

# Parallel Partition Revisited

Leonor Frias<sup>\*,\*\*</sup> and Jordi Petit<sup>\*\*\*</sup>

Departament de Llenguatges i Sistemes Informàtics  
Universitat Politècnica de Catalunya  
lfrias@lsi.upc.edu

**Abstract.** In this paper we consider parallel algorithms to partition an array with respect to a pivot. We focus on implementations for current widely available multi-core architectures. After reviewing existing algorithms, we propose a modification to obtain the minimal number of comparisons. We have implemented these algorithms and drawn an experimental comparison.

## 1 Introduction

The partitioning of an array is a basic building block of many key algorithms, as quicksort and quickselect. Partitioning an array with respect to a pivot  $x$  consists of rearranging its elements such that, for some splitting position  $s$ , all elements at the left of  $s$  are smaller than  $x$ , and all other elements are greater or equal than  $x$ . It is well known that an array of  $n$  elements can be partitioned sequentially and in-place using exactly  $n$  comparisons and  $m$  swaps, where  $m$  is the number of greater elements than  $x$  whose original position is smaller than  $s$ .

In this paper we consider the problem of partitioning an array in parallel, focusing on current widely available multi-core architectures.

Several algorithms have been proposed to partitioning in parallel [1,2,3,4,5]. In this paper, we consider a simple algorithm by Francis and Pannan [2], a fetch-and-add based algorithm by Tsigas and Zhang [3] and a variation of the former in the MCSTL library [5]. These algorithms, which we survey in Sect. 2, seem suitable for a practical multi-core implementation. However, in order to avoid too much synchronization, they perform more than  $n$  comparisons and  $m$  swaps. Though very different in nature, they can be divided into three main phases: a) A sequential setup of each processor's work, b) a parallel main phase in which most of the partitioning is done, and c) a cleanup phase, which is usually sequential.

In this paper we show that these algorithms disregard part of the work done in the main parallel phase when cleaning up. In order to overcome this drawback, we propose an alternative parallel cleanup phase that uses the whole comparison information of the parallel phase. A small static order-statistics tree is used to

---

\* Supported by grant number 2005FI 00856 of the *Agència de Gestió d'Ajuts Universitaris i de Recerca* with funds of the European Social Fund.

\*\* Partially supported by Spanish project ALINEX (ref. TIN2005-05446).

\*\*\* Partially supported by FET proactive integrated project 15964 (AEOLUS) and by Spanish project FORMALISM (ref. TIN2007-66523).

efficiently locate the elements to be swapped and to swap them in parallel. With this new method, we obtain scalable parallel partitioning algorithms that achieve an optimal number of comparisons. We provide a detailed analysis.

We have implemented and evaluated all these algorithms, both with their original cleanup and with our cleanup. Besides, the implementation is provided according to the specification of the `partition` function of the Standard Template Library (STL) of the C++ programming language [6]. Previously, only F&A implementation was available in the MCSTL library [5]. Our goal is to get a comparison of their behavior when executed on a currently inexpensive widely available parallel machine, namely a machine with two quad-core processors.

The paper is organized as follows. In Sect. 2, we present the considered algorithms. Then, in Sect. 3, we present our cleanup algorithm. Then, we present our implementation of the previous algorithms and the experimental results in Sect. 4 and 5 respectively. We sum up the conclusions of this work in Sect. 6.

## 2 Previous Work and a Variant

In this section, we present an overview of the partitioning algorithms we consider in this paper. In the following, the input consist of an array of  $n$  elements and a pivot.  $p$  processors are available and we assume  $p \ll n$ . Besides, we disregard some details as rounding issues for the sake of simplicity.

**Strided Algorithm.** The STRIDED algorithm by Francis and Pannan [2] works as follows:

1. **Setup:** The input is (conceptually) divided into  $p$  pieces of size  $n/p$ . The pieces are not made of consecutive elements, but one of every  $p$  elements instead. That is, the  $i$ -th piece is made up of elements  $i, i + p, i + 2p, \dots$ .
2. **Main phase:** Each processor  $i$ , in parallel, gets a piece, applies sequential partitioning on it, and returns its splitting position  $v_i$ .
3. **Cleanup:** Let  $v_{min} = \min\{v_i : 1 \leq i \leq p\}$  and  $v_{max} = \max\{v_i : 1 \leq i \leq p\}$ . It holds that all the elements at the left of  $v_{min}$  and at the right of  $v_{max}$  are already well placed with respect to the pivot. In order to complete the partition, sequential partition is applied to the range  $(v_{min}, v_{max})$ .

The main phase takes  $\Theta(n/p)$  parallel time. For random inputs, the cleanup phase is expected to take constant time. However, [2] did not state that in the worst-case it takes  $\Theta(n)$  time and thus, there is no speedup. E.g. If the pieces are made exclusively of either smaller or greater elements than the pivot and these are alternated, then,  $v_{min} \leq p$  and  $v_{max} \geq n - p$ , and  $|(v_{min}, v_{max})| = \Theta(n)$ .

**Blocked Algorithm.** Accessing elements with stride  $p$  as in STRIDED, can provoke a high cache miss ratio. We propose BLOCKED to overcome this problem. It uses blocks of  $b$  elements instead of individual elements. Each block in the piece is separated by stride  $p$  blocks. If  $b = 1$ , BLOCKED is equal to STRIDED.

**F&A Algorithms.** Heidelberg *et al.* [1] proposed a parallel partitioning algorithm in the PRAM model in which elements from both ends of the array are taken using fetch-and-add instructions. Fetch-and-add instructions (atomically increment a variable and return its original value) were introduced in [7] and are useful, for instance, to implement synchronization and mutual exclusion.

In a first approach, exactly one element is taken at a time and so, at the end of the parallel phase, the array is already partitioned. In this case,  $n$  fetch-and-add operations are used. In a second approach, the algorithm is generalized to blocks: a block of  $b$  elements is acquired at each fetch-and-add instruction. So, the number of fetch-and-add instructions is  $n/b$ . However, in this case, some sequential cleanup remains to be applied after the parallel phase.

Later, Tsigas and Zhang [3] presented a variant of the second approach for multiprocessors. More recently, a further variant has been included in the MC-STL library [5]. In the latter, the cleanup phase is partially done in parallel.

Let us now briefly describe these F&A algorithms:

1. **Setup:** Each processor takes two blocks, one from each end of the array. Namely, one left block and one right block.
2. **Main phase:** While there are blocks, each processor applies the so-called *neutralize* method to its two blocks. The neutralize method consists on applying the sequential partitioning algorithm to the array made by (conceptually) concatenating the right block to the left. However, the left and right pointers to the current elements cannot cross the borders of a block. When a left (right) block is completely processed (i.e. neutralized), a fresh left (right) block is acquired and processed.
3. **Cleanup:** At this point, at most  $p$  blocks remain unneutralized. Each author presents a different cleanup algorithm:
  - In [3], while unneutralized blocks remain, one block is taken from each end and neutralization is applied to them. Then, the unneutralized blocks are placed between the neutralized blocks. At most  $p$  blocks need to be swapped and this is done sequentially. Finally, sequential partition is applied to the range of blocks with unprocessed elements.
  - In [5], all unneutralized blocks are placed between the neutralized blocks. Then, the parallel partitioning algorithm is applied recursively to this range. The number of processors is divided by two in each call until there is only one processor or block. Finally, the remaining range is partitioned sequentially.

The main parallel phase takes  $\Theta(n/p)$  parallel time. The cleanup phase takes  $\Theta(bp)$  sequential time in [3]. Rather, in [5], it takes  $\Theta(b \log p)$  parallel time.

### 3 The New Parallel Cleanup Phase

In this section, we present our cleanup algorithm. It avoids extra comparisons and swaps the elements fully in parallel. We have applied it on the top of STRIDED, BLOCKED and F&A. First, we introduce the terminology. Then, we present the data structure on which the algorithm relies. It follows the cleanup

algorithm itself. Finally, we analyze the resulting STRIDED, BLOCKED and F&A algorithms.

**Terminology.** In the following, we shall use the following terms to describe our algorithm. A *subarray* is the basic unit of our algorithm and data structure. The *splitting position*  $v$  of an array is the position that would occupy the pivot after partitioning. A *frontier* separates a subarray in two consecutive parts that have different properties. A *misplaced element* is an element that must be moved by our algorithm. We denote by  $m$  the total number of misplaced elements and by  $M$  the total number of subarrays that may have misplaced elements.

*The Case of BLOCKED.* In this case, subarrays correspond with exactly one of the  $p$  pieces. Moreover,  $v$  can be easily known after the parallel phase. The frontier of a subarray corresponds with the position that would occupy the pivot after partitioning this subarray. Thus, a frontier defines a left and a right part. A misplaced element corresponds either to an element smaller than the pivot that is on the right of  $v$  (misplaced on the right) or to an element greater than the pivot that is on the left of  $v$  (misplaced on the left). The total number of subarrays that may have misplaced elements ( $M$ ) is at most  $p$ .

*The Case of F&A.* In this case, a subarray corresponds to one block. The frontier separates a processed part from an unprocessed part. The processed part of left blocks is the left part and the processed part of right blocks is the right part. Though  $v$  is unknown after the parallel phase, it holds that  $v$  is in some range  $V = [v_{\text{beg}}, v_{\text{end}}]$ . A misplaced element corresponds either with a processed element that is in  $V$  or with an unprocessed element that is not in  $V$ .

We consider now the array of elements made by left blocks and the array of elements made by right blocks. We denote by  $\beta$  the global border that separates the left and right blocks (i.e. the point where no more blocks could be obtained).

In left blocks, a misplaced element on the left is an unprocessed element in a position smaller than  $v_{\text{beg}}$  and a misplaced element on the right is a smaller element than the pivot in a position greater or equal than  $v_{\text{beg}}$ . Let  $\mu_l$  be the total number of misplaced elements in left blocks and rank those misplaced elements starting to count from the leftmost towards the right. In right blocks, a misplaced element on the left is an element greater than the pivot in a position smaller or equal than  $v_{\text{end}}$  and a misplaced element on the right is an unprocessed element in a position greater than  $v_{\text{end}}$ . Let  $\mu_r$  be the total number of misplaced elements in right blocks and rank those misplaced elements starting to count from the rightmost towards the left. Note that  $m = \mu_l + \mu_r$  and that the total number of subarrays that may have misplaced elements  $M$  is at most  $2p$ , because there at most  $p$  blocks that may contain unprocessed elements and there at most  $p$  blocks that may contain misplaced processed elements.

### 3.1 The Data Structure

We use a complete binary tree with  $M$  leaves (or the next power of two if  $M$  is not a power of two) to know which pairs of elements must be swapped. This tree

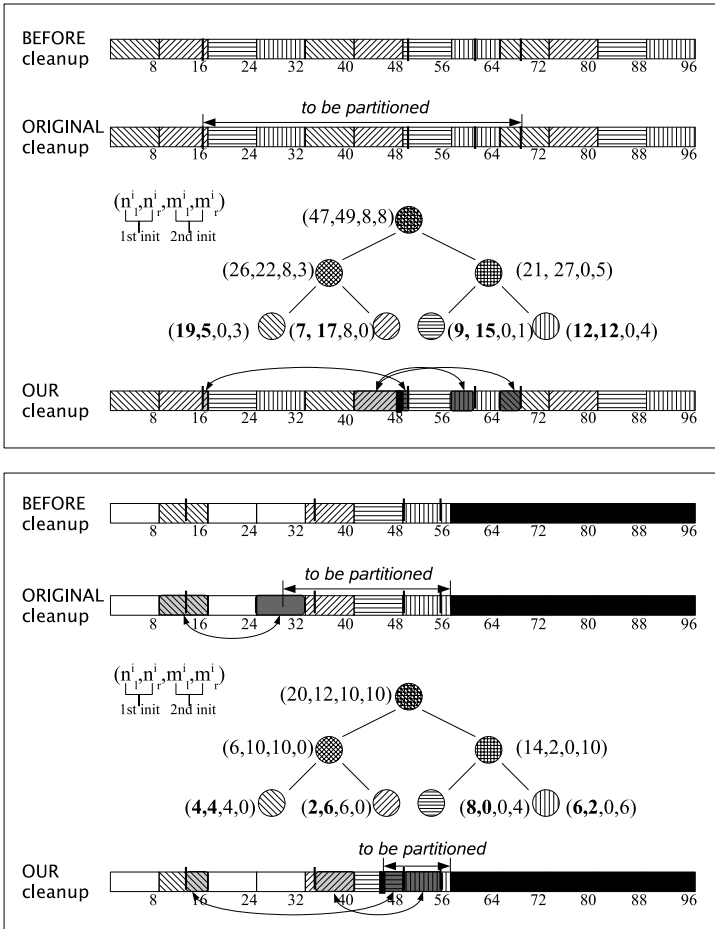


Fig. 1. Example of our data structure

is shared by all processors and is stored in an array (like a heap, which provides easy and efficient access to the nodes).

Each leaf stores information of the  $i$ -th subarray. Specifically, how many elements are misplaced to the left and to the right of its frontier ( $m_l^i$  and  $m_r^i$ ) and how many elements are in the left and in the right to its frontier ( $n_l^i$  and  $n_r^i$ ). The internal nodes accumulate the information of their children but do not add any new information. In particular, the root stores the information of the array made of all the subarrays in the leaves.

So, our tree data structure can be considered as a special kind of order-statistics tree in which the internal nodes have no information by themselves. An order-statistics tree (see e.g. [8, Sect. 14]) perform rank operations efficiently using the information of the size of the subtrees.

Figure 1 shows two instances of our tree data structure.

### 3.2 The Algorithm

*Tree Initialization Phase.* In this phase the tree is initialized. Specifically, two bottom-up traversals of the tree are needed. Only the first initialization of the leaves depends on the partition algorithm used in the main parallel phase.

1. **First initialization of the leaves.** In the case of BLOCKED, the leaves values  $n_l^i$  and  $n_r^i$  for each subarray  $i$  can be trivially computed during the parallel phase. In the case of F&A, the left (right) blocks that contain unprocessed elements can be easily known after the parallel phase. The left (right) blocks that contain misplaced elements but have already been processed can only be located between the left (right) unprocessed blocks. In order to locate the latter efficiently, we sort the unneutralized blocks with respect to the block position in the array. Then, we iterate on left (right) blocks (sequentially) to the left (right) of the border  $\beta$  until  $p$  neutralized blocks have been found or the leftmost (rightmost) unneutralized block has been reached.
2. **First initialization of the non-leaves.** Using a parallel reduce operation, each internal node computes its  $n_l^j$  and  $n_r^j$  values from its children. As a result, the root stores the number of left and right elements in the whole array. Thus, the splitting point  $v$  can be directly deduced.
3. **Second initialization of the leaves.** The leaves get the values  $m_l^i$  and  $m_r^i$  using  $n_l^i$ ,  $n_r^i$  and  $v$ . At this point, it may turn out that some subarrays have no misplaced elements. This does not disturb the correctness of our algorithm.
4. **Second initialization of the non-leaves.** The number of misplaced elements for the internal nodes are computed using a second parallel reduction operation on  $m_l^j$  and  $m_r^j$  fields.

*Parallel Swapping Phase.* In this phase, our tree data structure is queried so that the misplaced elements can be swapped in parallel and no comparisons are needed. This phase is independent of the specific partitioning algorithm.

The total number of misplaced elements is used to distribute the work equally among the processors. A range of ranks  $[r_i, s_i)$  of misplaced elements to swap is assigned to each processor. The elements are swapped in ascending rank. Specifically, the  $j$ -th misplaced left element is swapped with the  $j$ -th misplaced right element. To locate the first pair of elements to swap, respective rank queries are made to the tree. That is, a query is made for the  $r_i$  left misplaced element and another for the  $r_i$  right misplaced element. Misplaced elements are swapped as long as the rank  $s_i$  is not reached. If the rank  $s_i$  has not yet been reached but the current subarray has no more misplaced elements, the next subarray is fetched. Let  $c_i$  be the position in the tree corresponding to the current subarray. Then, the next subarray of left misplaced elements is in  $c_i + 1$  and the next subarray of right misplaced elements is in  $c_i - 1$ .

*Completion Phase.* This phase depends on the specific partitioning algorithm. In the case of BLOCKED, the whole array has already been partitioned and we are done. In the case of F&A, some unprocessed elements may remain. When this happens,  $V$  is not empty and includes exclusively all the unprocessed elements.

In order to obtain a valid partition, we apply  $\log p$  times our parallel partitioning algorithm recursively in  $V$  using blocks of half their original size until  $b$  elements or less remain. Sequential partition is applied to those remaining elements. Note that in each recursive call, the size of the problem is at most half the previous because at least  $p$  blocks have been fully processed.

### 3.3 Cost Analysis

**Theorem 1.** *BLOCKED and F&A perform exactly  $n$  comparisons when using our cleanup algorithm.*

*Proof.* The tree initialization and the parallel swapping phases perform no comparisons for both BLOCKED and F&A.

In the case of BLOCKED, the completion phase is empty. Therefore, no comparisons are performed during cleanup for BLOCKED and thus, comparisons are only performed during the main parallel phase, which are exactly  $n$ .

In the case of F&A, after the first main parallel phase  $n - |V|$  elements have been compared and  $|V|$  have remained unprocessed. In the next recursive step,  $V$  is the input. Besides, elements can only be compared during a certain parallel phase and at most once. All the elements must be eventually compared because our algorithm produces a valid partition. Thus, our cleanup algorithm makes exactly  $|V|$  comparisons, and  $n$  comparisons are needed as a whole.

**Lemma 1.** *The tree initialization phase takes  $\Theta(\log p)$  parallel time for BLOCKED and F&A.*

*Proof.* The algorithm-independent part takes  $\Theta(\log p)$  parallel time because all the work is done in parallel, and is dominated by the two parallel reductions, which can be performed in logarithmic parallel time [9].

In the case of BLOCKED, the algorithm-dependent part takes constant parallel time because each leaf can be initialized trivially and in parallel.

In the case of F&A, the algorithm-dependent part takes  $\Theta(\log p)$  parallel time, because  $2p$  elements are sorted and this takes  $\Theta(\log p)$  parallel time using  $p$  processors [10].

Thus, in both cases, the total cost is  $\Theta(\log p)$  parallel time.

**Lemma 2.** *The parallel swapping phase performs exactly  $m/2$  swaps and requires  $\Theta(m/p)$  parallel time. In the case of F&A, this parallel time is  $O(b)$ .*

*Proof.* There are  $m$  misplaced elements after the main parallel phase. The parallel swapping phase swaps pairs of misplaced elements so that their final position is not misplaced. Therefore,  $m/2$  swap operations are needed. Besides, the pairs are evenly divided among the  $p$  processors. Thus swapping all of them takes  $\Theta(m/p)$  parallel time. In the case of F&A,  $m \leq 2bp$ , thus parallel swapping takes  $O(b)$  parallel time.

**Theorem 2.** *The cleanup phase takes  $\Theta(m/p + \log p)$  parallel time for BLOCKED and the whole partition takes  $\Theta(n/p + \log p)$  parallel time.*

**Table 1.** Summary of costs for BLOCKED and F&A algorithms

BLOCKED						
	total comparisons		total swaps		parallel time	
	original	tree	original	tree	original	tree
main	$n$		$\leq n/2$		$\Theta(n/p)$	
cleanup	$v_{max} - v_{min}$	$0$	$m/2$	$m/2$	$\Theta(v_{max} - v_{min})$	$\Theta(m/p + \log p)$
total	$n + v_{max} - v_{min}$	$n$	$\leq \frac{n+m}{2} \leq \frac{n+m}{2}$		$O(n)$	$\Theta(n/p + \log p)$

F&A						
	comparisons		swaps		parallel time	
	original	tree	original	tree	original	tree
main	$n -  V $		$\leq \frac{n- V }{2}$		$\Theta(n/p)$	
cleanup	$\leq 2bp$	$ V $	$\leq 2bp$	$\leq m/2 +  V $	$\Theta(b \log p)$	$\Theta(\log^2 p + b)$
total	$\leq n + 2bp$	$n$	$\leq \frac{n- V }{2} + 2bp \leq \frac{n+m}{2} +  V $		$\Theta(n/p + b \log p)$	$\Theta(n/p + \log^2 p)$

*Proof.* From Lemmas 1 and 2 follows that the cleanup phase takes  $\Theta(m/p + \log p)$  parallel time for BLOCKED. Given that  $m = O(n)$ , the whole BLOCKED algorithm takes  $\Theta(n/p + \log p)$  parallel time in the average and in the worst-case.

**Theorem 3.** Consider  $p \leq b$ . The cleanup phase takes  $\Theta(\log^2 p + b)$  parallel time for F&A and the whole partition takes  $\Theta(n/p + \log^2 p + b)$  parallel time.

*Proof.* F&A takes  $T(n, p) = \Theta(n/p) + C(b, p)$ , where  $C(b, p)$  is the cost of our cleanup algorithm.  $C(b, p) = b + \log p + T'(b/2, \log p)$  parallel time, and  $T'$  is defined by the following recurrence:

$$T'(b, i) = \begin{cases} O(3b + \log p) + T'(b/2, i - 1) & \text{if } i > 1, \\ O(2b/p) & \text{otherwise.} \end{cases}$$

There are  $2\beta p$  blocks at the beginning of each recursive step and  $\log p - 1$  recursive steps are needed. Thus,  $C(b, p) = O(\log^2 p + b)$  parallel time.

Theorem 3 improves previous bounds for F&A (provided  $\log p \leq b$ , which is of practical relevance).

Table 1 summarizes worst-case results for BLOCKED and F&A algorithms.

## 4 Implementation

Since implementations of STRIDED were not available, we have resorted to implement it ourselves. We have also implemented our BLOCKED variant, which improves STRIDED cache performance. As for F&A, we have taken its implementation from MCSTL 0.7.3-beta and we have implemented it ourselves.

Our implementation of F&A and the one in MCSTL differ in the following: a) ours statically assigns the initial work and, so, avoids mutual exclusion here; b) ours does not use `volatile` variables and critical regions are slightly simpler; and c) ours avoids redundant comparisons using a better book-keeping.



On the other hand, we have implemented our cleanup algorithm on the top of the previous four algorithms.

The implementation is available at <http://www.lsi.upc.edu/~lfrias>. It uses C++ and OpenMP. Besides, it follows the specification of the `partition` function of the STL, so that it can be used instead of the sequential implementation.

## 5 Experimental Analysis

We have analyzed STRIDED, BLOCKED, and F&A with and without our cleanup algorithm.

The tests have been run on a machine with 4 GB of main memory and two sockets, each one with an Intel Xeon quad-core processor at 1.66 GHz with a shared L2 cache of 4 MB shared among two cores. Thus, there are 8 cores in total. We have used the GCC 4.2.0 compiler with the `-O3` optimization flag.

All tests have been repeated 100 times; figures show averages.

*Basic Evaluation.* We have first analyzed the speedup of partitioning in parallel a large number of random integers. The speedup is always measured with respect to the sequential partition algorithm of the STL. The block size  $b$  has been set to  $10^4$  (see the reason below).

Figure 2 shows the results. In this figure, `Strided` refers to our implementation of STRIDED, `BlockedStrided` refers to our implementation of BLOCKED, `F&A_MCSTL` refers to the MCSTL implementation of F&A, `F&A` refers to our own implementation of F&A. We add the suffix `_tree` to the previous labels to refer to the algorithm with our modified parallel cleanup phase.

These results show that F&A is better than BLOCKED, which is better than STRIDED. Whereas the speedup of STRIDED is nonexistent for more than two threads, BLOCKED performs reasonably well and our F&A implementation achieves some better results than the MCSTL F&A. Besides, using our cleanup phase maintains the same speedups for STRIDED and BLOCKED and improves slightly the speedup of F&A, making it almost perfect for up to 4 threads.

The awful performance of STRIDED is due to its high cache miss ratio; its behavior clearly contrasts with BLOCKED (which uses blocks of elements rather than individual elements).

In order to understand the loss of performance when using more than 5 threads we have devised two new experiments. The first one reproduces the previous experiment but uses a slower comparison function. Its results are shown in Fig. 3. In this case, all algorithms show similar behavior and excellent speedups with up to 8 threads. Specifically, there is not much of a difference whether our cleanup phase is used or not. Our second experiment has consisted in measuring the speedup of a trivial parallel program to compute the sum of two arrays. The resulting speedups (not shown) are also not optimal for the biggest number of threads. So, we can conclude that memory bandwidth is limiting the efficiency of the partitioning algorithms, which are demanding with regard to I/O.

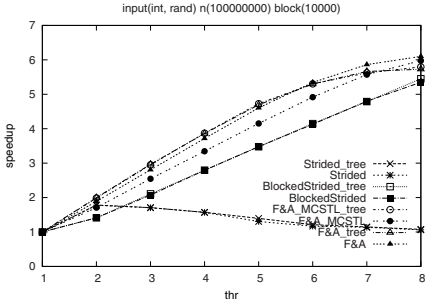


Fig. 2. Parallel partition speedup,  $n = 10^8$  and  $b = 10^4$

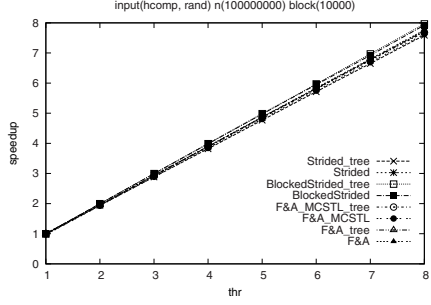


Fig. 3. Parallel partition speedup for costly  $<$ ,  $n = 10^8$  and  $b = 10^4$

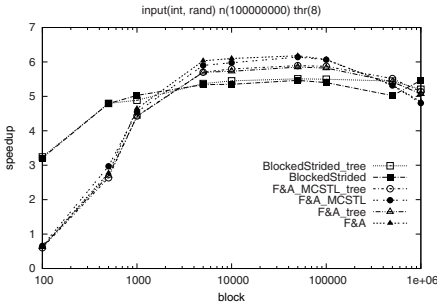


Fig. 4. Parallel partition with varying block size,  $n = 10^8$  and  $\text{num\_threads} = 8$

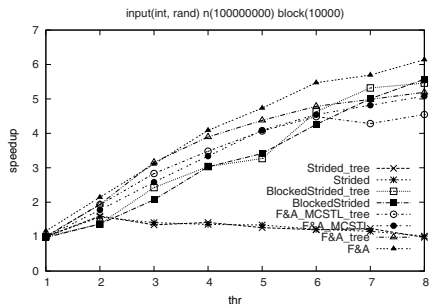


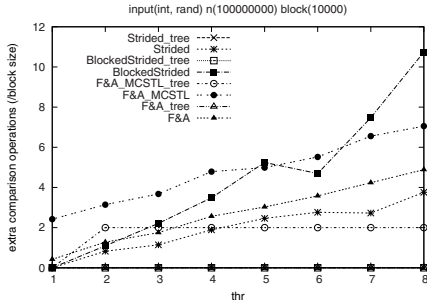
Fig. 5. Parallel quickselect speedup,  $n = 10^8$  and  $b = 10^4$

*Influence of Block Size.* Several algorithms rely on a block size parameter  $b$ . In order to determine its optimal value, we have run various tests, with 8 threads and different values of  $b$ . The results in Fig. 4 show that, except for very small block sizes, the performance is not much affected. Besides, given that for smaller input sizes, big block sizes are not convenient, our selection has been  $b = 10^4$ .

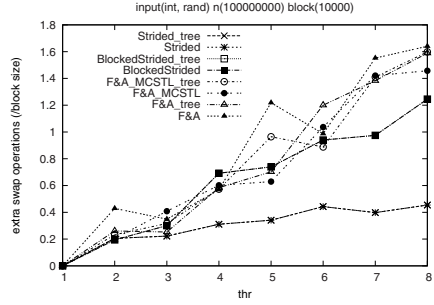
*Operations Count.* In order to analyze the behavior of the cleanup phase in more detail, we have counted swap and comparison operations. Figures 6 and 7 show respectively the number of extra comparisons and swaps with respect to the sequential implementation. They are depicted divided by the block size.

Figure 6 gives an experimental proof of Theorem 1. Combining our cleanup algorithm with the original MCSTL algorithm does not achieve the optimality in the number of comparisons, because this implementation makes extra comparisons whenever a new block is fetched in the main parallel phase. Specifically, our experiments show that two comparisons are repeated per block in the average.

Figure 7 shows that our cleanup algorithm does not need more swaps than the original cleanup algorithms. Essentially, the same number of extra swaps are needed. In the case of F&A, we could not give such an equality analytically.



**Fig. 6.** Number of extra comparisons,  $n = 10^8$  and  $b = 10^4$



**Fig. 7.** Number of extra swaps,  $n = 10^8$  and  $b = 10^4$

As a by product, these results show that the number of misplaced elements resulting from the parallel phase is really small, no matter the partitioning algorithm. In particular, STRIDED is the algorithm that performs less extra operations (for a random input of integers). However, its performance is the worst because of its bad cache usage.

*Application: Quickselect.* Quicksort and quickselect are typical applications of partitioning. As quicksort offers two (not exclusive) ways to be parallelized — parallelizing the partitioning and parallelizing the independent work by divide and conquer—, we found more interesting analysing quickselect. For this test, we have made that the STL `nth_element` function calls the parallel partitioning algorithms in this paper instead of the sequential (unless the array is small).

The results are shown in Fig. 5. These are coherent with those of partition but different given that the relative behavior between the algorithms changes slightly with the size of the input. First, our F&A implementation advantage increases. Second, our cleanup algorithm harms a little F&A based quickselect. Indeed, in our experiments we have observed that for a big number of threads and as input gets smaller, using our cleanup algorithm with F&A is counterproductive. Finally, Fig. 5 shows that the simple BLOCKED algorithm performs quite well.

## 6 Conclusions

In this paper we have presented, implemented and evaluated several parallel partitioning algorithms suitable for multi-core architectures.

From an algorithmic point of view, we have described a novel cleanup parallel algorithm that does not disregard comparisons made during the parallel phase. This cleanup has successfully been applied to three partitioning algorithms: STRIDED, BLOCKED (a cache-aware implementation of the former) and F&A. In the case of STRIDED and BLOCKED, a benefit of our cleanup is reducing its parallel time in the worst case from  $\Theta(n)$  to  $\Theta(n/p + \log p)$ . In the case of F&A, we have shown how to modify it to reduce its parallel time from

$\Theta(n/p + b \log p)$  to  $\Theta(n/p + \log^2 p)$ . Unlike their original versions, these algorithms perform the minimal number of comparisons when using our cleanup phase.

As automatic parallelization is still limited, and as parallel programming is hard and expensive, the use of parallel libraries is a simple way to benefit from multi-core processors. From this engineering perspective, we have contributed carefully designed implementations of the afore mentioned algorithms that are compliant with the specification of STL `partition`.

Finally, and from an experimental point of view, we have conducted an evaluation to compare those algorithms and implementations. According to our experiments, the partitioning algorithm of choice is F&A, because it scales nicely. Moreover, our implementation performs slightly better than the one in MCSTL. However, the results also show that, in practice, the benefits of our cleanup algorithm are limited. This happens because the number of misplaced elements after the parallel phase is very small.

Our experiments also show that I/O between the memory and the processor limits the performance achieved by parallel implementations as the number of threads increases. It remains to be further investigated how these results change for a bigger number of available cores or/and memory bandwidth.

## References

1. Heidelberg, P., Norton, A., Robinson, J.T.: Parallel quicksort using fetch-and-add. *IEEE Trans. Comput.* 39(1), 133–138 (1990)
2. Francis, R.S., Pannan, L.J.H.: A parallel partition for enhanced parallel quicksort. *Parallel Computing* 18(5), 543–550 (1992)
3. Tsigas, P., Zhang, Y.: 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), pp. 372–381 (2003)
4. Liu, J., Knowles, C., Davis, A.: A cost optimal parallel quicksorting and its implementation on a shared memory parallel computer. In: Pan, Y., Chen, D.-x., Guo, M., Cao, J., Dongarra, J. (eds.) *ISPA 2005. LNCS*, vol. 3758, pp. 491–502. Springer, Heidelberg (2005)
5. Singler, J., Sanders, P., Putze, F.: The Multi-Core Standard Template Library. In: Kermarrec, A.-M., Bougé, L., Priol, T. (eds.) *Euro-Par 2007. LNCS*, vol. 4641, pp. 682–694. Springer, Heidelberg (2007)
6. International Standard ISO/IEC 14882: Programming languages — C++. 1st edn. American National Standard Institute (1998)
7. Gottlieb, A., Grishman, R., Kruskal, C.P., McAuliffe, K.P., Rudolph, L., Snir, M.: The NYU ultracomputer - designing a MIMD, shared-memory parallel machine. In: *ISCA 1998: 25 years of the international symposia on Computer architecture (selected papers)*, pp. 239–254. ACM Press, New York, NY, USA (1998)
8. Cormen, T.H., Leiserson, C., Rivest, R., Stein, C.: *Introduction to algorithms*, 2nd edn. The MIT Press, Cambridge (2001)
9. Kumar, V.: *Introduction to Parallel Computing*. Addison-Wesley, Boston, MA, USA (2002)
10. JáJá, J.: *An introduction to parallel algorithms*. Addison-Wesley, Redwood City, CA, USA (1992)