

# Resource Oblivious Sorting on Multicores

Richard Cole <sup>\*</sup> and Vijaya Ramachandran <sup>\*\*</sup>

<sup>1</sup> Computer Science Dept., Courant Institute, NYU, New York, NY 10012. Email: [cole@cs.nyu.edu](mailto:cole@cs.nyu.edu).

<sup>2</sup> Dept. of Computer Sciences, UT, Austin, TX 78712. Email: [v1r@cs.utexas.edu](mailto:v1r@cs.utexas.edu).

**Abstract.** We present a new deterministic sorting algorithm that interleaves the partitioning of a sample sort with merging. Sequentially, it sorts  $n$  elements in  $O(n \log n)$  time cache-obliviously with an optimal number of cache misses. The parallel complexity (or critical path length) of the algorithm is  $O(\log n \log \log n)$ , which improves on previous bounds for deterministic sample sort. Given a multicore computing environment with a global shared memory and  $p$  cores, each having a cache of size  $M$  organized in blocks of size  $B$ , our algorithm can be scheduled effectively on these  $p$  cores in a cache-oblivious manner.

We improve on the above cache-oblivious processor-aware parallel implementation by using the Priority Work Stealing Scheduler (PWS) that we presented recently in a companion paper [12]. The PWS scheduler is both processor- and cache-oblivious (i.e., resource oblivious), and it tolerates asynchrony among the cores. Using PWS, we obtain a resource oblivious scheduling of our sorting algorithm that matches the performance of the processor-aware version. Our analysis includes the delay incurred by false-sharing. We also establish good bounds for our algorithm with the randomized work stealing scheduler.

## 1 Introduction

We present a new parallel sorting algorithm, which we call *Sample, Partition, and Merge Sort (SPMS)*. It has a critical path length of  $O(\log n \log \log n)$  and performs optimal  $O(n \log n)$  operations with optimal sequential cache misses. More importantly, using the PWS scheduler for multicores developed and analyzed in [12], and new algorithmic techniques given in this paper, we can schedule it resource-obliviously on a multicore while maintaining these performance bounds. We present background information on multicores, cache-efficiency and resource-obliviousness in Section 2.

The core of the sorting algorithm is a recursive multi-way merging procedure. A notable and novel aspect of this procedure is that it creates its recursive subproblems using a sample sort methodology. We view the sorting algorithm as interleaving a merge sort with a sample sort in a natural way.

**Previous Work.** Sorting is a fundamental algorithmic problem, and has been studied extensively. For our purposes, the most relevant results are sequential

---

<sup>\*</sup> This work was supported in part by NSF Grant CCF-0830516.

<sup>\*\*</sup> This work was supported in part by NSF Grant CCF-0850775 and CCF-0830737.

cache-oblivious sorting, for which provably optimal algorithms are known [15], optimal sorting algorithms addressing pure parallelism [3, 11], and recent work on multicore sorting [5, 4, 6, 16].

The existing multicore algorithms take two main approaches. The first is merge sort [4, 6, 5], either simple or the pipelined method from [11]. The second is deterministic sampling [16]: this approach splits the input into subsets, sorts the subsets, samples the sorted subsets, sort the sample, partitions about a subsample, and recursively sorts the resulting sets. Our algorithm can be viewed as applying this approach to the problem of merging a *suitable number* of sorted sets, which eliminates the need for the first two steps, resulting in significant speed-up.

More specifically, the algorithm in [6] is a simple multicore mergesort; it has polylog parallel time, and good, though not optimal cache efficiency; it is cache-oblivious for private caches (the model we consider in this paper). The algorithm in [4] achieves the optimal caching bound on an input of length  $n$ , with  $O(\log n)$  parallel time (modulo dependence on cache parameters), but it is both cache-aware and core-aware; this algorithm is based on [11]. The algorithm in [5] is cache oblivious with  $O(\log^2 n)$  parallel time, but due to an additive term the cache performance is not optimal on a multicore. The algorithm in [16] is designed for a BSP-style version of a cache aware, multi-level multicore. It uses a different collection of parameters, and so it is difficult to compare with it directly.

**Roadmap.** In Section 2 we present some background on multicores, and then state our main sorting result. In Section 3 we give a high level description of our parallel sorting algorithm, omitting the details needed to have a resource oblivious implementation. In Section 4, we review the computation model, the work stealing scheduler PWS, and the class of BP algorithms, as developed in [12]. In Section 5, we return to the sorting algorithm, describing the details needed for a resource oblivious implementation, and this is followed by its analysis. Due to space constraints, our matching lower bound for cache optimality (adapted from [2]), along with most of the proofs, are deferred to the full paper [13].

## 2 Statement of our Results

Before stating our main result, we give some background, as developed in [12].

**Multicore with Private Caches.** We model a multicore as consisting of  $p$  *cores* (or processors) with an arbitrarily large main memory, which serves as a shared memory. Additionally, each core has a private cache of size  $M$ . Data in the main memory is organized in blocks of size  $B$ , and the initial input of size  $n$  is in main memory, in  $n/B$  blocks. When a core  $C$  needs a data item  $x$  that is not in its private cache, it reads in the block  $\beta$  that contains  $x$  from main memory. This new block replaces an existing block in the private cache, which is evicted using an optimal cache replacement policy (LRU suffices for our algorithms, but we do not elaborate further). If another core  $C'$  modifies an entry in  $\beta$ , then  $\beta$  is

*invalidated* in  $C$ 's cache, and the next time core  $C$  needs to access data in block  $\beta$ , an updated copy of  $\beta$  is brought into  $C$ 's cache.

**Cache and Block Misses.** We distinguish between two types of cache-related costs incurred in a parallel execution.

The term *cache miss* denotes a read of a block from shared-memory into core  $C$ 's cache, when a needed data item is not currently in the cache, either because the block was never read by core  $C$ , or because it was evicted from  $C$ 's cache to make room for new data. This is the standard type of cache miss that occurs, and is accounted for, in sequential cache complexity analysis.

The term *block miss* denotes an update by a core  $C' \neq C$  to an entry in a block  $\beta$  that is in core  $C$ 's cache; this entails core  $C'$  acquiring block  $\beta$ ; if  $C$  has a subsequent write, it needs to reacquire the block. This type of 'cache miss' does not occur in a sequential computation, and is a problematic situation that can occur quite often, especially in the resource oblivious setting that we seek.

**Resource Obliviousness.** We have claimed that our multicore algorithms are resource oblivious: we mean that the algorithm is specified without any mention of the multicore parameters ( $p$ ,  $M$  and  $B$ ) and further, the PWS scheduler we use schedules tasks on available idle cores, without reference to the multicore parameters. Since multicores with a wide range of parameters are expected to appear on most desktops, such a resource oblivious feature in multicore algorithms appears to be helpful in supporting the portability of program codes. The PWS scheduler uses *work-stealing* [9, 7], where load balance is achieved by cores stealing tasks from other cores as needed.

Our main result, the SPMS sorting algorithm and its resource-oblivious performance, has the bounds stated below in Theorem 1. The algorithm proceeds in *rounds*, where a round, roughly speaking, corresponds to a parallel step. Our analysis uses the following parameters. We suppose that each core performs a single operation in  $O(1)$  time, a cache miss takes at most  $b$  time, a steal request takes at most  $s$  time (whether successful or not), and the scheduler's work at the start of each round takes at most  $S$  time. We consider a multicore with  $p$  cores, each having a private cache of size  $M$  organized in blocks of size  $B$ , with all caches sharing an arbitrarily large global memory. The input, of size  $n \geq Mp$  (this restriction ensures that both cores and caches can be fully utilized), is in the shared memory at the start of the computation, and SPMS is scheduled under PWS. Then:

**Theorem 1.** *On an input of length  $n$ , assuming  $M \geq B^2$  (the 'tall cache'), for  $p \leq \frac{n}{\max\{\log \log n, M\}}$ , the sorting algorithm SPMS takes parallel time*

$$O\left(\frac{1}{p} \left( n \log n + b \cdot \frac{n \log n}{B \log M} \right) + (b + s + S) \log n \log \log n + b\beta(n, p, B) \right).$$

*The fourth term,  $\beta(n, p, B) = O(B \log n \log \log(n/p))$ , is the block miss cost, and is bounded by the optimal sequential cache complexity provided  $p \leq \frac{n}{B^2 \log \log M \log M}$  (i.e., with a slightly 'taller' cache —  $M \geq B^2 \log B \log \log B$  suffices). This cost may also be reduced to match the optimal sequential cache complexity without*

this additional restriction on  $p$  if system support is provided for the locking of a block during writes, and limiting the minimum task size to be at least  $B$ .

If we ignore the block miss cost for the moment, this bound represents the optimal work bound, plus the optimal cache miss bound, plus the critical path length times the cost of one cache miss plus one steal plus one scheduling event. Further, we note that there is no need for a global clock, or tight synchronization on the part of the cores, though the scheduler certainly imposes a significant degree of synchronization. The computation is entirely resource-oblivious in that the algorithm makes no mention of  $p$ ,  $M$  or  $B$ , and PWS services idle cores without any reference to the number available or their cache parameters.

Extending the analysis of the randomized work stealer in [8, 1], we can obtain:

**Theorem 2.** *On an input of length  $n$ , assuming  $M \geq B^2$ , for  $p \leq \frac{n}{\max\{\log \log n, M\}}$ , the sorting algorithm SPMS when scheduled by the randomized work stealer, and taking into account both cache and block misses, takes expected parallel time*

$$O\left(\frac{1}{p} \left(n \log n + b \cdot \frac{n \log n}{B \log M}\right) + \frac{M}{B} \cdot \frac{b}{s} \left(bB \frac{\log n}{\log B} + (b + s) \log n \log \log n\right)\right).$$

(The analysis of this result can be found in [12, 14]).

**Discussion.** Our sorting algorithm is optimal in all respects except for the critical pathlength. The sorting algorithm for PEM in [4] achieves optimal  $O(\log n)$  parallel steps, but is both cache- and core-aware. Achieving the same bound in a resource-oblivious manner appears considerably more challenging, and it is not clear if it is possible. We leave this as a topic for further research.

Another challenging topic is to extend our results to resource-oblivious scheduling on a multi-level caching hierarchy. Given the conflicting requirements of private and shared caches noted in [5, 10], it appears that some mechanism of supplying scheduler hints within the algorithm, and having a scheduler that uses the machine parameters effectively is needed. One such approach is used in [10]; however, that scheduler is *not* resource-oblivious, in contrast to our results.

In comparing our PWS scheduling to the PEM and Multi-BSP models, we note that these models both compute in a bulk-synchronous manner. We can easily adapt our results to work on either PEM or multi-BSP with the same performance as achieved with the PWS scheduler. However, our PWS framework adapts much more gracefully to differences in speeds among the cores than these bulk-synchronous models. Thus, if we have a few cores that execute faster than others (perhaps because they have smaller cache miss cost due to the cache layout), then PWS would enable the faster cores to take over (i.e. steal) work from the slower cores, balancing the work across the cores more effectively.

### 3 SPMS, A New Deterministic Sample, Partition, and Merge Sort

The heart of the algorithm is a procedure for computing a merging subproblem  $MS$ , whose input comprises  $r$  sorted lists  $L_1, L_2, \dots, L_r$ , of total length  $m$ , with  $m \leq r^c$ , where  $c \geq 6$  is a constant.

The sorting algorithm simply calls the merging procedure with  $r = m = n$ .

The merging algorithm performs two successive collections of recursive  $\sqrt{r}$ -way merges, each merge being on lists of total length at most  $r^{c/2}$ . To enable this, suitable samples of the input lists will be sorted by a logarithmic time procedure, which then allows the original problem to be partitioned into smaller subproblems that are merged recursively. More precisely:

**Step 1.** Partition the original problem  $MS$  into  $k = O(m/r^{\frac{c}{2}-1})$  disjoint merging subproblems,  $M_1, M_2, \dots, M_k$ , each comprising  $r$  sorted lists, with each subproblem having at most  $r^{\frac{c}{2}}$  items in its  $r$  sorted lists. In addition, the items in  $M_i$  precede those in  $M_{i+1}$ , for  $1 \leq i < k$ .

**Step 2.** For each subproblem  $M_i$ , group its lists into disjoint subsets of  $\sqrt{r}$  lists, and then in parallel merge the lists in each group. As  $M_i$  contains at most  $r^{\frac{c}{2}}$  items, this bound applies to each of the individual groups too. Thus the  $\sqrt{r}$ -way merge in each group can be performed recursively. The output, for each subproblem  $M_i$ , is a collection of  $\sqrt{r}$  sorted lists of total length at most  $r^{\frac{c}{2}}$ .

**Step 3.** For each subproblem  $M_i$ , recursively merge the  $\sqrt{r}$  sorted lists computed in Step 2.

**Step 1 details.** The basic idea is to take a deterministic sample  $S$  of the input set comprising every  $r^{\frac{c}{2}}$ -th item in each list, to sort  $S$ , and to partition the  $r$  input lists about the items in  $S$  thereby forming smaller  $r$ -way merging subproblems. Some of these subproblems may have size as large as  $r^{\frac{c}{2}+1}$ , rather than the desired  $r^{\frac{c}{2}}$ . Any such subproblems are partitioned further, as needed, via samples  $S'$  of size  $m'/r^{\frac{c}{2}-1}$  for each subproblem of size  $m' \geq r^{\frac{c}{2}}$ . The samples  $S$  and  $S'$  are sorted by performing all pairwise comparisons. More precisely:

**Step 1.1.** Let  $S$  comprise every  $r^{\frac{c}{2}}$ -th item in each of the input lists. Extract  $S$  from the input lists and then sort  $S$ , using a simple logarithmic time, quadratic work algorithm.

**Step 1.2.** Partition the  $r$  input lists  $L_1, L_2, \dots, L_r$  about  $S$ , creating subproblems  $M'_1, M'_2, \dots, M'_{k'}$ , where  $k' = |S| + 1$ , and  $M'_i$  contains  $r$  sublists holding the items between the  $(i-1)$ th and  $i$ th items in  $S$ .

**Step 1.3.** Further partition any subproblem  $M'_i$  of size more than  $r^{\frac{c}{2}}$ , creating an overall collection of merging subproblems  $M_1, M_2, \dots, M_k$ , each of size at most  $r^{\frac{c}{2}}$ , with the further property that the items in  $M_i$  precede those in  $M_{i+1}$ , for  $1 \leq i < k$ . This is done using a sample comprising every  $r^{\frac{c}{2}-1}$ -th item in  $M'_i$ .

**Lemma 1.** *The merging algorithm, on an input of  $r$  sorted lists of total length  $m \leq r^c$ , uses  $O(m \log r)$  operations and  $O(\log r \log \log r)$  parallel time, if  $c \geq 6$ .*

*Proof.* The parallel run time  $T(r, m)$  is given by:  $T(r, m) \leq \log r + 2T(\sqrt{r}, r^{c/2}) = O(\log r \log \log r)$ .

Clearly, Steps 1.1 and 1.2 take  $O(m)$  operations. To see the same bound applies to Step 1.3, we argue as follows. Each subproblem  $M'_i$  of size  $m'$  generates a sorting task of size  $m'/r^{\frac{c}{2}-1} \leq r^{\frac{c}{2}+1}/r^{\frac{c}{2}-1} = r^2$ . Performing all these sorting tasks requires at most  $r^2 \cdot \sum m'/r^{\frac{c}{2}-1} \leq r^2 \cdot m/r^{\frac{c}{2}-1} \leq m$  operations, if  $c \geq 6$ .

Let  $W(r, m)$  be the operation count for a collection of merging problems of total size  $m$ , where each comprises the merge of  $r$  lists of combined size at most  $r^c$ . Then we have:  $W(r, m) \leq m + 2W(r^{1/2}, m) = O(m \log r)$ .

**Corollary 1.** *The sorting algorithm, given an input of size  $n$ , performs  $O(n \log n)$  operations and has parallel time complexity  $O(\log n \log \log n)$ , if  $c \geq 6$ .*

## 4 The Computation Model and PWS Scheduling

Before giving the resource oblivious implementation, we need to review the computation model and the PWS scheduling environment, mainly as developed in [12], although we make some changes here to address some new algorithmic features in SPMS.

The building blocks for our algorithms are computations on balanced binary trees such as for prefix sums. Such a computation is carried out by tasks: initially there is one task at the root of the tree; it forks two subtasks for each of its subtrees, and when they are done, it resumes and concludes the computation at the root. We will also use a tree of forking tasks to initiate a collection of parallel recursive calls, as needed in the merging and sorting algorithms.

Initially the root task for such a tree is given to a single core. Subtasks are acquired by other cores via task stealing. To this end, each core  $C$  has a task queue. It adds forked tasks to the bottom of the queue, while tasks are stolen from the top of the queue. So in particular, when  $C$ , on executing  $\tau$ , generates forked tasks  $\tau_1$  and  $\tau_2$ , it places the larger of  $\tau_1$  and  $\tau_2$  on its queue,  $\tau_2$  say, and continues with the execution of  $\tau_1$ . This is a small generalization from [12], where the two forked tasks were assumed to be of the same size. When  $C$  completes  $\tau_1$ , if  $\tau_2$  is still on its queue, it resumes the execution of  $\tau_2$ , and otherwise there is nothing on its queue so it seeks to steal a new task. Except for one routine, our algorithm will be constructed from BP trees [12], which are trees of equal-sized forking nodes with an  $O(1)$  operation computation at each node, and with the leaf nodes having either an  $O(1)$  operation task or a recursive computation as their task. There is a mirror image tree for the joins which also performs  $O(1)$  operations at each node. We will often talk of a subtree of the BP tree, when we really intend a subtree plus the mirror image subtree.

Let  $\tau$  be a task associated with such a subtree. As in [12], by the size of  $\tau$ ,  $|\tau|$ , we mean the amount of data  $\tau$  accesses in its computation. In contrast to [12], sometimes we will use the *virtual size* of  $\tau$ ,  $vs(\tau)$ ; always  $vs(\tau) \geq |\tau|$ . Efficiency is ensured by the following BP tree property: if  $\tau'$  is forked by  $\tau$ , then  $vs(\tau') \leq \frac{1}{2}vs(\tau)$ .

To help with the scheduling, each node in a BP tree receives the integer priority  $\log vs(\tau)$ . These are strictly decreasing from parent to child. We will use the Priority Work-Stealing Scheduler (PWS) [12], which only allocates tasks of highest priority in response to steal requests. As noted in [12], task priorities are strictly decreasing on each task queue, and thus there will be at most one steal of a task of priority  $d$  from each core, and so at most  $p$  steals of tasks of priority  $d$ , for any  $d$ . This is key to bounding the overhead of the PWS scheduler.

As noted in Section 2, the I/O cost of an individual task  $\tau$  is measured in *cache misses*, which we upper bound by how many blocks the core executing  $\tau$  has to read into its cache, assuming none are present at the start of  $\tau$ 's execution,

and *block misses*, which capture the cost of multiple cores writing to the same block.

As we shall see, each BP task  $\tau$  in the sort algorithm incurs  $O(\text{vs}(\tau)/B + \sqrt{\text{vs}(\tau)})$  cache misses when executed sequentially. Each task incurs only  $O(B)$  block miss delay: for most of the tasks this follows from [12] because they engage in consecutive writes to a linear array; we will show that the remaining class of tasks will also incur only  $O(B)$  block miss delay in their writing.

We will use the following bounds derived in [12] for a collection of parallel BP computations of total size  $n$  and sequential cache complexity  $Q$ , when scheduled under PWS. Here, the maximum (virtual) size of any root task in the BP collection is  $x$ , and any task of size  $s$  incurs  $O(s/B + \sqrt{s})$  cache misses and shares  $O(1)$  blocks with other tasks:

**Fact 1** *For the I/O cost for a computation of the type stated above:*

1. *The cache miss bound is  $O(Q + p \cdot (\frac{\min\{M,x\}}{B} + \log \min\{x, B\} + \sqrt{x}))$ .*
2. *The block miss bound is  $O(p \cdot \min\{B, x\} \cdot (1 + \log \min\{x, \frac{n}{p}\}))$ .*

The merging algorithm *MS* is built by combining (collections of parallel) BP computations, first by sequencing, and second by allowing the leaves of a BP tree to be recursive calls to the merging algorithm. This generalizes the above tree computation to a dag which is a series-parallel graph. The formal definition of such an ‘HBP’ computation is given in [12]. While we do not need the details of a general HBP computation here, we do need to define priorities carefully in view of the possible differences in the sizes of the recursive subproblems generated by a call to *MS*. We define these priorities in a natural way so that they are strictly decreasing along any path in the *MS* computation dag, and all tasks with the same priority have roughly the same size, as detailed in the full paper [13].

The recursive subproblems generated in Step 2 of *MS* need not be of the same size, so this portion of the algorithm does not exactly fit the BP framework of [12]. To handle this, we will first determine the cache-miss overhead for the natural parallel implementation of the algorithm, which we call the *ideal PWS costing*, and then add in the additional cache-miss cost for the PWS schedule. (The cost of block misses is discussed later.)

**Definition 1.** *The ideal costing assumes that a BP computation uses at most  $2n/M$  cores, one for each distinct subtree of size  $M$  and one for each node ancestral to these subtrees.*

The cache miss cost in the ideal costing is  $O(M/B)$  per subtree, plus  $O(1)$  for each ancestral node, for a total of  $O(n/B)$  cache misses.

We generalize this BP analysis to *MS* and *SPMS* by analyzing the algorithm in terms of parallel collection of tasks, each task of virtual size  $M$ . The cost of each task collection is bounded in turn: each task is costed as if it was allocated to a distinct core. As we will see, each such collection has total virtual size  $O(n)$ , and hence incurs  $O((n/B) + \frac{n}{M}\sqrt{M})$  cache misses, which is  $O(n/B)$  if  $M \geq B^2$ .

To analyze the cache miss cost of the PWS scheduling of *MS*, we separate the cost of steals of small tasks  $\tau$  (those with  $\text{vs}(\tau) \leq M$ ), which we bound later, and

consider the impact of steals of large tasks. To bound this cost for a large stolen task  $\tau$ , we overestimate by supposing that no small tasks are stolen from  $\tau$ . Then (the possibly overestimated)  $\tau$  executes one or more of the size  $M$  subtrees that are executed as distinct tasks in the ideal PWS costing, plus zero or more nodes ancestral to these subtree. Clearly,  $\tau$ 's cache miss cost is at most the sum of the cache miss costs of the corresponding distinct tasks in the ideal PWS costing. An analogous claim holds for the root task. Summing over the costs of the root task and of the large stolen tasks, yields that their total cost is bounded by the ideal PWS costing. This also bounds the processor-aware, cache-oblivious cost of MS, since the cost for block misses is minimal when there are no steals. We bound the cost of steals of small tasks using Fact 1, and results we derive here.

A final point concerns the management of local variables in recursive calls. We assume that if a task  $\tau$  stolen by a core  $C$  has local variables, then the space allocated by the memory manager for these variables does not share any blocks with space allocated to other cores. Further, if the data resides in cache till the end of  $C$ 's execution of  $\tau$ , then the now unneeded local variables are not written back to the shared memory.

## 5 The Cache-Oblivious Parallel Implementation of SPMS

To achieve efficient oblivious performance, the merging algorithm MS needs to be implemented using tasks achieving the optimal ideal costing, as defined above. Many of the steps in MS are standard BP computations; their ideal costing is  $O(n/B)$  and their PWS overhead can be bounded directly using Fact 1. However, here we address three types of computations in MS that do not fall within the framework in [12].

1. The recursion may well form very unequal sized subproblems. However, to achieve a small cache miss cost, the PWS scheduling requires forking into roughly equal sized subtasks. Accordingly we present the method of *grouping unequal sized tasks*, which groups subproblems in a task tree so as to achieve the balanced forking needed to obtain good cache-miss efficiency.
2. Balancing I/O for reads and writes in what amount to generalized transposes, which we call *transposing copies*. This issue arises in the partitioning in Steps 1.2 and 1.3. The challenge faced by a multicore implementation results from the delay due to the block misses, as discussed earlier.
3. One collection of tasks for sorting the samples in Step 1 uses non-contiguous writes. Fortunately, they use relatively few writes. We develop the *sparse writing* technique to cope.

**Grouping Unequal Sized Tasks.** We are given  $k$  ordered tasks  $\tau_1, \tau_2, \dots, \tau_k$ , where each  $\tau_i$  accesses  $O(|\tau_i|/B)$  blocks in its computation (they are all recursive merges). Let  $t_i \leq t_{ave}^2$  for all tasks, where  $t_{ave}$  is the average size of the tasks.

The tasks need to be grouped in a height  $O(\log k)$  binary tree, called the *u-tree*, with leaves holding the tasks in their input order. The u-tree is used for the forking needed to schedule the tasks. The u-tree will use virtual sizes for its scheduling subtasks and has the bounds given below. See [13] for more details.



**Lemma 2.** *The ideal PWS costing for scheduling the  $u$ -tree plus the cost of executing tasks  $\tau_i$  of size  $M$  or less is  $O(\sum_{i=1}^k t_i/B)$ , where  $t_i = |\tau_i|$ .*

**The Transposing Copy.** The problem, given a vector  $A$  consisting of the sequence  $A_{11}, \dots, A_{1k}, \dots, A_{h1}, \dots, A_{hk}$  of subvectors, is to output the transposed sequence  $A_{11}, \dots, A_{h1}, \dots, A_{1k}, \dots, A_{hk}$ , where we are given that the average sequence length  $l = |A|/hk \geq h$ .

This is done by creating  $\lceil \frac{|A_{ij}|}{l} \rceil$  tasks of virtual size  $l$  to carry out the copying of  $A_{ij}$ . The tasks are combined in column major order, i.e. in the order corresponding to destination locations. The full paper proves the following bound.

**Lemma 3.** *Let  $\tau$  be a task copying lists of combined size  $s$  in the transposing copy. Then  $\tau$  incurs  $O(s/B + \sqrt{s})$  cache misses.*

**Sparse Writing.** Let  $A$  be an  $s \times s$  array in which each of locations  $c \cdot s, 1 \leq c \leq s$  is written exactly once, but not in any particular order. A sequential execution incurs at most  $s$  cache misses.

Now consider a BP execution of this computation in which each leaf is responsible for one write. We claim that the I/O cost for all writes to  $A$  is  $O(s^2/B + B)$  regardless of the ordering of the writes. We establish this bound as follows.

If  $s \geq B$ , each write into  $A$  incurs one cache miss, for a total cost of  $O(s) \leq O(s^2/B)$  cache misses. There are no block misses in this case.

If  $s < B$ , there are only  $s$  accesses, but these can incur block misses. Let  $i$  be the integer satisfying  $s \cdot i \leq B < s \cdot (i + 1)$ . Then, at most  $i$  writes occur within a single block. The worst-case configuration in terms of block miss cost occurs when there are  $i$  different sets of  $s/i$  writes, with each set having one write per block. Each such write to a block may incur a block wait cost equal to that of  $\Theta(i)$  cache misses. Hence the overall delay in this case is at most that of  $O(i^2 \cdot s/i) = O(s \cdot i) = O(B)$  cache misses.

## 5.1 Details of Step 1 in SPMS

Now, we describe parts of the algorithm in detail.

Each substep (except one) uses a BP computation or a collection BP computations running in parallel. We characterize the complexity of each size  $x$  (collection of) computations. Clearly it will have depth  $O(\log x)$ , and unless otherwise specified will incur  $O(x/B)$  cache misses.

We begin with some notation. Let  $L_1, L_2, \dots, L_r$  be the  $r$  sorted input lists of total length  $m \leq r^c$ . The  $r$  lists are stored in sequential order. Let  $S = \{e_1, e_2, \dots, e_s\}$  comprise every  $r^{\frac{c}{2}}$ th item in the sequence of sorted lists; recall that  $S$  is sorted in Step 1.1, and then used to partition the  $r$  lists in Step 1.2.

**Step 1.1.** Sort  $S$ .

**1.1.1.** Construct arrays  $S_1, S_2, \dots, S_s$ ; each  $S_i$  is an array of length  $s$  which contains a copy of the sequence of elements in  $S$ .

(a) Compact the  $s$  samples within the list sequence  $L_1, \dots, L_r$ , using prefix sums for the compaction. The result is an array  $S_1[1..s]$  containing the  $s$  samples. This uses a sequence of 2 BP computations of size  $m$ .

(b) Form arrays  $S_i$ ,  $2 \leq i \leq s$ , where each  $S_i$  is a copy of  $S_1$ . This is a BP computation of size  $s^2 \leq m$ .

**1.1.2.** In parallel for each  $i$ , compute rank of  $e_i$  in  $S_i$ .

First, for each  $S_i$ , compare  $e_i$  to each element in  $S_i$ . Then count the number of  $e_j \geq e_i$ , the desired rank of  $e_i$  in  $S_i$ . This uses two BP computations, and over all  $i$ ,  $1 \leq i \leq s$ , they have combined size  $O(s^2)$ .

**1.1.3.** Create the sorted array  $S[1..s]$  where  $S[i]$  contains the element  $e_j$  with rank  $\rho_j = i$ .

The simple way to implement this step is for each element  $e_i$  to index itself into location  $S[\rho_i]$ . This will incur  $s$  cache misses, which can be argued is acceptable with a tall cache, since  $s^2 = O(m)$ . But this implementation could incur  $s \cdot B$  block misses, which is excessive. To reduce the block miss cost, we split this step into two substeps:

(a) Initialize an all-zero auxiliary array  $A'[1..m]$  and write each  $e_i$  into location  $A'[\rho_i \cdot r^{c/2}]$ .

This is the sparse writing setting analyzed earlier, and results in  $O(s^2/B + B) = O(m/B + B)$  cache and block misses in a depth  $\log s$  computation.

(b) Compact array  $A'$  into  $S[1..s]$ , which gives the desired sorted array of the samples. This is a prefix sums computation, a BP computation of size  $O(m)$ .

**Step 1.2.** Partition  $L_1, L_2, \dots, L_r$  about  $S$ . (Details in [13].)

**Step 1.3.** For each subproblem  $M'_i$  with  $|M'_i| > r^{\frac{d}{2}}$  create a task to further partition  $M'_i$ . It is analogous to Steps 1.1 and 1.2 except that it uses a sample  $S'$  of size  $m'_i/r^{\frac{d}{2}-1}$ , where  $m'_i = |M'_i|$ .

**Ideal PWS Costing.** Summarizing the above discussion of the cache-miss costs for the merge ( $MS$ ) gives the following bound for the number of cache misses in the ideal PWS costing.

**Lemma 4.** *In the ideal PWS costing, the merging algorithm  $MS$ , in performing a collection of merging tasks of total size  $n \geq Mp$ , in which each task comprises the merge of  $r$  lists of combined length at most  $r^c$ , incurs  $O(\lceil \frac{n}{B} \rceil \lceil \frac{\log r}{\log M} \rceil)$  cache-misses, if  $c \geq 6$  and  $M \geq B^2$ .*

*Proof.* As argued in the description of the algorithm, for each merging problem of size  $m = \Omega(M)$ , Substep 1 incurs  $O(m/B + \sqrt{m}) = O(m/B)$  cache-misses, as  $M \geq B^2$ ; smaller subproblems fit in cache and so incur  $O(\lceil m/B \rceil)$  cache-misses.

Now let  $C(r, n)$  be the cache-miss count for performing such a collection of merges for problems of merging  $r$  lists each of combined size at most  $r^c$ . Then, as the algorithm uses  $O(n)$  space, we have, for a suitable constant  $\gamma > 1$ : for  $n \leq \gamma M$ :  $C(r, n) = \lceil n/B \rceil$ , and for  $n \geq \gamma M$ :  $C(r, n) \leq \frac{n}{B} + 2C(r^{1/2}, n)$ .

For a processor-aware, cache-oblivious implementation of SPMS, there is only a constant number of block misses per task executed by a core, costing  $O(bB)$

per task, and the number of tasks in a  $p$ -core processor-aware implementation is  $O(p \cdot \frac{\log n}{\log(n/p)})$ . Thus, the block miss cost is dominated by the cache miss cost under our assumption that  $n \geq Mp$  and  $M \geq B^2$ . Hence, with the above analysis and the parallel time bounds for the basic SPMS algorithm, as well as for the BP computations in the implementations given in this section, we obtain the result that SPMS can be scheduled on  $p$  cores, for  $p \leq \frac{n}{\max\{M, \log \log n\}}$ , to obtain optimal speed-up and cache-oblivious cache-efficiency. Note that in such a processor-aware schedule, there is no need for steals, and hence there is no further overhead beyond the cache-miss, block miss, and depth bounds that we have established for the computation.

We next establish that when scheduled under PWS, this implementation also achieves similar bounds resource-obliviously.

**The Analysis of the PWS overhead.** In addition to the results in Fact 1, the companion paper [12] shows that:

1. The cost of each up-pass is bounded by that of the corresponding downpass in BP and HBP algorithms.

2. The idle work (the time spent by a core when it is not computing nor writing on a cache or block miss), in a (parallel collection of) BP computations, aside the waiting already accounted for in the up-pass, is bounded by  $O(p \cdot ((s + S + b) \log x + b \min\{x, B\}))$  where  $x$  is the size of the largest task in the collection,  $s$  bounds the time for a steal,  $S$  bounds the time to initiate a scheduling round, and  $b$  bounds the time for a cache miss.

3. Additional cache miss costs due to small tasks taking over the work of large tasks on an up-pass are bounded by the cache miss costs in the downpass.

And, as already noted in this paper, for the present algorithm:

4. The delay due to block misses for a stolen task  $\tau$  is bounded by the time for  $O(B)$  cache misses. This follows from results in Fact 1 for block misses, and from our method for Step 1.1.3, described earlier.

**Lemma 5.** *Let  $M \geq B^2$ . The delay  $\text{BM}_M(n, r^c)$  due to block misses in the merging algorithm for a collection of merging problems each of size at most  $r^c$ , and of total size  $n \geq Mp$ , is bounded by:  $pB \log r^c (\log \log \frac{n}{p} - \log \log B)$  if  $r^c \geq B$  and  $B \leq \frac{n}{p} < r^c$ ,  $pB \log r^c (\log \log \frac{r^c}{B} - \log \log B)$  if  $r^c \geq B$  and  $\frac{n}{p} \geq r^c$ , and by  $pr^c \log r^c$  if  $n/p, r^c < B$ .*

*Proof.* Using the bounds in Fact 1 for block misses, and since  $M \geq B^2$  the top level BP computation causes the following number,  $\text{BMT}(n, r^c)$ , of block misses:  $pB \log \frac{n}{p}$  if  $\frac{n}{p} \leq r^c$  and  $r^c \geq B$ ,  $pB \log r^c$  if  $\frac{n}{p} > r^c$  and  $r^c \geq B$ , and  $pr^c \log r^c$  if  $r^c < B$ .

Since  $\text{BM}_M(n, r^c) \leq \text{BMT}(n, r^c) + 2\text{BM}_M(n, r^{c/2})$ , induction confirms the claimed bound.

Similar arguments bound the cache miss costs and idle time (see [13]). Adding these costs, plus those from Lemma 4 yields our main result.

**Theorem 1.** When run on  $p \leq \min\{\frac{n}{\log \log n}, \frac{n}{M}\}$  cores using the PWS scheduler on an input of size  $n \geq Mp$ , where  $M \geq B^2$ , the merging algorithm runs in time

$$O\left(\frac{n \log n}{p} + \frac{bn}{Bp} \frac{\log n}{\log M} + \log n \log \log n(s + S + b) + bB \log n \log \log \frac{n}{p}\right).$$

The same bound applies to the sorting algorithm.

## References

- [1] U. A. Acar, G. E. Blelloch, and R. D. Blumofe. The data locality of work stealing. *Theory of Computing Systems*, 35(3), 2002. Springer.
- [2] A. Aggarwal and J. S. Vitter. The input/output complexity of sorting and related problems. *CACM*, 31:1116–1127, 1988.
- [3] M. Ajtai, J. Komlos, and E. Szemerédi. An  $O(n \log n)$  sorting network. *Combinatorica*, 3:1–19, 1983.
- [4] L. Arge, M. T. Goodrich, M. Nelson, and N. Sitchinava. Fundamental parallel algorithms for private-cache chip multiprocessors. In *ACM SPAA*, pages 197–206, 2008.
- [5] G. Blelloch, R. Chowdhury, P. Gibbons, V. Ramachandran, S. Chen, and M. Kozuch. Provably good multicore cache performance for divide-and-conquer algorithms. In *ACM-SIAM SODA*, pages 501–510, 2008.
- [6] G. Blelloch, P. Gibbons, and H. Simhadri. Brief announcement: Low depth cache-oblivious sorting. In *ACM SPAA*. ACM, 2009.
- [7] R. Blumofe and C. E. Leiserson. Scheduling multithreaded computations by work stealing. *JACM*, pages 720–748, 1999.
- [8] R. D. Blumofe and C. E. Leiserson. Scheduling multithreaded computations by work stealing. *Journal of the ACM*, 46(5):720–748, 1999.
- [9] F. Burton and M. R. Sleep. Executing functional programs on a virtual tree of processors. In *ACM FPLCA*, pages 187–194, 1981.
- [10] R. A. Chowdhury, F. Silvestri, B. Blakeley, and V. Ramachandran. Oblivious algorithms for multicores and network of processors. In *IEEE IPDPS*, 2010.
- [11] R. Cole. Parallel merge sort. *SIAM J Comput*, 17(4), 1988.
- [12] R. Cole and V. Ramachandran. Efficient resource oblivious scheduling of multicore algorithms. Manuscript, 2010.
- [13] R. Cole and V. Ramachandran. Resource oblivious sorting on multicores. TR-10-13, Dept of Comp Sci, UT-Austin, 2010.
- [14] R. Cole and V. Ramachandran. Resource oblivious sorting on multicores. Submitted, 2010.
- [15] M. Frigo, C. E. Leiserson, H. Prokop, and S. Ramachandran. Cache-oblivious algorithms. In *IEEE FOCS*, 1999.
- [16] L. G. Valiant. A bridging model for multi-core computing. In *Proc. of the 16th Annual ESA*, volume 5193, pages 13–28, 2008.