytech4 - Option
parallélisme - 2014

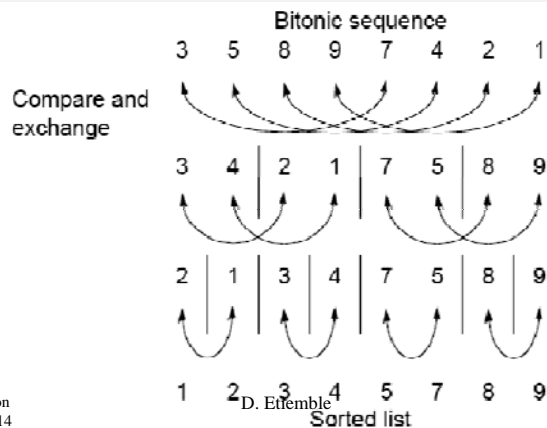
D. Etienne

13

Sorting a bitonic sequence

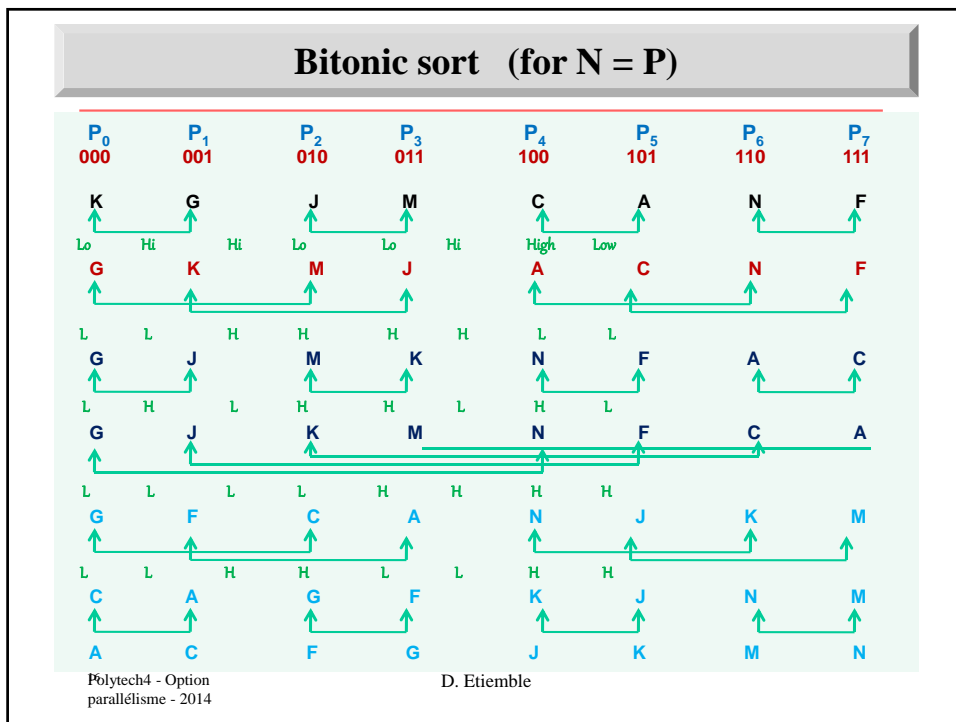
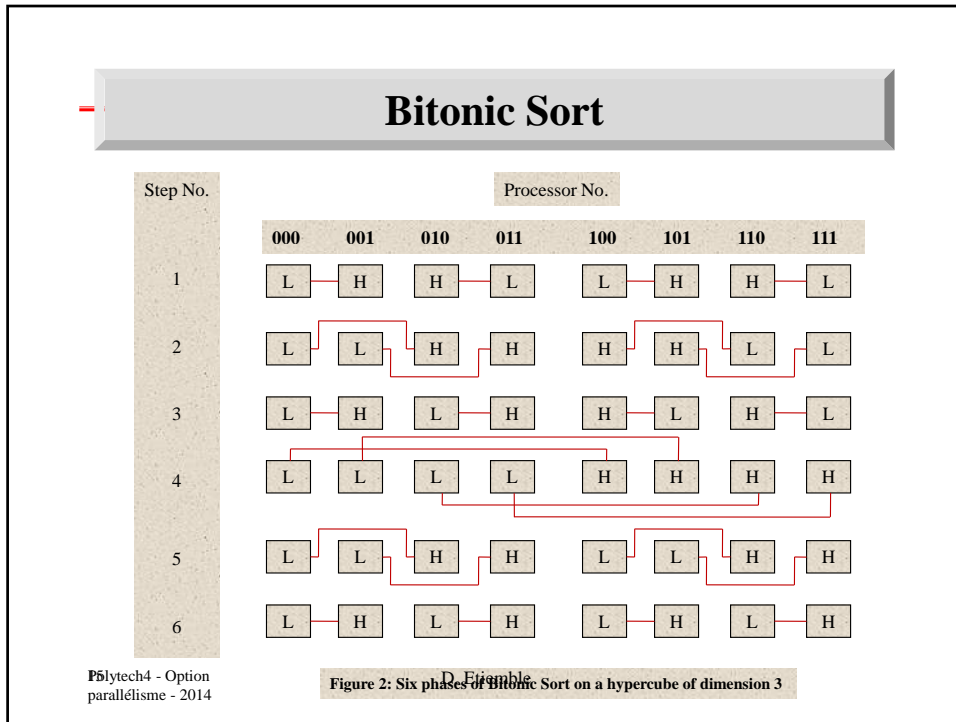
Compare-and-exchange moves smaller numbers of each pair to left and larger numbers of pair to right.

Given a bitonic sequence, recursively performing '*binary split*' will sort the list.



Polytech4 - Option
parallélisme - 2014

D. Etienne

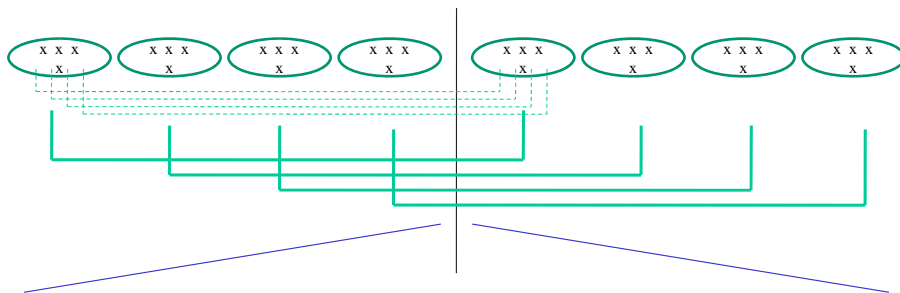


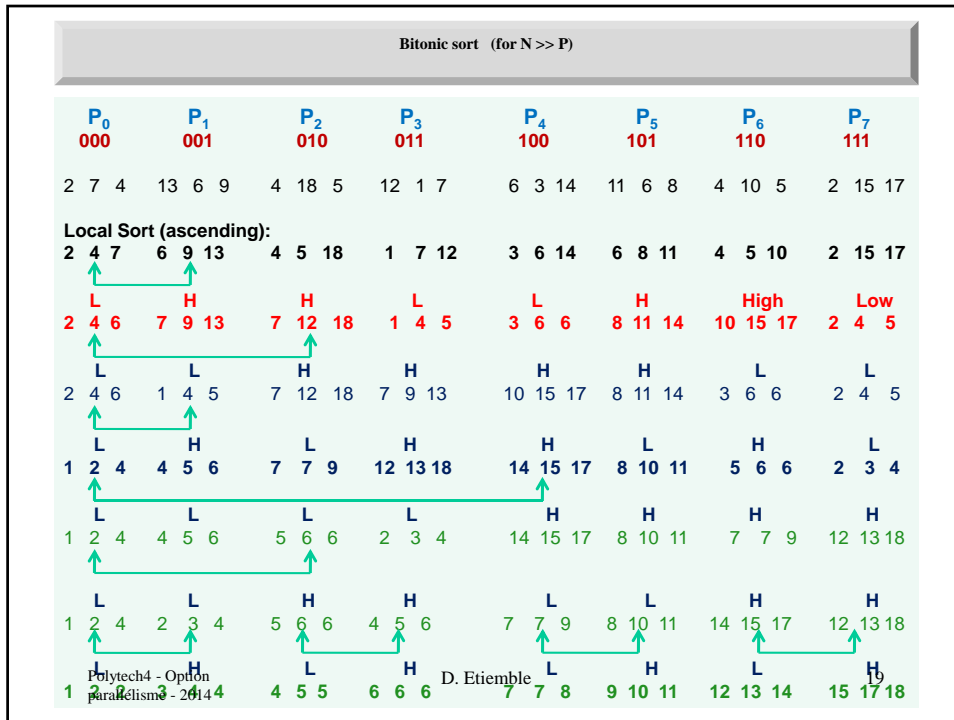
Number of steps (P=n)

In general, with $n = 2^k$, there are k phases, each of 1, 2, 3, ..., k steps.
Hence the total number of steps is:

$$T_{par}^{bitonic} = \sum_{i=1}^{i=\log n} i = \frac{\log n(\log n + 1)}{2} = O(\log^2 n)$$

Bitonic sort (for N >> P)





Number of steps (for $N \gg P$)

$$T_{par}^{bitonic} = \text{Local Sort} + \text{Parallel Bitonic Merge}$$

$$= \frac{N}{P} \log \frac{N}{P} + 2 \frac{N}{P} (1 + 2 + 3 + \dots + \log P)$$

$$= \frac{N}{P} \left\{ \log \frac{N}{P} + 2 \left(\frac{\log P (1 + \log P)}{2} \right) \right\}$$

$$= \frac{N}{P} (\log N - \log P + \log P + \log^2 P)$$

$$T_{par}^{bitonic} = \frac{N}{P} (\log N + \log^2 P)$$

Polytech4 - Option parallélisme - 2014

D. Etiemble

Sequential Quicksort

17 14 65 4 22 63 11

Unordered list of values

17 14 65 4 22 63 11

Choose pivot value

Sequential Quicksort

14 4 11 17 65 22 63

Low list
(≤ 17)

High list
(> 17)

Sequential Quicksort



Recursively apply quicksort
to low list



Recursively apply quicksort
to high list

Sequential Quicksort



Sorted list of values

Attributes of Sequential Quicksort

- Average-case time complexity: $\Theta(n \log n)$
- Worst-case time complexity: $\Theta(n^2)$
 - Occurs when low, high lists maximally unbalanced at every partitioning step
- Can make worst-case less probable by using sampling to choose pivot value
 - Example: “Median of 3” technique

Quicksort Good Starting Point for Parallel Algorithm

- Speed
 - Generally recognized as fastest sort in average case
 - Preferable to base parallel algorithm on fastest sequential algorithm
- Natural concurrency
 - Recursive sorts of low, high lists can be done in parallel

Definitions of “Sorted”

- Definition 1: Sorted list held in memory of a single processor
- Definition 2:
 - Portion of list in every processor’s memory is sorted
 - Value of last element on P_i ’s list is less than or equal to value of first element on P_{i+1} ’s list
- We adopt Definition 2: Allows problem size to scale with number of processors

Parallel Quicksort

75, 91, 15, 64, 21, 8, 88, 54

P_0

50, 12, 47, 72, 65, 54, 66, 22

P_1

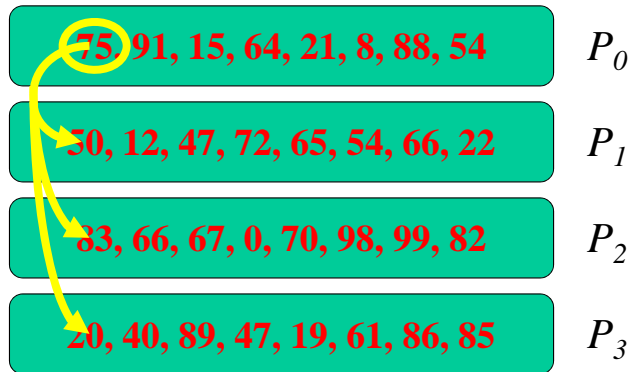
83, 66, 67, 0, 70, 98, 99, 82

P_2

20, 40, 89, 47, 19, 61, 86, 85

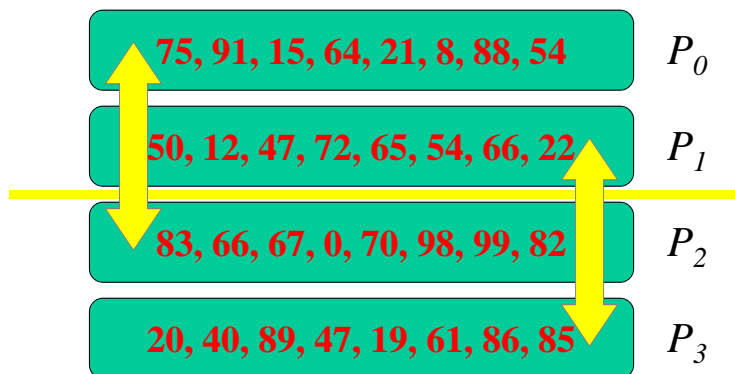
P_3

Parallel Quicksort



Process P_0 chooses and broadcasts
randomly chosen pivot value

Parallel Quicksort



Exchange “lower half” and “upper half” values”

Parallel Quicksort

Lower "half"	75, 15, 64, 21, 8, 54, 66, 67, 0, 70	P_0
	50, 12, 47, 72, 65, 54, 66, 22, 20, 40, 47, 19, 61	P_1
	83, 98, 99, 82, 91, 88	P_2
Upper "half"	89, 86, 85	P_3

After exchange step

Polytech4 - Option parallélisme - 2014
D. Etiemble
31

Parallel Quicksort

Lower "half"	75, 15, 64, 21, 8, 54, 66, 67, 0, 70	P_0
	50, 12, 47, 72, 65, 54, 66, 22, 20, 40, 47, 19, 61	P_1
	83, 98, 99, 82, 91, 88	P_2
Upper "half"	89, 86, 85	P_3

Processes P0 and P2 choose and broadcast randomly chosen pivots

Polytech4 - Option parallélisme - 2014
D. Etiemble
32

Parallel Quicksort

Lower
"half"

75, 15, 64, 21, 8, 54, 66, 67, 0, 70

P_0

Upper
"half"

50, 12, 47, 72, 65, 54, 66,
22, 20, 40, 47, 19, 61

P_1

83, 98, 99, 82, 91, 88

P_2

89, 86, 85

P_3

Exchange values

Polytech4 - Option
parallélisme - 2014
D. Etiemble
33

Parallel Quicksort

Lower "half"
of lower "half"

15, 21, 8, 0, 12, 20, 19

P_0

Upper "half"
of lower "half"

50, 47, 72, 65, 54, 66, 22, 40,
47, 61, 75, 64, 54, 66, 67, 70

P_1

Lower "half"
of upper "half"

83, 82, 91, 88, 89, 86, 85

P_2

Upper "half"
of upper "half"

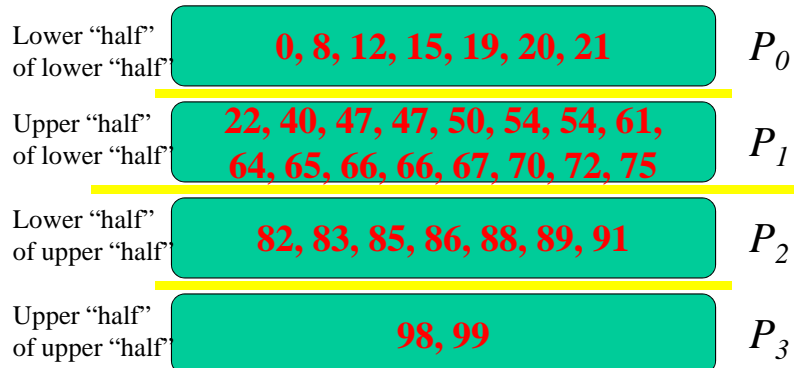
98, 99

P_3

Exchange values

Polytech4 - Option
parallélisme - 2014
D. Etiemble
34

Parallel Quicksort



Each processor sorts values it controls

Hyperquicksort

- Start where parallel quicksort ends: each process sorts its sublist
- First "sortedness" condition is met
- To meet second, processes must still exchange values
- Process can use median of its sorted list as the pivot value
- This is much more likely to be close to the true median

Hyperquicksort

75, 91, 15, 64, 21, 8, 88, 54

P_0

50, 12, 47, 72, 65, 54, 66, 22

P_1

83, 66, 67, 0, 70, 98, 99, 82

P_2

20, 40, 89, 47, 19, 61, 86, 85

P_3

Number of processors is a power of 2

Hyperquicksort

8, 15, 21, 54, 64, 75, 88, 91

P_0

12, 22, 47, 50, 54, 65, 66, 72

P_1

0, 66, 67, 70, 82, 83, 98, 99

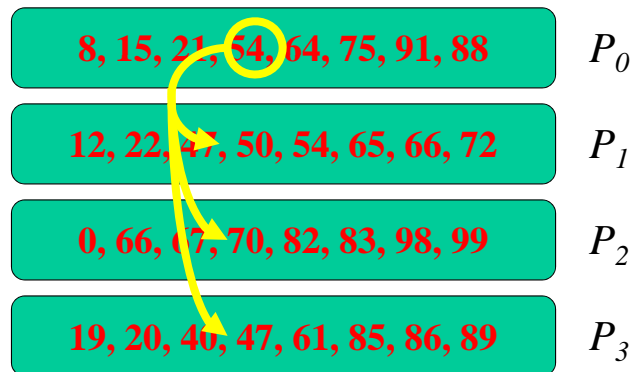
P_2

19, 20, 40, 47, 61, 85, 86, 89

P_3

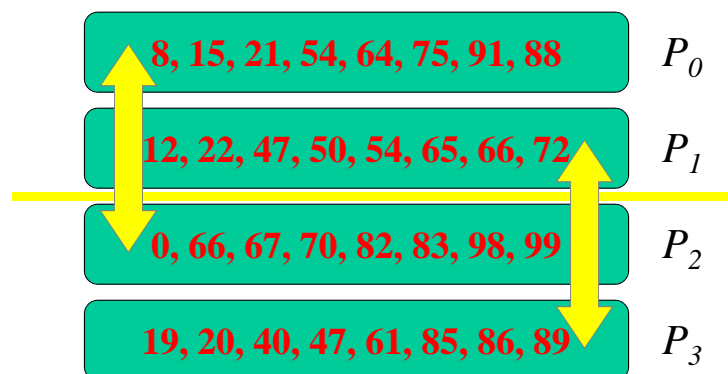
Each process sorts values it controls

Hyperquicksort



Process P_0 broadcasts its median value

Hyperquicksort



Processes will exchange “low”, “high” lists

Hyperquicksort

0, 8, 15, 21, 54 P_0

12, 19, 20, 22, 40, 47, 47, 50, 54 P_1

64, 66, 67, 70, 75, 82, 83, 88, 91, 98, 99 P_2

61, 65, 66, 72, 85, 86, 89 P_3

Processes merge kept and received values.

Hyperquicksort

0, 8, 15, 21, 54 P_0

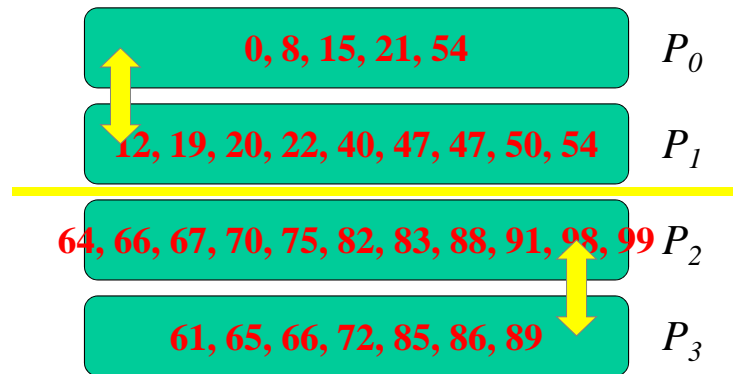
12, 19, 20, 22, 40, 47, 47, 50, 54 P_1

64, 66, 67, 70, 75, 82, 83, 88, 91, 98, 99 P_2

61, 65, 66, 72, 85, 86, 89 P_3

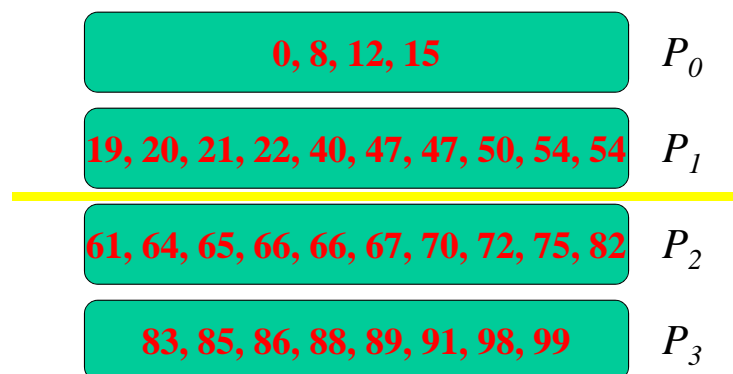
Processes P_0 and P_2 broadcast median values.

Hyperquicksort



Communication pattern for second exchange

Hyperquicksort



After exchange-and-merge step

Complexity Analysis Assumptions

- Average-case analysis
- Lists stay reasonably balanced
- Communication time dominated by message transmission time, rather than message latency

Complexity Analysis

- Initial quicksort step has time complexity $\Theta((n/p) \log(n/p))$
- Total comparisons needed for $\log p$ merge steps: $\Theta((n/p) \log p)$
- Total communication time for $\log p$ exchange steps: $\Theta((n/p) \log p)$

Another Scalability Concern

- Our analysis assumes lists remain balanced
- As p increases, each processor's share of list decreases
- Hence as p increases, likelihood of lists becoming unbalanced increases
- Unbalanced lists lower efficiency
- Would be better to get sample values from all processes before choosing median

Parallel Sorting by Regular Sampling (PSRS Algorithm)

- Each process sorts its share of elements
- Each process selects regular sample of sorted list
- One process gathers and sorts samples, chooses pivot values from sorted sample list, and broadcasts these pivot values
- Each process partitions its list into p pieces, using pivot values
- Each process sends partitions to other processes
- Each process merges its partitions

PSRS Algorithm

75, 91, 15, 64, 21, 8, 88, 54

P_0

50, 12, 47, 72, 65, 54, 66, 22

P_1

83, 66, 67, 0, 70, 98, 99, 82

P_2

Number of processors does not
have to be a power of 2.

PSRS Algorithm

8, 15, 21, 54, 64, 75, 88, 91

P_0

12, 22, 47, 50, 54, 65, 66, 72

P_1

0, 66, 67, 70, 82, 83, 98, 99

P_2

Each process sorts its list using quicksort.

PSRS Algorithm

8, 15, 21, 54, 64, 75, 88, 91 P_0

12, 22, 47, 50, 54, 65, 66, 72 P_1

0, 66, 67, 70, 82, 83, 98, 99 P_2

Each process chooses p regular samples.

PSRS Algorithm

8, 15, 21, 54, 64, 75, 88, 91 P_0

12, 22, 47, 50, 54, 65, 66, 72 P_1

0, 66, 67, 70, 82, 83, 98, 99 P_2

15, 54, 75, 22, 50, 65, 66, 70, 83

One process collects p^2 regular samples.

PSRS Algorithm

8, 15, 21, 54, 64, 75, 88, 91

P_0

12, 22, 47, 50, 54, 65, 66, 72

P_1

0, 66, 67, 70, 82, 83, 98, 99

P_2

15, 22, 50, 54, 65, 66, 70, 75, 83

One process sorts p^2 regular samples.

PSRS Algorithm

8, 15, 21, 54, 64, 75, 88, 91

P_0

12, 22, 47, 50, 54, 65, 66, 72

P_1

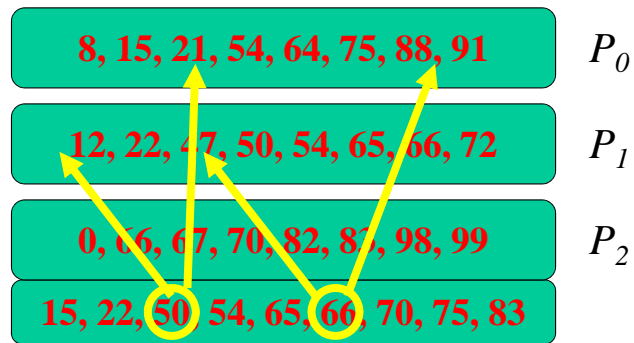
0, 66, 67, 70, 82, 83, 98, 99

P_2

15, 22, 50, 54, 65, 66, 70, 75, 83

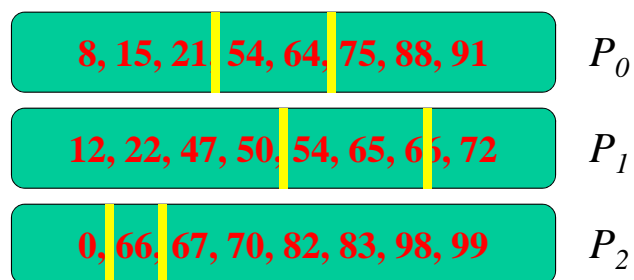
One process chooses $p-1$ pivot values.

PSRS Algorithm



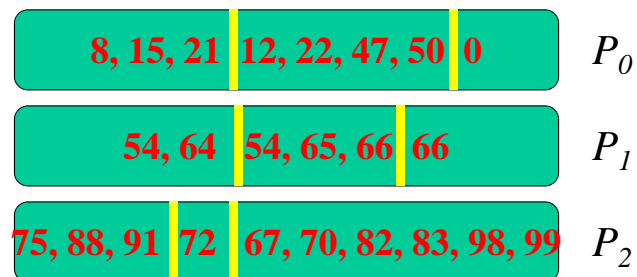
One process broadcasts $p-1$ pivot values.

PSRS Algorithm



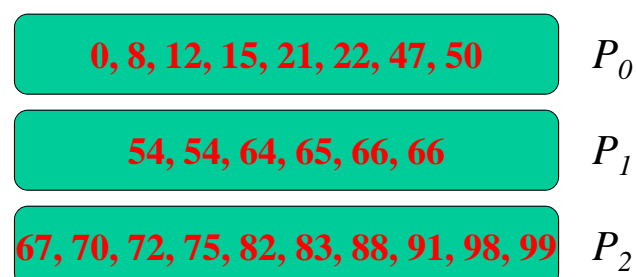
Each process divides list, based on pivot values.

PSRS Algorithm



Each process sends partitions to correct destination process.

PSRS Algorithm



Each process merges p partitions.

Assumptions

- Each process ends up merging close to n/p elements
- Experimental results show this is a valid assumption
- Processor interconnection network supports p simultaneous message transmissions at full speed
- 4-ary hypertree is an example of such a network

Time Complexity Analysis

- Computations
 - Initial quicksort: $\Theta((n/p)\log(n/p))$
 - Sorting regular samples: $\Theta(p^2 \log p)$
 - Merging sorted sublists: $\Theta((n/p)\log p)$
 - Overall: $\Theta((n/p)(\log(n/p) + \log p) + p^2 \log p)$
- Communications
 - Gather samples pivots: $\Theta(p^2)$
 - Broadcast $p-1$ pivots: $\Theta(p \log p)$
 - All-to-all exchange: $\Theta(n/p)$
 - Overall: $\Theta(n/p + p^2)$

Summary

- Three parallel algorithms based on quicksort
- Keeping list sizes balanced
 - Parallel quicksort: poor
 - Hyperquicksort: better
 - PSRS algorithm: excellent
- Average number of times each key moved:
 - Parallel quicksort and hyperquicksort: $\log p / 2$
 - PSRS algorithm: $(p-1)/p$

Analysis of Parallel Quicksort

- Execution time dictated by when last process completes
- Algorithm likely to do a poor job balancing number of elements sorted by each process
- Cannot expect pivot value to be true median
- Can choose a better pivot value

Sorting on Specific Networks

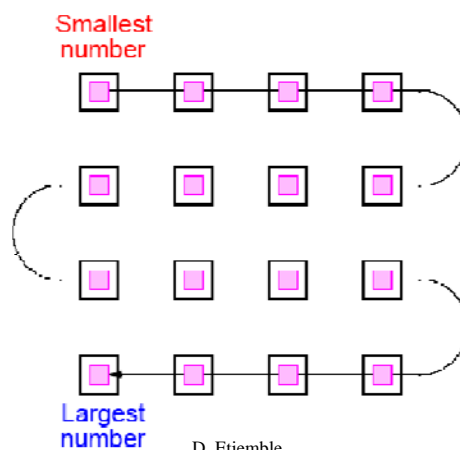
- Two network structures have received special attention: **mesh** and **hypercube**
Parallel computers have been built with these networks.
- However, it is of less interest nowadays because networks got faster and clusters became a viable option.
- Besides, network architecture is often hidden from the user.
- MPI provides libraries for mapping algorithms onto meshes, and one can always use a mesh or hypercube algorithm even if the underlying architecture is not one of them.

l@lytech4 - Option
parallélisme - 2014

D. Etiemble

Two-Dimensional Sorting on a Mesh

The layout of a sorted sequence on a mesh could be row by row or *snakelike*:

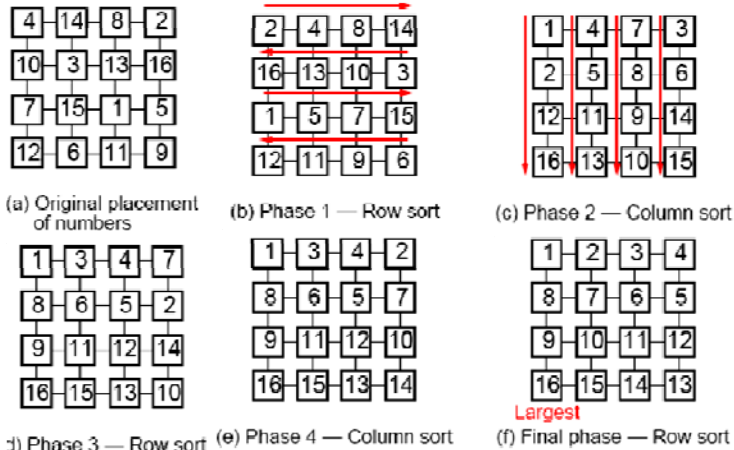


l@lytech4 - Option
parallélisme - 2014

D. Etiemble

Shearsort

Alternate row and column sorting until list is fully sorted.
 Alternate row directions to get snake-like sorting:



(a) Original placement of numbers (b) Phase 1 — Row sort (c) Phase 2 — Column sort
 (d) Phase 3 — Row sort (e) Phase 4 — Column sort (f) Final phase — Row sort
 Polytech4 - Option parallélisme - 2014 D. Etiemble 65

Shearsort – Time complexity

On a $n \times n$ Mesh, it takes $2 \log n$ phases to sort n^2 numbers.
 Therefore:

$$T_{par}^{shearsort} = O(n \log n) \quad \text{on a } n \times n \text{ mesh}$$

Since sorting n^2 numbers sequentially takes $T_{seq} = O(n^2 \log n)$;

$$Speedup_{shearsort} = \frac{T_{seq}}{T_{par}} = O(n) \quad (\text{for } P = n^2)$$

$$\text{However, efficiency} = \frac{1}{n}$$

Polytech4 - Option parallélisme - 2014 D. Etiemble 66

Rank Sort

Number of elements that are smaller than each *selected element* is counted. This count provides the position of the selected number, its “rank” in the sorted list.

- First $a[0]$ is read and compared with each of the other numbers, $a[1] \dots a[n-1]$, recording the number of elements less than $a[0]$.

Suppose this number is x . This is the index of $a[0]$ in the final sorted list.

- The number $a[0]$ is copied into the final sorted list $b[0] \dots b[n-1]$, at location $b[x]$. Actions repeated with the other numbers.

Overall sequential time complexity of rank sort: $T_{seq} = O(n^2)$
(not a good sequential sorting algorithm!)

Sequential code

```
for (i = 0; i < n; i++) {      /* for each number */
    x = 0;
    for (j = 0; j < n; j++)    /* count number less than it */
        if (a[i] > a[j]) x++;

    b[x] = a[i];              /* copy number into correct place */
}
```

**This code needs to be fixed if duplicates exist in the sequence.*

sequential time complexity of rank sort: $T_{seq} = O(n^2)$

Parallel Rank Sort (P=n)

One number is assigned to each processor.
 P_i finds the final index of $a[i]$ in $O(n)$ steps.

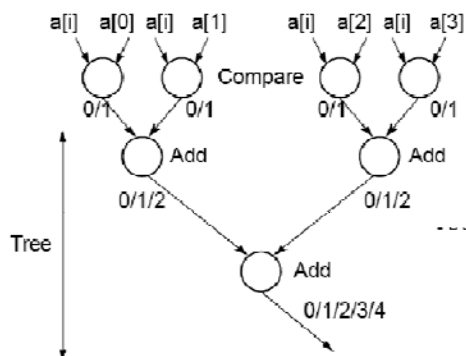
```
forall (i = 0; i < n; i++) { /* for each no. in parallel*/
    x = 0;
    for (j = 0; j < n; j++) /* count number less than it */
        if (a[i] > a[j]) x++;
    b[x] = a[i]; /* copy no. into correct place */
}
```

Parallel time complexity, $O(n)$, as good as any sorting algorithm so far. Can do even better if we have more processors.

Parallel time complexity: $T_{par} = O(n)$ (for $P=n$)

Parallel Rank Sort with $P = n^2$

Use n processors to find the rank of one element. The final count, i.e. rank of $a[i]$ can be obtained using a binary addition operation (global sum \rightarrow MPI_Reduce())

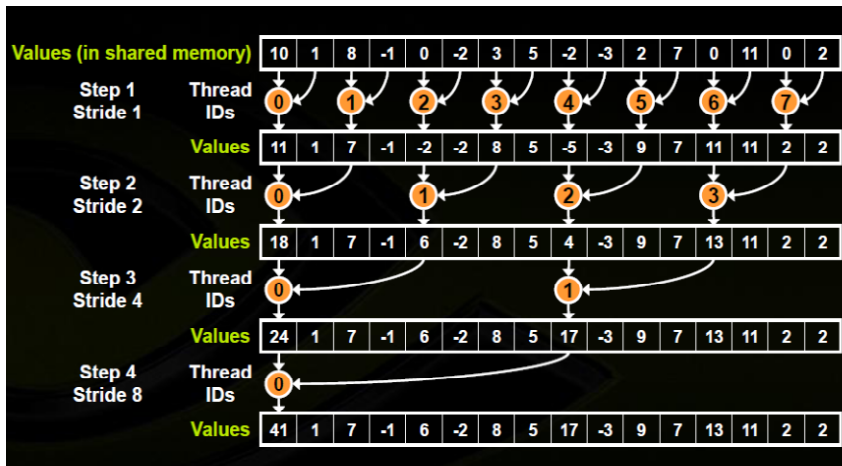


Time complexity
 (for $P=n^2$):
 $T_{par} = O(\log n)$
 Can we do it in $O(1)$?

Réduction sur GPU



Entrelacement



Polytech4 - Option
parallélisme - 2014

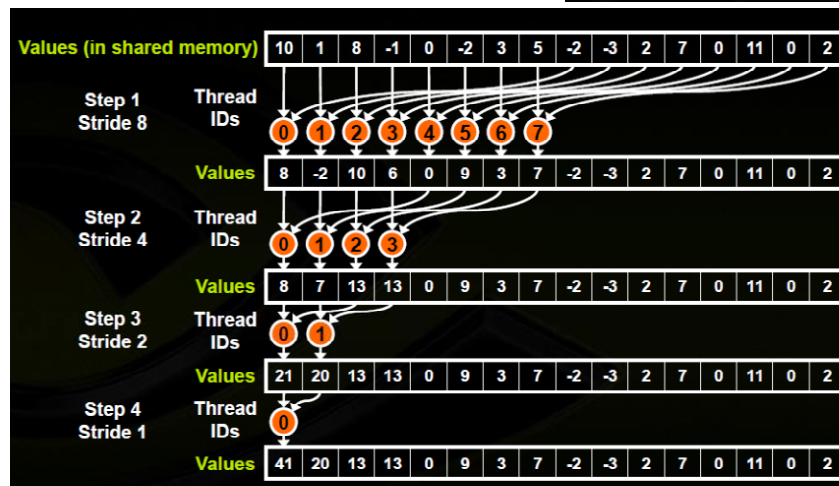
D. Etiemble

71

Réduction sur GPU



Continu



Polytech4 - Option
parallélisme - 2014

D. Etiemble

72

Algorithme Scan (GPU)

$$\text{scan}(A) = [l, a_0, (a_0 \oplus a_1), \dots, (a_0 \oplus a_1 \oplus \dots \oplus a_{n-2})]$$

