

Parallel Partition Revisited

Leonor Frias and Jordi Petit

Dep. de Llenguatges i Sistemes Informàtics, Universitat Politècnica de Catalunya

WEA 2008

Overview

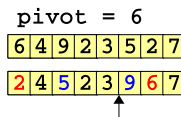
Partitioning an array with respect to a pivot is a basic building block of key algorithms such as *quicksort* and *quickselect*.

Overview

Partitioning an array with respect to a pivot is a basic building block of key algorithms such as *quicksort* and *quickselect*.

Given a pivot, rearrangement s.t for some splitting position s ,

- elements at the left of s are \leq pivot
- elements at the right of s are \geq pivot

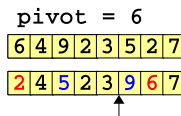


Overview

Partitioning an array with respect to a pivot is a basic building block of key algorithms such as *quicksort* and *quickselect*.

Given a pivot, rearrangement s.t for some splitting position s ,

- elements at the left of s are \leq pivot
- elements at the right of s are \geq pivot



Sequential cost:

- n comparisons
- m swaps

Overview

Partitioning an array with respect to a pivot is a basic building block of key algorithms such as *quicksort* and *quickselect*.

Overview

Partitioning an array with respect to a pivot is a basic building block of key algorithms such as *quicksort* and *quickselect*.

Nowadays, multi-core computers are ubiquitous.

Overview

Partitioning an array with respect to a pivot is a basic building block of key algorithms such as *quicksort* and *quickselect*.

Nowadays, **multi-core computers** are ubiquitous.

Several suitable parallel partitioning algorithms for these architectures exists.

- Algorithms by Francis and Pannan, Tsigas and Zang and MCSTL.

Overview

Partitioning an array with respect to a pivot is a basic building block of key algorithms such as *quicksort* and *quickselect*.

Nowadays, **multi-core computers** are ubiquitous.

Several suitable parallel partitioning algorithms for these architectures exists.

HOWEVER, they perform **more operations** than the sequential algorithm.

Overview

Partitioning an array with respect to a pivot is a basic building block of key algorithms such as *quicksort* and *quickselect*.

Nowadays, **multi-core computers** are ubiquitous.

Several suitable parallel partitioning algorithms for these architectures exists.

HOWEVER, they perform **more operations** than the sequential algorithm.

IN THIS PAPER:

- Show how to modify these algorithms so that they achieve **a minimal number of comparisons**.
- Provide implementations and a detailed **experimental** comparison.

Outline

- 1 Previous work
- 2 Algorithm
- 3 Experiments
- 4 Conclusions
- 5 References

Partitioning in parallel: overview

General pattern

- 1 Sequential **setup** of each processor's work
- 2 Parallel **main phase** in which most of the partitioning is done
- 3 **Cleanup** phase

p processors used to partition an array of n elements ($p \ll n$).

Partitioning in parallel: STRIDED (1)

STRIDED algorithm by Francis and Pannan.

- 1 Setup: Division into p pieces, elements in a piece with stride p

`pivot = 40, p = 4`

12	91	15	3	71	86	25	47	30	35	64	19	2	39	85	17	53	46	10	27	95	54	5	59
----	----	----	---	----	----	----	----	----	----	----	----	---	----	----	----	----	----	----	----	----	----	---	----

Partitioning in parallel: STRIDED (1)

STRIDED algorithm by Francis and Pannan.

- 1 Setup: Division into p pieces, elements in a piece with stride p

`pivot = 40, p = 4`

12	91	15	3	71	86	25	47	30	35	64	19	2	39	85	17	53	46	10	27	95	54	5	59
----	----	----	---	----	----	----	----	----	----	----	----	---	----	----	----	----	----	----	----	----	----	---	----

- 2 Main phase: Sequential partitioning in each piece

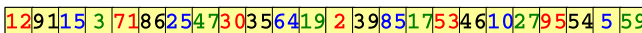
12	39	15	3	2	35	25	27	30	86	5	19	71	19	11	10	17	53	46	85	47	95	54	64	59
))	(())		((

Partitioning in parallel: STRIDED (1)

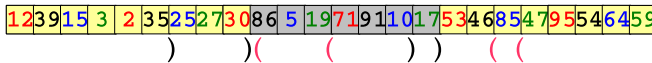
STRIDED algorithm by Francis and Pannan.

- 1 Setup: Division into p pieces, elements in a piece with stride p

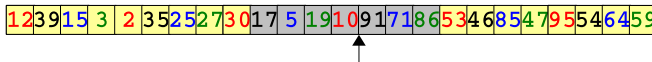
`pivot = 40, p = 4`



- 2 Main phase: Sequential partitioning in each piece



- 3 Cleanup: Sequential partitioning in the not correctly partitioned range



Partitioning in parallel: STRIDED (2)

STRIDED Analysis:

- Main phase: $\Theta(n/p)$ parallel time
- Cleanup phase: $O(1)$ expected but can be $\Theta(n)$

1	2	9	1	1	5	7	1	3	8	6	2	5	4	7	3	0	6	4	3	5	5	9	2	3	9	8	6	8	7	1	3	4	6	5	5	7	3	5	5	4	2	2	5	9
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Partitioning in parallel: STRIDED (2)

STRIDED Analysis:

- Main phase: $\Theta(n/p)$ parallel time
- Cleanup phase: $O(1)$ expected but can be $\Theta(n)$

12	9	1	1	5	7	1	3	8	6	2	5	4	7	3	0	6	4	3	5	9	2	3	9	8	6	8	7	1	3	4	6	5	5	7	3	5	5	4	2	2	5	9
----	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Partitioning in parallel: STRIDED (2)

STRIDED Analysis:

- Main phase: $\Theta(n/p)$ parallel time
- Cleanup phase: $O(1)$ expected but can be $\Theta(n)$

BESIDES, it has poor cache locality.

Partitioning in parallel: BLOCKED

We can generalize STRIDED to blocks to improve cache locality.

`pivot = 40, p = 4, b = 3`

129115 3 7186 254730 356419 2 3985 175346 102795 54 5 59

123915 17 3 86 252730 5 3519 2 9185 715346 104795 546459

) () () () (

123915 17 3 10 252730 5 3519 2 9185 715346 864795 546459



If $b = 1$, BLOCKED is equal to STRIDED.

Partitioning in parallel: F&A (1)

Processors take elements from both ends of the array as they are needed.

Fetch-and-add instructions are used to *acquire* the elements.

Blocks of elements are used to avoid too much synchronization.

References:

- PRAM model: Heidelberg *et al.*
- *real* machines: Tsigas and Zhang and MCSTL library

Partitioning in parallel: F&A (2)

- ① Setup: Each processor takes one left block and one right block

1	2	9	1	1	5	3	7	1	8	6	2	5	4	7	3	0	3	5	6	4	1	9	2	3	9	8	5	1	7	5	3	4	6	1	0	2	7	9	5	5	4	5	5	9
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Partitioning in parallel: F&A (2)

- 1 Setup: Each processor takes one left block and one right block

12	9	1	15	3	7	1	8	6	2	5	4	7	3	0	3	5	6	4	1	9	2	3	9	8	5	1	7	5	3	4	6	1	0	2	7	9	5	5	4	5	5	9
----	---	---	----	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

- 2 Main phase: Sequential partitioning in sequence made by left block + right block. When one block border is reached and so *neutralized*, another block is acquired.

1	2	5	15	3	7	1	8	6	2	5	4	7	3	0	3	5	6	4	1	9	2	3	9	8	5	1	7	5	3	4	6	1	0	2	7	9	5	5	4	9	1	5	9
---	---	---	----	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Partitioning in parallel: F&A (2)

- 1 Setup: Each processor takes one left block and one right block

12	9	15	3	7	18	6	25	4	7	3	0	3	5	6	4	1	9	2	3	9	8	5	1	7	5	3	4	6	1	0	2	7	9	5	5	4	5	5	9
----	---	----	---	---	----	---	----	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

- 2 Main phase: Sequential partitioning in sequence made by left block + right block. When one block border is reached and so *neutralized*, another block is acquired.

12	5	15	3	17	8	6	25	3	9	3	0	3	5	2	7	1	9	2	4	7	8	5	7	1	5	3	4	6	1	0	6	4	9	5	5	4	9	1	5	9
----	---	----	---	----	---	---	----	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Partitioning in parallel: F&A (2)

- 1 Setup: Each processor takes one left block and one right block

129115371862547303564192398517534610279554559

- 2 Main phase: Sequential partitioning in sequence made by left block + right block. When one block border is reached and so *neutralized*, another block is acquired.

- 3 Cleanup: At most p blocks remain not completely partitioned (unneutralized). The unpartitioned elements must be moved to the *middle*

- Tsigas and Zhang do it sequentially.
- MCSTL moves the blocks in parallel and applies recursively parallel partition to this range.

125153271925393035171086648571534624795549159



Partitioning in parallel: F&A (3)

F&A Analysis:

- Main phase: $\Theta(n/p)$ parallel time
- Cleanup phase:
 - Tsigas and Zhang: $O(bp)$
 - MCSTL: $\Theta(b \log p)$

New Parallel Cleanup Phase

Existing algorithms disregard part of the work done in the main parallel phase when cleaning up.

New Parallel Cleanup Phase

Existing algorithms disregard part of the work done in the main parallel phase when cleaning up.

We present a new cleanup algorithm.

New Parallel Cleanup Phase

Existing algorithms disregard part of the work done in the main parallel phase when cleaning up.

We present a new cleanup algorithm.

- It avoids redundant comparisons.
- The elements are swapped fully in parallel.

We apply it on the top of STRIDED, BLOCKED and F&A algorithms.

Terminology (1)

Our algorithm is described in terms of

- **Subarray**
- **Frontier**: Defines two parts (left and right) in a subarray
- **Misplaced element**

Their realization depends on the algorithm used in the main parallel phase.

Terminology (1)

Our algorithm is described in terms of

- **Subarray**
- **Frontier**: Defines two parts (left and right) in a subarray
- **Misplaced element**

Their realization depends on the algorithm used in the main parallel phase.

m : total number of misplaced elements

M : total number of subarrays that may have misplaced elements.

Terminology (1)

Our algorithm is described in terms of

- **Subarray**
- **Frontier**: Defines two parts (left and right) in a subarray
- **Misplaced element**

Their realization depends on the algorithm used in the main parallel phase.

m : total number of misplaced elements

M : total number of subarrays that may have misplaced elements.

Terminology for BLOCKED



- **Subarray**: each of the p pieces.
- **Frontier**: position that would occupy the pivot after partitioning the array.
- **Misplaced elements**: as in the sequential algorithm.
- $M \leq p$

Terminology for F&A

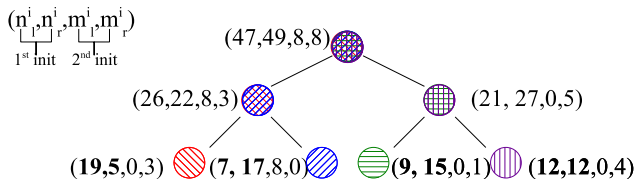
We deal separately and analogously with left and right blocks.



- **Subarray**: one block.
- **Frontier**: separates the processed part in a block from the unprocessed part.
- **Misplaced elements**: unprocessed elements not in the *middle* and processed elements that are in the *middle*.
- $M \leq 2p$ (p unneutralized blocks which could be all misplaced and almost full)

Data Structure

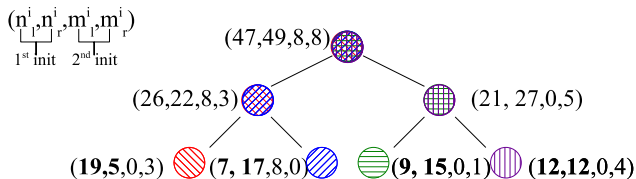
Shared arrayed binary tree with M leaves.



- **Leaves:** information on the subarrays
- **Internal nodes:** accumulate children information

Data Structure

Shared arrayed binary **tree** with M leaves.



- **Leaves:** information on the subarrays
- **Internal nodes:** accumulate children information

Use: deciding pairs of elements to be swapped without doing new comparisons

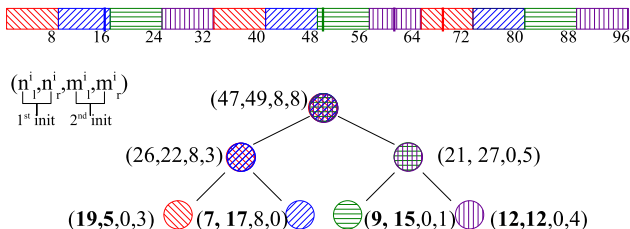
Algorithm (1)

Tree initialization

- 1 First initialization of the leaves: Computation of $n_{l,r}^i$.
- 2 First initialization of the non-leaves: Computation of $n_{l,r}^j, v$.
- 3 Second initialization of the leaves: Computation of $m_{l,r}^i$.
- 4 Second initialization of the non-leaves: Computation of $m_{l,r}^j$.

Tree initialization for BLOCKED

Computation of $n_{l,r}^i$: trivially
(the layout is deterministic, b and i are known)

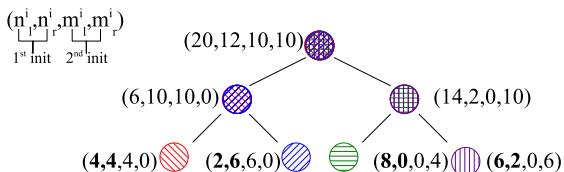


Tree initialization for F&A

Computation of $n_{l,r}^i$: Trivially once the subarrays are known.

Determination of the subarrays:

- The unneutralized blocks are known after the parallel phase.
- To locate the misplaced neutralized blocks, the unneutralized blocks are sorted by address and then, traversed.



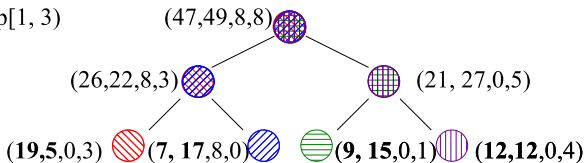
Algorithm (2)

Parallel swapping

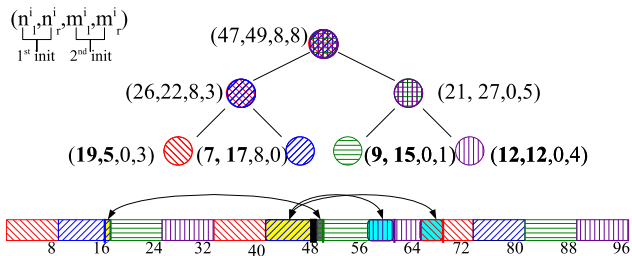
Independent of the algorithm in the main parallel phase.

The misplaced elements to swap are **divided equally** among the processors.

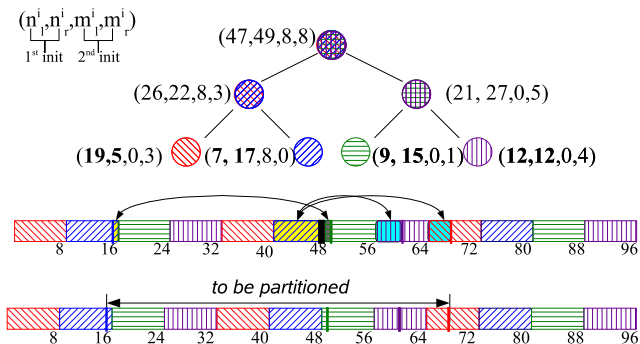
swap[1, 3)



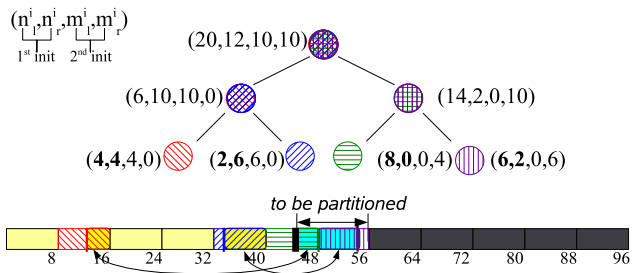
Parallel swapping for BLOCKED



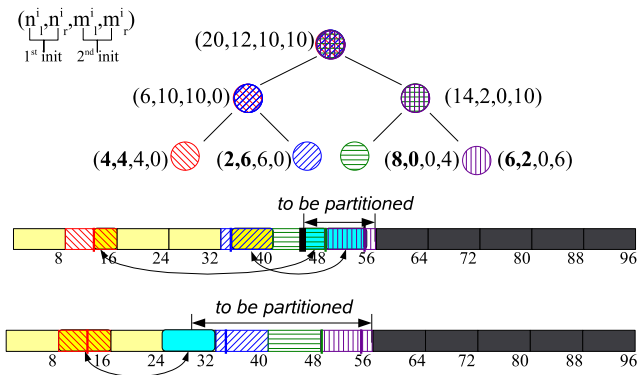
Parallel swapping for BLOCKED



Parallel swapping for F&A



Parallel swapping for F&A



Algorithm (3)

Completion

BLOCKED : The array is already partitioned.

F&A : The array is partitioned except for the elements in the *middle* (not yet processed).

- Apply **recursively parallel partitioning** in the *middle*.

We provide a **better cost bound** making recursion on *b*
($b \leftarrow b/2$ for $\log p$ times) instead of p .

Analysis: comparisons & swaps

BLOCKED				
	comparisons		swaps	
	original	tree	original	tree
main	n		$\leq n/2$	
cleanup	$v_{max} - v_{min}$	0	$m/2$	$m/2$
total	$n + v_{max} - v_{min}$	n	$\leq \frac{n+m}{2}$	$\leq \frac{n+m}{2}$

F&A				
	comparisons		swaps	
	original	tree	original	tree
main	$n - V $		$\leq \frac{n- V }{2}$	
cleanup	$\leq 2bp$	$ V $	$\leq 2bp$	$\leq m/2 + V $
total	$\leq n + 2bp$	n	$\leq \frac{n- V }{2} + 2bp$	$\leq \frac{n+m}{2} + V $

Analysis: worst-case time

BLOCKED		
	parallel time	
	original	tree
main	$\Theta(n/p)$	
cleanup	$\Theta(v_{\max} - v_{\min})$	$\Theta(m/p + \log p)$
total	$\Theta(n)$	$\Theta(n/p + \log p)$

F&A		
	parallel time	
	original	tree
main	$\Theta(n/p)$	
cleanup	$\Theta(b \log p)$	$\Theta(\log^2 p + b)^1$
total	$\Theta(n/p + b \log p)$	$\Theta(n/p + \log^2 p)$

¹better provided that $\log p \leq b$

Implementation

Algorithms: STRIDED, BLOCKED, F&A (MCSTL & own)

- With original cleanup
- With our cleanup

Languages: C++, OpenMP

STL partition interface.

Setup

Machine

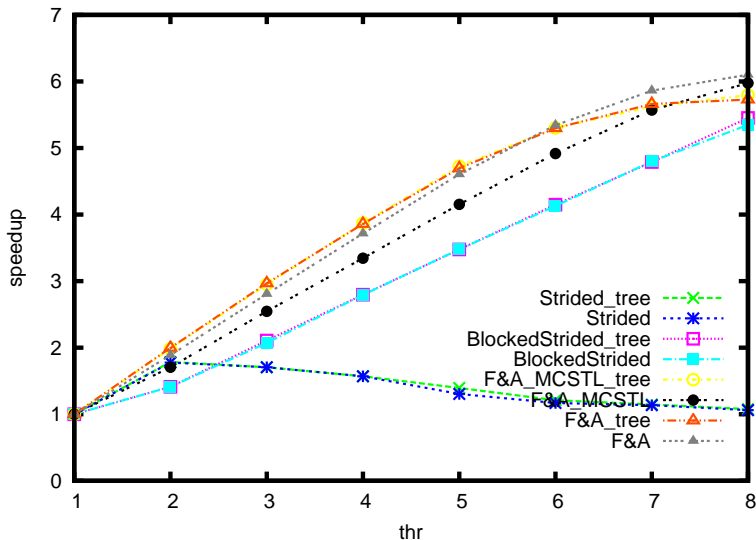
- 4 GB of main memory
- 2 sockets x Intel Xeon quad-core processor at 1.66 GHz with a shared L2 cache of 4 MB shared among two cores

Compiler: GCC 4.2.0, -O3 optimization flag.

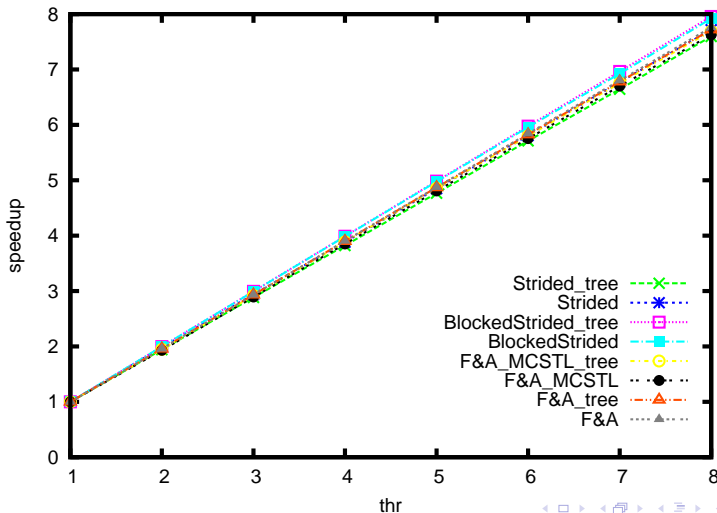
Measurements:

- 100 repetitions
- Speedups with respect to the sequential algorithm in the STL

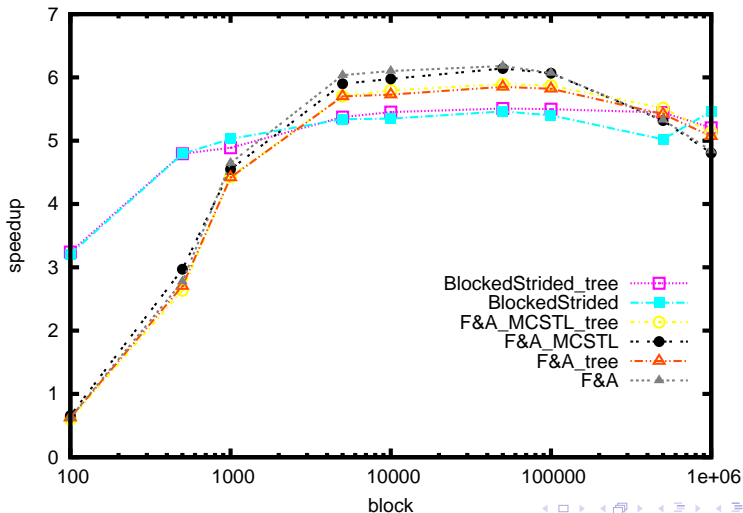
Parallel partition speedup, $n = 10^8$ and $b = 10^4$



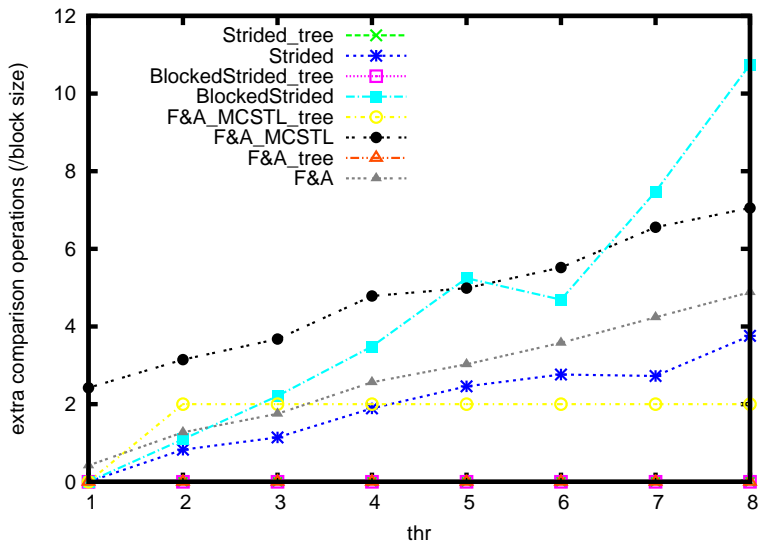
Parallel partition speedup for costly $<$, $n = 10^8$ and $b = 10^4$



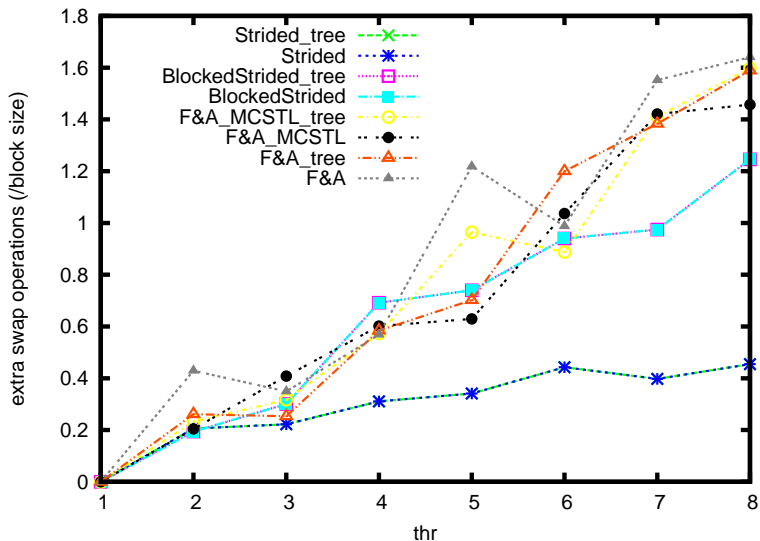
Parallel partition with varying block size, $n = 10^8$ and $\text{num_threads} = 8$



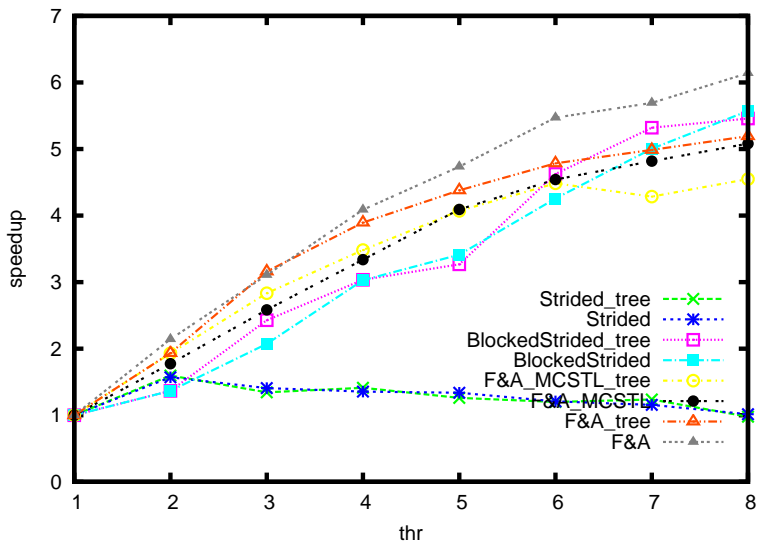
Number of extra comparisons, $n = 10^8$ and $b = 10^4$



Number of extra swaps, $n = 10^8$ and $b = 10^4$



Parallel quickselect speedup, $n = 10^8$ and $b = 10^4$



Conclusions (1)

We have presented, implemented and evaluated several parallel **partitioning algorithms** suitable for **multi-core** architectures.

Conclusions (1)

We have presented, implemented and evaluated several parallel **partitioning algorithms** suitable for **multi-core** architectures.

Algorithmic contributions:

- Novel cleanup algorithm NOT disregarding any **comparisons** made in the parallel phase.
- Applied to STRIDED, BLOCKED and F&A partitioning algorithms.
 - **STRIDED and BLOCKED** : worst-case parallel time from $\Theta(n)$ to $\Theta(n/p + \log p)$.
- Better cost bound for **F&A** changing recursion parameters.

Conclusions (2)

Implementation contributions: carefully designed implementations following **STL partition** specifications.

Detailed **experimental** comparison. Conclusions:

- Algorithm of choice: **F&A** (**ours** was best).
- Benefits **in practice** of the cleanup algorithm very **limited**.
- I/O limits performance as the number of threads increases.

Thank you for your attention

More information:

<http://www.lsi.upc.edu/~lfrias>.

References



R. S. Francis and L. J. H. Pannan.

A parallel partition for enhanced parallel quicksort.
Parallel Computing, 18(5):543–550, 1992.



P. Heidelberger, A. Norton, and John T. Robinson.

Parallel quicksort using fetch-and-add.
IEEE Trans. Comput., 39(1):133–138, 1990.



J. Singler, P. Sanders, and F. Putze.

The Multi-Core Standard Template Library.

In *Euro-Par 2007: Parallel Processing*, volume 4641 of *Lecture Notes in Computer Science*, pages 682–694, Rennes, France, 2007. Springer Verlag.



P. Tsigas and Y. Zhang.

A simple, fast parallel implementation of quicksort and its performance evaluation on SUN enterprise 10000.

In 11th Euromicro Workshop on Parallel, Distributed and Network-Based Processing (PDP 2003), pages 372–381, 2003.