

Cache-Oblivious Priority Queue and Graph Algorithm Applications*

Lars Arge[†] Michael A. Bender[‡] Erik D. Demaine[§]
Dept. of Computer Science Dept. of Computer Science Lab. for Computer Science
Duke University SUNY Stony Brook MIT
large@cs.duke.edu bender@cs.sunysb.edu edemaine@mit.edu

Bryan Holland-Minkley[¶] J. Ian Munro^{||}
Dept. of Computer Science Dept. of Computer Science
Duke University University of Waterloo
bhm@cs.duke.edu imunro@uwaterloo.ca

Submitted to SIAM Journal on Computing

May 16, 2003

Abstract

We develop an optimal cache-oblivious priority queue data structure, supporting insertion, deletion, and deletemin operations in $O(\frac{1}{B} \log_{M/B} \frac{N}{B})$ amortized memory transfers, where M and B are the memory and block transfer sizes of any two consecutive levels of a multilevel memory hierarchy. In a cache-oblivious data structure, M and B are not used in the description of the structure. Our structure is as efficient as several previously developed external memory (cache-aware) priority queue data structures, which all rely crucially on knowledge about M and B . Priority queues are a critical component in many of the best known external memory graph algorithms, and using our cache-oblivious priority queue we develop several cache-oblivious graph algorithms.

1 Introduction

As the memory systems of modern computers become more complex, it is increasingly important to design algorithms that are sensitive to the structure of memory. One of the essential features

* An extended abstract version of this paper was presented at the 2002 ACM Symposium on Theory of Computation (STOC'02) [12].

[†]Supported in part by the National Science Foundation through ESS grant EIA-9870734, RI grant EIA-9972879, CAREER grant CCR-9984099, and ITR grant EIA-0112849.

[‡]Supported in part by the National Science Foundation through ITR grant EIA-0112849 and by HRL Laboratories and Sandia National Laboratories.

[§]Supported in part by the National Science Foundation through ITR grant EIA-0112849.

[¶]Supported in part by the National Science Foundation through ITR grant EIA-0112849.

^{||}Supported in part by the Natural Science and Engineering Council through grant RGPIN 8237-97 and the Canada Research Chair in Algorithm Design.

of modern memory systems is that they are made up of a hierarchy of several levels of cache, main memory, and disk. While traditional theoretical computational models have assumed a “flat” memory with uniform access time, the access times of different levels of memory can vary by several orders of magnitude in current machines. For example, level-two cache is often around 100 times faster than main memory, while main memory is around 1,000,000 times faster than disks. In order to amortize the large access time of memory levels far away from the processor, memory systems often transfer data between memory levels in large blocks. Thus it is becoming increasingly important to obtain high data locality in memory access patterns.

The standard approach to obtaining good locality is to design algorithms parameterized by several aspects of the memory hierarchy, such as the size of each memory level, and the speed and block sizes of memory transfers between levels. Unfortunately, this parameterization often leads to complex algorithms that are tuned to particular architectures. As a result these algorithms are inflexible and not portable. To avoid the complexity of many parameters, a lot of research has been done on simpler *two-level* memory models. Recently, a promising new line of research has aimed at developing memory-hierarchy-sensitive algorithms that avoid any memory-specific parameterization whatsoever. It has been shown that if such a so-called *cache-oblivious* algorithm works optimally on a two-level hierarchy then it works optimally on *all* levels of any multilevel memory hierarchy—cache-oblivious algorithms automatically tune to arbitrary memory architectures. It seems surprising that data locality can be achieved without using the parameters describing the structure of the memory hierarchy, but nevertheless it is possible; cache-oblivious algorithms have been developed for fundamental problems such as sorting and searching.

In this paper we develop an optimal cache-oblivious priority queue. Previously known memory hierarchy efficient priority queues all rely crucially on parameters describing the hierarchy. We use the priority queue to develop several cache-oblivious graph algorithms. These are the first such algorithms and the complexity of most of our algorithms matches the complexity of the best known two-level cache-aware algorithms.

1.1 Background and previous results

Traditionally, most algorithmic work has been done in the *Random Access Machine* (RAM) model of computation, which models a “flat” memory with uniform access time. Recently, some attention has turned to the development of theoretical models and algorithms for modern complicated hierarchical memory systems; refer e.g. to [3, 4, 5, 7, 48, 54]. Developing models that are both simple and realistic is a challenging task since a memory hierarchy is described by many parameters. A typical hierarchy consists of a number of memory levels, with memory level ℓ having size M_ℓ and being composed of M_ℓ/B_ℓ blocks of size B_ℓ . In any memory transfer from level ℓ to $\ell - 1$, an entire block is moved atomically. Each memory level also has an associated *replacement strategies*, which is used to decide what block to remove in order to make room for a new block being brought into that level of the hierarchy. Further complications are the limited *associativity* of some levels of the hierarchy, meaning that a given block can only be loaded into a limited number of memory positions, as well as the complicated *prefetching strategies* deployed by many memory systems. In order to avoid the complications of multilevel memory models, a body of work has focused on two-level memory hierarchies; refer e.g. to [10, 54]. Most of this work has been done in the context of problems involving massive datasets, because the extremely long access times of disks compared to the other levels of the hierarchy means that I/O between main memory and disk often is the bottleneck in such problems.

1.1.1 Two-level I/O model

In the two-level *I/O model* (or *external memory model*) introduced by Aggarwal and Vitter [6], the memory hierarchy consists of an internal memory of size M , and an arbitrarily large external memory partitioned into blocks of size B . The efficiency of an algorithm in this model (a so-called *I/O* or *external memory* algorithm) is measured in terms of the number of block transfers it performs between these two levels (here called *memory transfers*). An algorithm has complete control over placement of blocks in main memory and on disk. The simplicity of the I/O model has resulted in the development of a large number of external memory algorithms and techniques. See e.g. [10, 54] for recent surveys.

The number of memory transfers needed to read N contiguous elements from disk is $scan(N) = \Theta(\frac{N}{B})$ (the *linear* or *scanning* bound). Aggarwal and Vitter[6] showed that $\Theta(\frac{N}{B} \log_{M/B} \frac{N}{B})$ memory transfers are necessary and sufficient to sort N elements. In this paper, we use $sort(N)$ to denote $\frac{N}{B} \log_{M/B} \frac{N}{B}$ (the *sorting* bound). The number of memory transfers needed to search for an element among a set of N elements is $\Omega(\log_B N)$ (the *searching* bound) and this bound is matched by the B-tree, which also supports updates in $O(\log_B N)$ memory transfers [15, 31, 35, 34]. An important consequence of these bounds is that, unlike in the RAM model, one cannot sort optimally with a search tree—inserting N elements in a B-tree takes $O(N \log_B N)$ memory transfers, which is a factor of $(B \log_B N) / (\log_{M/B} \frac{N}{B})$ from optimal. Finally, permuting N elements according to a given permutation takes $\Theta(\min\{N, sort(N)\})$ memory transfers and for all practical values of N, M and B this is $\Theta(sort(N))$ [6]. This represents another fundamental difference between the RAM and I/O model, since N elements can be permuted in $O(N)$ time in the RAM model.

The I/O model can also be used to model the two-level hierarchy between cache and main memory; refer e.g. to [49, 37, 38, 14, 44, 47]. Some of the main shortcomings of the I/O model on this level is the lack of explicit application control of placement of data in the cache and the low associativity of many caches. However, as demonstrated by Sen et al. [49], I/O model results can often be used to obtain results in more realistic two-level cache main memory models.

1.1.2 Cache-oblivious model

One of the main disadvantages of two-level memory models is that they force the algorithm designer to focus on a particular level of the hierarchy, resulting in algorithms that are less flexible to different-scale problems. Nevertheless, the I/O model has been successful because it is convenient to consider only two levels of the hierarchy. Very recently, a new model that combines the simplicity of the I/O mode with the realism of more complicated hierarchical models was introduced by Frigo et al. [33]. The idea in the *cache-oblivious model* is to design and analyze algorithms in the I/O model but without using the parameters M and B in the algorithm description. It is assumed that $M \geq B^2$ (the *tall cache* assumption) and that when an algorithm accesses an element that is not stored in main memory, the relevant block is automatically fetched into memory with a *memory transfer*. If the main memory is full, the *ideal* block in main memory is elected for replacement based on the future characteristics of the algorithm, that is, an *optimal* paging strategy is assumed. While this model may seem unrealistic, Frigo et al. [33] showed that it can be simulated by essentially any memory system with only a small constant-factor overhead. For example, the least-recently-used (LRU) block-replacement strategy approximates the omniscient strategy within a constant factor [33, 50]. The main advantage of the cache-oblivious model is that it allows us to reason about a simple two-level memory model, but prove results about an unknown, multilevel memory hierarchy; because an analysis of an algorithm in the two-level model holds for any block and main memory size, it holds for *any* level of the memory hierarchy. As a consequence, if the algorithm is

optimal in the two-level model, it is optimal on *all* levels of a multilevel memory hierarchy.

Frigo et al. [33] developed optimal cache-oblivious algorithms for matrix multiplication, matrix transposition, Fast Fourier Transform, and sorting. Subsequently several authors developed dynamic cache-oblivious B-trees with a search and update cost of $O(\log_B N)$ matching the standard (cache-aware) B-tree [18, 24, 19, 43, 17]. Recently, several further results have been obtained [20, 52, 21, 22, 23, 16, 2]

1.1.3 Priority queues

A priority queue maintains a set of elements each with a priority (or key) under the operations *insert*, *delete*, and *deletemin*, where a *deletemin* operation finds and deletes the minimum key element in the queue. The heap is a standard implementation of a priority queue and a balanced search tree can of course also easily be used to implement a priority queue. In the I/O model, a priority queue based on a B-tree would support all operations in $O(\log_B N)$ memory transfers. The standard heap can also be easily modified (to have fanout B) so that all operations are supported in the same bound (see e.g. [38]). The existence of a cache-oblivious B-tree immediately implies the existence of an $O(\log_B N)$ cache-oblivious priority queue.

As discussed, the use of an $O(\log_B N)$ search tree (or priority queue) to sort N elements results in an I/O model algorithm that is a factor of $(B \log_B N)/(\log_{M/B}(N/B))$ from optimal. To sort optimally we need a data structure supporting the relevant operations in $O(\frac{1}{B} \log_{M/B} \frac{N}{B})$ memory transfers. Note that for reasonable values of N , M , and B , this bound is less than 1 and we can therefore only obtain it in an amortized sense. To obtain such a bound, Arge developed the *buffer tree technique* [11]. The main idea in this technique is to perform operations in a lazy (or batched) manner using main memory sized buffers attached to the nodes of a data structure. Arge showed how to use the buffer tree technique on a B-tree in order to obtain a priority queue supporting all operations in $O(\frac{1}{B} \log_{M/B} \frac{N}{B})$ amortized memory transfers [11]. This structure seems hard to make cache-oblivious since it (apart from the memory sized buffers) involves periodically finding the $\Theta(M)$ smallest key elements in the structure and storing them in internal memory. Efficient I/O model priority queues have also been obtained by using the buffer technique on heap structures [32, 36]. The heap structure by Fadel et al. [32] seems hard to make cache-oblivious because it requires collecting $\Theta(M)$ insertions and performing them all on the structure at the same time. The structure by Kumar and Schwabe [36] avoids this by using the buffer technique on a tournament tree data structure. However, their structure still requires memory sized buffers. Apart from this they only obtained $O(\frac{1}{B} \log_2 N)$ bounds, mainly because their structure was designed to also support an update operation. Finally, Brodal and Katajainen [25] developed a priority queue structure based on a M/B -way merging scheme. Intuitively, this structure also seems hard to make cache-oblivious. However, very recently, and after the appearance of the conference version of this paper, Brodal and Fagerberg [22] managed to develop a merging based cache-oblivious priority queue based on ideas they developed in [21].

1.1.4 I/O model graph algorithms

The super-linear lower bound on permutation in the I/O model has important consequences for the I/O-complexity of graph algorithms, because the solution of almost any graph problem involves somehow permuting the V vertices or E edges of the graph. Thus $\Omega(\min\{V, \text{sort}(V)\})$ is in general a lower bound on the number of memory transfers needed to solve most graph problems. Refer to [9, 27, 40]. As mentioned, this bound is $\Omega(\text{sort}(V))$ in all practical cases. Still, even though a large number of I/O model graph algorithms have been developed (see [54, 55] and references

therein), not many algorithms match this bound. Below we review the results most relevant to our work.

Like for PRAM graph algorithms [45], list ranking—the problem of ranking the elements in a linked lists stored unordered in memory—is the most fundamental I/O model graph problem. Using PRAM techniques, Chiang et al. [27] developed the first efficient I/O model list ranking algorithm. Using an I/O-efficient priority queue, Arge [11] showed how to solve the problem in $O(\text{sort}(V))$ memory transfers. The list ranking algorithm and PRAM techniques can be used in the development of $O(\text{sort}(V))$ algorithms for most problems on trees, such as computing an Euler Tour, Breadth-First-Search (BFS), Depth-First-Search (DFS), and computing a centroid decomposition [27]. The best known DFS and BFS algorithms for general directed graphs use $O(V + \frac{EV}{BM})$ [27] or $O((V + E/B) \log_2 V + \text{sort}(E))$ [26] memory transfers. For undirected graphs improved $O(V + \text{sort}(E))$ and $O(\sqrt{\frac{V \cdot E}{B}} + \text{sort}(E))$ BFS algorithms have been developed [40, 39]. The best known algorithms for computing the connected components and the minimal spanning forest of a general undirected graph both use $O(\text{sort}(E) \cdot \log_2 \log_2(\frac{VB}{E}))$ or $O(V + \text{sort}(E))$ memory transfers [40, 13].

1.2 Our results

Problem	Our cache-oblivious result	Previous best cache-aware result
Priority queue	$O(\frac{1}{B} \log_{M/B} \frac{N}{B})$	$O(\frac{1}{B} \log_{M/B} \frac{N}{B})$ [11]
List ranking	$O(\text{sort}(V))$	$O(\text{sort}(V))$ [27, 11]
Tree algorithms	$O(\text{sort}(V))$	$O(\text{sort}(V))$ [27]
Directed BFS and DFS	$O((V + E/B) \log_2 V + \text{sort}(E))$	$O((V + E/B) \log_2 V + \text{sort}(E))$ [26]
		$O(V + \frac{EV}{BM})$ [27]
Undirected BFS	$O(V + \text{sort}(E))$	$O(V + \text{sort}(E))$ [40]
		$O(\sqrt{\frac{V \cdot E}{B}} + \text{sort}(E))$ [39]
Minimal spanning forest	$O(\text{sort}(E) \cdot \log_2 \log_2 V)$	$O(\text{sort}(E) \cdot \log_2 \log_2 \frac{VB}{E})$ [13]
	$O(V + \text{sort}(E))$	$O(v + \text{sort}(E))$ [13]

Figure 1: Summary of our results (Priority queue bounds are amortized).

The main result of this paper is an optimal cache-oblivious priority queue. Our structure supports insert, delete, and deletemin operations in $O(\frac{1}{B} \log_{M/B} \frac{N}{B})$ amortized memory transfers and $O(\log N)$ amortized computation time; it is described in Section 2. Our structure is based on a novel combination of several new ideas with ideas used in previous recursively defined cache-oblivious algorithms and data structures [33, 18], the buffer technique of Arge [11, 36], and the M/B -way merging scheme utilized by Brodal and Katajainen [25]. When the conference version of this paper appeared, our structure was the first cache-oblivious priority queue to obtain the same bounds as in the cache-aware case.

In the second part of the paper, Section 3, we use our priority queue to develop several cache-oblivious graph algorithms. Previously, no such algorithms were known. We first show how to solve the list ranking problem in $O(\text{sort}(V))$ memory transfers. Using this result we develop $O(\text{sort}(V))$ algorithms for fundamental problems on trees, such as the Euler Tour, BFS, and DFS problems. The complexity of all of these algorithms matches the complexity of the best known cache-aware

algorithms. Next we consider DFS and BFS on general graphs. Using modified versions of the data structures used in the $O((V + E/B) \log_2 V + \text{sort}(E))$ DFS and BFS algorithms for directed graphs [26], we make these algorithms cache-oblivious. We also discuss how the $O(V + \text{sort}(E))$ BFS algorithm for undirected graphs [40] can be made cache-oblivious. Finally, we develop two cache-oblivious algorithms for computing a minimal spanning forest (MSF), and thus also for computing connected components, of an undirected graph using $O(\text{sort}(E) \cdot \log_2 \log_2 V)$ and $O(V + \text{sort}(E))$ memory transfers, respectively. The two algorithms can be combined to compute the MSF in $O(\text{sort}(E) \cdot \log_2 \log_2 \frac{V}{V'} + V')$ memory transfers for any V' independent of B and M . Figure 1 summarizes our results. We believe our priority queue and the developed algorithms will prove useful in the development of many other cache-oblivious algorithms.

2 Priority Queue

In this section we describe our new optimal cache-oblivious priority queue. In Section 2.1 we define the data structure and in Section 2.2 we describe the supported operations.

2.1 Structure

2.1.1 Levels

Our priority queue data structure consists of $\Theta(\log \log N)$ levels whose sizes vary from N to some small size c beneath a constant threshold c_t . The size of a level corresponds (asymptotically) to the number of elements that can be stored within it. The i 'th level from above has size $N^{(2/3)^{i-1}}$ and for convenience we refer to the levels by their size. Thus the levels from largest to smallest are level N , level $N^{2/3}$, level $N^{4/9}$, \dots , level $X^{9/4}$, level $X^{3/2}$, level X , level $X^{2/3}$, level $X^{4/9}$, \dots , level $c^{9/4}$, level $c^{3/2}$, and level c . In general, a level can contain any number of elements less than or equal to its size, except level N , which always contains $\Theta(N)$ elements. Intuitively, smaller levels store elements with smaller keys or elements that were more recently inserted. In particular, the minimum key element and the most recently inserted element are always in the smallest (lowest) level c . Both insertions and deletions are initially performed on the smallest level and may propagate up through the levels.

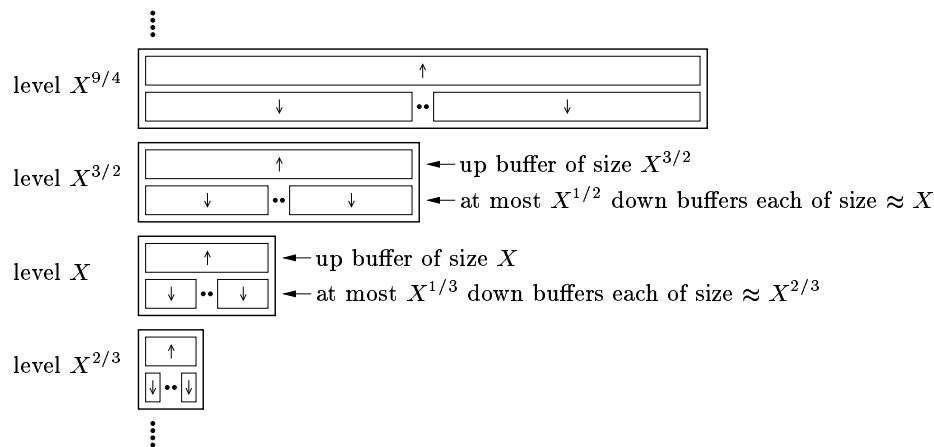


Figure 2: Levels $X^{2/3}$, X , $X^{3/2}$, and $X^{9/4}$ of the priority queue data structure.

2.1.2 Buffers

Elements are stored in a level in a number of *buffers*, which are also used to transfer elements between levels. Level X consists of one *up buffer* u^X that can store up to X elements, and at most $X^{1/3}$ *down buffers* $d_1^X, \dots, d_{X^{1/3}}^X$ each containing between $\frac{1}{2}X^{2/3}$ and $2X^{2/3}$ elements. Thus the maximum capacity of level X is $3X$. Refer to Figure 2. Note that the size of a down buffer at one level matches the size (up to a constant factor) of the up buffer one level down.

We maintain three invariants about the relationships between the elements in buffers of various levels:

Invariant 1 *At level X , elements are sorted among the down buffers, that is, elements in d_i^X have smaller keys than elements in d_{i+1}^X , but the elements within d_i^X are unordered.*

The element with largest key in each down buffer d_i^X is called a *pivot element*. Pivot elements mark the boundaries between the ranges of the keys of elements in down buffers.

Invariant 2 *At level X , the elements in the down buffers have smaller keys than the elements in the up buffer.*

Invariant 3 *The elements in the down buffers at level X have smaller keys than the elements in the down buffers at the next higher level $X^{3/2}$.*

The three invariants ensure that the keys of the elements in the down buffers get larger as we go from smaller to larger levels of the structure. Furthermore, an order exists between the buffers on one level: keys of elements in the up buffer are larger than keys of elements in down buffers. Therefore, down buffers are drawn below up buffers on Figure 2. However, the keys of the elements in an up buffer are unordered relative to the keys of the elements in down buffers one level up. Intuitively, up buffers store elements that are “on their way up”, that is, they have yet to be resolved as belonging to a particular down buffer in the next (or higher) level. Analogously, down buffers store elements that are “on their way down”—these elements are partitioned into several clusters so that we can quickly find the cluster of smallest key elements of size roughly equal to the next level down. In particular, the element with overall smallest key is in the first down buffer at level c .

2.1.3 Layout

We store the priority queue in a linear array as follows. The levels are stored consecutively from smallest to largest with each level occupying a single region of memory. For level X we reserve space for exactly $3X$ elements; X for the up buffer and $2X^{2/3}$ for each possible down buffer. The up buffer is stored first, followed by the down buffers stored in an arbitrary order but linked together to form an ordered linked list. Thus $3 \sum_{i=0}^{\log_{3/2} \log_c N} N^{(2/3)^i} = O(N)$ is an upper bound on the total size of the array.

Lemma 1 *A cache-oblivious priority queue on N elements uses $O(N)$ space.*

2.2 Operations

To implement the priority queue operations we will use two general operations, *push* and *pull*. Push inserts X elements into level $X^{3/2}$, and pull removes the X elements with smallest keys from

level $X^{3/2}$ and returns them in sorted order. This way, deletemin corresponds to pulling an element from the smallest level, and insert corresponds to pushing an element into the smallest level. More generally, whenever an up buffer on level X overflows we push the X elements in the buffer one level up, and whenever the down buffers on level X become too empty we pull X elements from one level up.

2.2.1 Push

To push X elements into level $X^{3/2}$, we first sort the X elements cache-obliviously using $O(1 + \frac{X}{B} \log_{M/B} \frac{X}{B})$ memory transfers and $O(X \log_2 X)$ time [33]. Next we distribute the elements in the sorted list into the $X^{1/2}$ down buffers of level $X^{3/2}$ by scanning through the list and simultaneously visiting the down buffers in (linked) order. More precisely, we append elements to the end of the current down buffer $d_i^{X^{3/2}}$, and advance to the next down buffer $d_{i+1}^{X^{3/2}}$ as soon as we encounter an element with larger key than the pivot of $d_i^{X^{3/2}}$. Elements with keys larger than the pivot of the last down buffer are inserted in the up buffer $u^{X^{3/2}}$. Scanning through the X elements take $O(1 + X/B)$ memory transfers and $O(X)$ time. Even though we do not scan through every down buffer, we perform at least one memory transfer for each of the $X^{1/2}$ buffers. Thus the total cost of distributing the X elements is $O(X/B + X^{1/2})$ memory transfers and $O(X + X^{1/2}) = O(X)$ time.

During the distribution of elements a down buffer may run full, that is, contain $2X$ elements. In this case, we split the buffer into two down buffers each containing X elements. We can perform the split in $O(1 + X/B)$ memory transfers and $O(X)$ time by first finding the median of the elements in the buffer in $O(1 + X/B)$ transfers and $O(X)$ time [33], and then partitioning the elements into the two new buffers in a simple scan. We place the new buffer in any free down buffer spot for the level and update the linked list accordingly. If the level already has the maximum number $X^{1/2}$ of down buffers, we first remove the last down buffer $d_{X^{1/3}}^X$ by inserting the less than $2X$ elements in $d_{X^{1/3}}^X$ into the up buffer. Because X elements must have been inserted since the last time the buffer split, the amortized splitting cost per element is $O(1/X + 1/B)$ transfers and $O(1)$ time. In total, the amortized number of memory transfers and time used on splitting buffers while distributing the X elements is $O(1 + X/B)$ and $O(X)$, respectively.

If the up buffer runs full during the above process, that is, contains more than $X^{3/2}$ elements, we recursively *push* all of these elements into the next level up. Note that after such a recursive push, $X^{3/2}$ elements have to be inserted (pushed) into the up buffer of level $X^{3/2}$ before another recursive push is needed.

The invariants are all maintained during a push of X elements into level $X^{3/2}$: As we sort the elements to distribute them among the down buffers, it is clear we maintain Invariant 1. Only elements larger than the pivot of the last down buffer are placed in the up buffer, so that Invariant 2 is maintained. Similarly, only elements smaller than the pivot of the last down buffer are placed in the down buffers, so Invariant 3 is maintained.

Ignoring the cost of recursive push operations for the moment we have:

Lemma 2 *A push of X elements from level X into level $X^{3/2}$ can be performed in $O(X^{1/2} + \frac{X}{B} \log_{M/B} \frac{X}{B})$ memory transfers and $O(X \log_2 X)$ time amortized, not counting the cost of any recursive push operations, while maintaining Invariants 1–3.*

2.2.2 Pull

To describe how to pull the X smallest keys elements from level $X^{3/2}$, we first assume that the down buffers contain at least $\frac{3}{2}X$ elements. In this case the first three down buffers $d_1^{X^{3/2}}$, $d_2^{X^{3/2}}$,

and $d_3^{X^{3/2}}$ together contain the smallest between $\frac{3}{2}X$ and $6X$ elements (Invariants 1 and 2). We find and remove the X smallest elements simply by sorting these elements using $O(1 + \frac{X}{B} \log_{M/B} \frac{X}{B})$ memory transfers and $O(X \log_2 X)$ time. The remaining between $X/2$ and $5X$ elements are left in one, two, or three down buffers of size between $X/2$ and $2X$. These buffers can easily be constructed in $O(1 + X/B)$ transfers and $O(X)$ time. Thus we use $O(1 + \frac{X}{B} \log_{M/B} \frac{X}{B})$ memory transfers and $O(X \log_2 X)$ time in total. It is easy to see that Invariants 1–3 are maintained.

In the case where the down buffers contain fewer than $\frac{3}{2}X$ elements, we first *pull* the $X^{3/2}$ elements with smallest keys from the next level up. Because these elements do not necessarily have smaller keys than the, say U , elements in the up buffer $u^{X^{3/2}}$, we then sort this up buffer and merge the two sorted lists. Then we insert the U elements with largest keys into the up buffer, and distribute the remaining between $X^{3/2}$ and $X^{3/2} + \frac{3}{2}X$ elements into $X^{1/2}$ down buffers of size between X and $X + \frac{3}{2}X^{1/2}$ each (such that the $O(1/X + 1/B)$ amortized down buffer split bound is maintained). It is easy to see that this maintains the three invariants. Afterwards, we can find the X minimal key elements as above. Note that after a recursive pull, $X^{3/2}$ elements have to be deleted (pulled) from the down buffers of level $X^{3/2}$ before another recursive pull is needed. Note also that a pull on level $X^{3/2}$ does not affect the number of elements in the up buffer $u^{X^{3/2}}$. Since we distribute elements into the down and up buffers after a recursive pull using one sort and one scan of $X^{3/2}$ element, the cost of doing so is dominated by the cost of the recursive pull operation itself. Ignoring these costs for the moment we have:

Lemma 3 *A pull of X elements from level $X^{3/2}$ down to level X can be performed in $O(1 + \frac{X}{B} \log_{M/B} \frac{X}{B})$ memory transfers and $O(X \log_2 X)$ time amortized, not counting the cost of any recursive pull operations, while maintaining Invariants 1–3.*

2.2.3 Total cost

As mentioned, an insert or a deletemin is performed by performing a push or pull, respectively, on the smallest level. This may require recursive pushes or pulls on higher levels. We now examine the total cost of an insert or deletemin, taking into account the recursive push or pull operations.

To maintain that our structure always uses $\Theta(N)$ space, we rebuild the structure bottom-up after every $N/2$ operations (often referred to as global rebuilding [41]): we find the largest value $c < c_t$ such that $c^{3/2^i} = N$ for some integer i . Then we build levels $c, c^{3/2}, \dots, N^{2/3}$ such that all up buffers are empty and such that level X has exactly $X^{1/3}$ down buffers of size $X^{2/3}$. The remaining $\Theta(N)$ elements are placed in level N such that it has $\Theta(N^{1/3})$ down buffers of size $N^{2/3}$. Because of the periodical rebuilding, the largest level of the structure always contains $\Theta(N)$ elements; it contains that many elements just after a rebuilding, and the structure will be rebuilt again before it can lose more than half of these elements. Thus the structure always contains $\Theta(\log \log N)$ levels. We can easily perform the rebuilding in a sorting and a scanning step using a total of $O(\frac{N}{B} \log_{M/B} \frac{N}{B})$ memory transfers and $O(N \log_2 N)$ time, or $O(\frac{1}{B} \log_{M/B} \frac{N}{B})$ transfers and $O(\log_2 N)$ time per operation amortized. Note that after a rebuilding, X elements have to be pushed into or pulled from level X before recursive pushes or pulls are needed.

To analyze the amortized cost of an insertion or deletemin operation we consider the total cost of $N/2$ such operations: we charge a push of X elements from level X up to level $X^{3/2}$ to level X . Since X elements have to be inserted in the up buffer u^X of level X between such pushes, and as elements can only be inserted in u^X when elements are inserted (pushed) into level X (pulls from level X and from level $X^{3/2}$ into level X do not affect the number of elements in u^X), $O(N/X)$ pushes are charged to level X during the $N/2$ operations. Similarly, we charge a pull of X elements from level $X^{3/2}$ down to level X to level X . Since between such pulls (at least) X elements have

to be deleted from the down buffers of level X by pulls on X (a push of elements into level X only increase the number of deletions needed and a push of elements from level X to level $X^{3/2}$ does not affect the number of elements in down buffers of level X), $O(N/X)$ pulls are charged to level X during the $N/2$ operations.

By Lemma 2 and 3 we know that a push or pull charged to level X uses $O(X^{1/2} + \frac{X}{B} \log_{M/B} \frac{X}{B})$ memory transfers. We can reduce this cost to $O(\frac{X}{B} \log_{M/B} \frac{X}{B})$ by more carefully examining the costs for differently sized levels. First consider a push or pull of $X \geq B^2$ elements into or from level $X^{3/2} \geq B^3$. In this case we trivially have that $O(X^{1/2} + \frac{X}{B} \log_{M/B} \frac{X}{B}) = O(\frac{X}{B} \log_{M/B} \frac{X}{B})$. If $B^{4/3} \leq X < B^2$, the $X^{1/2}$ term in the push bound can dominate and we have to analyze the cost of a push more carefully. In this case we are working on a level $X^{3/2}$ where $B^2 \leq X^{3/2} < B^3$. There is only one such level. Recall that the $X^{1/2}$ cost was from distributing X sorted elements into the less than $X^{1/2}$ down buffers of level $X^{3/2}$. More precisely, a block of each buffer may have to be loaded and written back without transferring a full block of elements into the buffer. However, because $X^{1/2} \leq B$ and $M = \Omega(B^2)$ (the tall-cache assumption), a block for each of the buffers can fit into main memory. Consequently, if a fraction of the main memory is used to keep a partially filled block of each buffer of level $X^{3/2}$ ($B^2 \leq X^{3/2} \leq B^3$) in memory at all times, and only full block are written to disk, the $X^{1/2}$ cost would be eliminated. In addition, if all of the levels of size less than B^2 (of total size $O(B^2)$) are also kept in memory, all transfer costs associated with them would be eliminated. The optimal paging strategy is able to keep the relevant blocks in memory at all times and thus eliminate these costs.

Since each of the $O(N/X)$ push and pull operations charged to level X ($X > B^2$) uses $O(\frac{X}{B} \log_{M/B} \frac{X}{B})$ memory transfers amortized, the total amortized transfer cost of an insert or deletemin operation in the sequence of $N/2$ such operations is $O(\sum_{i=0}^{\infty} \frac{1}{B} \log_{M/B} (N^{(2/3)^i}/B)) = O(\frac{1}{B} \log_{M/B} \frac{N}{B})$. Similarly the total amortized time cost is $O(\sum_{i=0}^{\infty} \log_2(N^{(2/3)^i})) = O(\log_2 N)$.

2.2.4 Supporting deletes

Using ideas from [11, 36] we can easily support a delete in $O(\frac{1}{B} \log_{M/B} \frac{N}{B})$ memory transfers and $O(\log_2 N)$ time amortized, provided that we are given the key of the element to be deleted. To perform a deletion, we simply insert a special element in the priority queue, with value and key equal to the element to be deleted. At some point during the sort performed during a push, pull, or rebuild, this special element and the element to be deleted will be compared, as they have the same key. When this happens we remove both elements from the structure. Note that the structure cannot contain more than a constant fraction of special elements or elements to be deleted, as all such elements will be compared and removed when we rebuild the structure.

To analyze a delete operation, we first note that it behaves exactly like an insertion, except that the special delete element and the element to be deleted are both removed when they “meet”. We incur the standard insertion cost when inserting a special delete element in the structures. The removal of the two elements on a level X eventually contributes to a pull operation being performed on level $X^{3/2}$; the cost incurred by the two deletions is upper bounded by the cost of two deletemin operations. Thus in total we have obtained the following:

Theorem 1 *Using $O(M)$ main memory, a set of N elements can be maintained in a linear-space cache-oblivious priority queue data structure supporting each insert, deletemin, and delete operation in $O(\frac{1}{B} \log_{M/B} \frac{N}{B})$ amortized memory transfers and $O(\log_2 N)$ amortized computation time.*

3 Graph Algorithms

In this section we discuss how our cache-oblivious priority-queue can be used to develop several cache-oblivious graph algorithms. We first consider the simple list ranking problem and algorithms on trees, and then we go on and consider BFS, DFS and minimal spanning tree algorithms for general graphs.

3.1 List ranking

In the list ranking problem we are given a linked list of V nodes stored as an unordered sequence. More precisely, we have an array with V nodes, each containing the position of the next node in the list (an edge). The goal is to determine the *rank* of each node v , that is, the number of edges from v to the end of the list. In a more general version of the problem, each edge has a weight and the goal is to find for each node v the sum of the weights of edges from v to the end of the list. Refer to Figure 3.

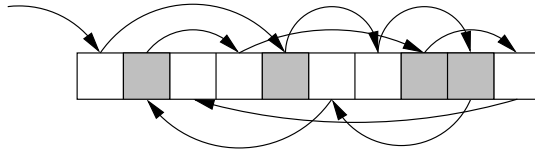


Figure 3: List ranking problem. An independent set of size 4 is marked. There are two forward lists (on top) and two backwards lists (on bottom).

Based on ideas from efficient PRAM algorithms [8, 29], Chiang et al. [27] designed an $O(\text{sort}(V))$ I/O model list ranking algorithm. The main idea in the algorithm is as follows. An *independent set* of $\Theta(V)$ nodes (nodes without edges to each other) is found, nodes in the independent set are “bridging out” (edges incident to nodes in this set are contracted), the remaining list is recursively ranked, and finally the contracted nodes are reintegrated into the list (their ranks are computed). The main innovation in the algorithm by Chiang et al. [27] was an $O(\text{sort}(V))$ memory transfer algorithm for computing an independent set of size V/c for some constant $c > 0$. The rest of the non-recursive steps of the algorithm can easily be performed in $O(\text{sort}(V))$ memory transfers using a few scans and sorts of the nodes of the list as follows. To bridge out the nodes in the independent set, we first identify nodes with a successor in the independent set. We do so by creating a copy of the list of nodes, sorting it by successor position, and simultaneously scanning the two lists. During this process, we can also mark each predecessor of an independent set node v with the position of the successor w of v , as well as with the weight of the edge (v, w) . Next we in a simple scan create a new list where the two edges incident to each independent set node v have been replaced with an edge from the predecessor of v to the successor of v . The new edge has weight equal to the sum of the two old edges. Finally, we create the list to be ranked recursively by removing the independent set nodes and “compressing” the remaining nodes, that is, storing them in an array of size $V(1 - 1/c)$. We do so by scanning through the list, while creating a list of the nodes not in the independent set, as well as a list that indicates the old and new position of each node (that is, the position of each node in the old and new array). Then we update the successor fields of the first list by sorting it by successor position, and simultaneously scanning it and the second list. After having ranked this list recursively, we reintegrate the removed nodes in the list, while computing their ranks. The rank of an independent set node v is simply the rank of its successor w minus the weight of edge (v, w) . The reintegration of the independent set nodes can be performed in a

few scans and sorts similar to the way we bridged out the independent set. Overall the number of memory transfers used to rank a V node list is $T(V) = O(\text{sort}(V)) + T(V/c) = O(\text{sort}(V))$.

All that remains in order to obtain a cache-oblivious list ranking algorithm is to develop a cache-oblivious independent set algorithm. Under different assumptions about the memory and block size, Chiang et al. [27] developed several independent set algorithms based on 3-coloring; in a 3-coloring, every node is colored with one of three colors such that adjacent nodes have different colors. The independent set (of size at least $V/3$) then consists of the set of nodes with the most popular color. Arge [11] and Kumar and Schwabe [36] later removed the main memory and block assumptions.

One way of computing a 3-coloring is as follows [11, 27]: We call an edge (v, w) a *forward* edge if v appears before w in the (unordered) sequence of nodes—otherwise it is called a *backward* edge. First we imagine splitting the list into two sets consisting of forward running segments (forward lists) and backward running segments (backward lists). Each node is included in at least one of these sets, and nodes at the head or tail of a segment (nodes at which there is a reversal of the direction) will be in both sets. Refer to Figure 3. Next we color the nodes in the forward lists red or blue by coloring the head nodes red and the other nodes alternately red and blue. Similarly, the nodes in the backward lists are colored green and blue, with the head nodes being colored green. In total, every node is colored with one color, except for the heads/tails, which have two colors. It is easy to see that we obtain a 3-coloring if we color each head/tail node red unless it was initially colored blue and green, in which case we color it green [27].

In the above algorithm we can cache-obliviously color the forward lists as follows (the backwards lists can be colored similarly). In a single scan we identify the head nodes and for each such node v we insert a red element in a cache-oblivious priority queue with key equal to the position of v in the unordered list. We then repeatedly extract the minimal key element e from the queue. If e corresponds to a node v , we access v in the list, color it the same color as e , and insert an element corresponding to its successor in the queue. The inserted element is colored in the opposite color of e . After processing all elements in the queue we have colored all forward lists. Since we use a cache-oblivious priority queue we can perform the $O(V)$ priority queue operations in $O(\text{sort}(V))$ memory transfers. Apart from this, we also perform what appears to be random accesses to the $O(V)$ nodes in the list. However, since we only process the forward list nodes in position order, the accesses overall end up corresponding to a scan of the list. Thus they only require $O(V/B)$ transfers. Thus overall we obtain the following:

Theorem 2 *The list ranking problem on a V node list can be solved cache-obliviously in $O(\text{sort}(V))$ memory transfers.*

3.2 Algorithms on trees

Many efficient PRAM algorithms on undirected trees use Euler Tour techniques [51, 53]. An Euler Tour of a graph is a cycle that traverses each edge exactly once. Not every graph has an Euler Tour but a tree where each undirected edge has been replaced with two directed edges does; refer to Figure 4. When we in the following refer to an Euler Tour of an undirected tree, we mean a tour in the graph obtained when replacing each edge in the tree with two directed edges.

To cache-obliviously compute an Euler Tour of an undirected tree, that is, to compute an ordered list of the edges along the tour, we use ideas from similar PRAM algorithms. Consider imposing a (any) cyclic order on the nodes adjacent to each node v in the tree. In [46] it is shown that an Euler Tour is obtained if we traverse the tree such that a visit to v from u (through the incoming edge (u, v)) is followed by a visit to the node w following u in the cyclic order (through

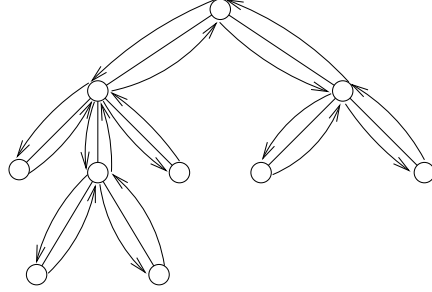


Figure 4: An undirected tree and an Euler Tour of the corresponding directed graph.

the outgoing edge (v, w)). Thus we can compute the successor edge of each edge e , that is, the edge following e in the Euler Tour, as follows: we first construct a list of incoming edges to each node v sorted according to the cyclic order. If two edges (u, v) and (w, v) are stored next to each other in this list, the successor edge for the (incoming) edge (u, v) is simply the (outgoing) edge (v, w) . Therefore we can compute all successor edges in a scan of the list. Given a list of all edges augmented with their successor edge, we can then compute the Euler Tour simply by list ranking the list and the sorting the edges by their rank. Thus overall we compute an Euler Tour of an undirected tree using $O(\text{sort}(V))$ memory transfers.

Using our Euler Tour algorithm, we can easily compute a depth-first search (DFS) numbering of the nodes of a tree starting at a source node s [46]. First note that if we start a walk of the Euler Tour in the source node s it is actually a DFS tour of the tree. To compute the numbering, we therefore first classify each edge as being either a *forward* or a *backward* edges; an edge is a forward edge if it is traversed before its reverse in the tour. After numbering the edges along the tour and sorting the list of edges such that reverse edges appear consecutively, we can classify each edge in a simple scan of the list. Then we assign each forward edge weight 1 and each backward edge weight 0. The DFS number of a node v is then simply the sum of the weights on the edges from s to v . Thus we can obtain the DFS numbering by solving the general version of list ranking. Since we only use Euler Tours computation, list ranking, sorting, and scanning, we in total compute the DFS numbering cache-obliviously in $O(\text{sort}(V))$ memory transfers.

Using an Euler Tour, list ranking, and sorting, we can also compute a breadth-first search (BFS) numbering of the nodes of a tree cache-obliviously in $O(\text{sort}(V))$ memory transfers in a similar way [46]. Using standard PRAM ideas, we can also e.g. compute the centroid decomposition of a tree in $O(\text{sort}(V))$ memory transfers [51, 53, 27]. The centroid of a tree is the node that, if removed, minimizes the size of the larger of the remaining subtrees. The centroid decomposition of a tree is a recursive partition of a tree into subtrees around the centroid.

Theorem 3 *The Euler Tour, BFS, DFS, and centroid decomposition problems on a tree with V nodes can be solved cache-obliviously in $O(\text{sort}(V))$ memory transfers.*

3.3 DFS and BFS

We now consider the DFS and BFS numbering problems for general graphs. We first describe a cache-oblivious DFS algorithm for directed graphs and then we modify it to compute a BFS numbering. Finally we develop an improved BFS algorithm for undirected graphs.

3.3.1 Depth-First Search

In the RAM model, directed DFS can be solved in linear time using a stack S containing vertices v that have not yet been visited but have an edge (w, v) incident to a visited vertex w , as well as an array A containing an element for each vertex v , indicating if v has been visited or not. The top vertex v of S is repeatedly popped, marked as visited in A , and all unvisited vertices adjacent to v are pushed onto S ; for each vertex adjacent to v an access is made to A to determine if the vertex has already been visited. It is easy to realize that if the stack S is implemented using a doubling array then a *push* or *pop* requires $O(1/B)$ cache-oblivious memory transfers amortized, since the optimal paging strategy can always keep the last block of the array (accessed by both push and pop) in main memory. However, each access to A may require a separate memory transfer resulting in $\Omega(E)$ memory transfers in total.

In the I/O model, Chiang et al. [27] modified the above algorithm to obtain an $O(V + \frac{E}{B} \frac{V}{M})$ algorithm. In their algorithm all visited vertices (marked vertices in array A) are stored in main memory. Every time the number of visited vertices grows larger than the main memory, all visited vertices and all their incident edges are removed from the graph. Since this algorithm relies crucially on knowledge of the main memory size, it seems hard to make it cache-oblivious. Buchsbaum et al. [26] described another $O((V + \frac{E}{B}) \log_2 V + \text{sort}(E))$ I/O model algorithm. This algorithm uses a number of data structures: V priority queues, a stack, and a so-called *buffered repository tree*. As discussed above, a stack can trivially be made cache-oblivious. In the following we first describe how to make the buffered repository tree cache-oblivious. Since our priority queue requires $O(M)$ space, we cannot immediately use our structure for the V priority queues. Instead we use a structure called a *buffered priority tree* similar to, and described after, the buffered repository tree. Finally, we discuss how the use of these structures leads to a cache-oblivious version the algorithm by Buchsbaum et al. [26].

Buffered repository tree. A buffered repository tree (BRT) maintains $O(E)$ elements with keys in the range $[1..V]$ under operations *insert* and *extract*. The insert operation inserts a new element, while the extract operation reports and deletes all elements with a certain key.

Our cache-oblivious version of the BRT consists of a static binary tree with the keys 1 through V in sorted order in the leaves. A buffer is associated with each node and leaf of the tree. The buffer of a leaf v contains elements with key v and the buffers of the internal nodes are used to perform insertions in a batched manner. We perform an insertion simply by inserting the new element into the root buffer. To perform an extraction of elements with key v we traverse the path from the root to the leaf containing v . At each node μ on this path we scan the associated buffer and report and delete elements with key v . During the scan we also distribute the remaining elements among the two buffers associated with the children of μ ; we distribute an element with key w to the buffer of the child of μ on the path to w . We place elements inserted in a buffer during the same scan consecutively in memory (but not necessarily right after the other elements in the buffer). This way the buffer of a node μ can be viewed as consisting of a linked list of *buckets* of elements in consecutive memory locations, with the number of buckets being equal to the number of buffer emptyings that have been performed on the parent of μ since the last emptying of μ 's buffer. To avoid performing a memory transfer on each insertion, we implement the root buffer slightly different, namely as a doubling array (like a stack). Since only the root buffer is implemented this way, the optimal paging strategy can keep the last block of the array in main memory and we obtain an $O(1/B)$ amortized root buffer insertion bound. Refer to Figure 5 for an illustration of a BRT.

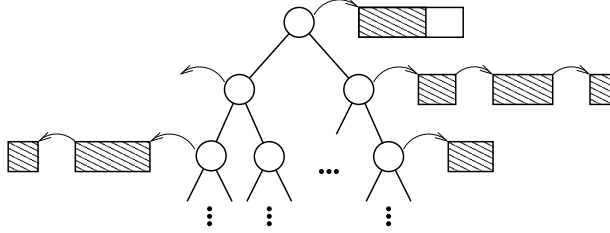


Figure 5: Buffered repository tree (BRT). Each non-root node has a buffer of elements stored as a linked list of buckets. The root node buffer is implemented using a doubling array.

Lemma 4 *A cache-oblivious buffered repository tree uses $\Theta(B)$ main memory space and supports insert and extract operations in $O(\frac{1}{B} \log_2 V)$ and $O(\log_2 V)$ memory transfers amortized, respectively.*

Proof: As discussed, an insertion in the root buffer requires $O(1/B)$ memory transfers amortized. During an extract operation, we use $O(X/B + K)$ memory transfers to empty the buffer of a node μ containing X elements in K buckets (since we access each element and each bucket). We charge the X/B -term to the insert operations that inserted the X elements in the BRT. Since each element is charged at most once on each level of the tree, an insert is charged $O(\frac{1}{B} \log_2 V)$ transfers. We charge the K -term to the extract operation that created the K buckets. Since an extract operation creates 2 buckets on each level of the tree, it is charged a total of $O(\log_2 V)$ memory transfers. \square

Buffered Priority Tree. A buffered priority tree is initially constructed on $O(E_v)$ elements with keys in the range $[1..V]$, and maintains these elements under *buffered delete* operations. Given E' elements currently in the structure, a buffered delete operation deletes the E' elements and reports and deletes the minimal element in the structure.

The buffered priority tree is implemented similarly to the buffered repository tree. It consists of a binary tree with keys 1 through V in sorted order in the leaves, with each node and leaf having a buffer associated with it. Initially, the $O(E_v)$ elements are stored in the buffers of the leaves. Unlike the buffered repository tree, the buffers are used to store elements intended to be deleted from the structure, and we maintain that each node v contains a count of the number of (live) elements stored in the tree rooted in v . To perform a buffered delete we first insert the E' elements in the buffer of the root and update the counter of the root. To find the minimal element, we then follow the leftmost path such that each node on the path has at least one live element beneath it. We can determine if a node is root in an empty tree using its counter. In each node we also scan the associated buffer and distribute the elements among the two buffers associated with its children, exactly as in the buffered repository tree. During the distribution, we also update the counters in the two children. When we reach a leaf l , we report and delete one of the elements stored in l 's buffer. Finally, we decrement the counters on the path from the root to l .

Lemma 5 *Using no permanent main memory, a cache-oblivious buffered priority tree supports buffered deletes of E' elements in $O((\frac{E'}{B} + 1) \log_2 V)$ memory transfers amortized. It can be constructed in $O(\text{sort}(E_v))$ memory transfers.*

Proof: The number of transfers needed to construct the tree is dominated by the $O(\text{sort}(E_v))$ memory transfers needed to sort the $O(E_v)$ elements and construct the leaf buffers. After that the tree can be constructed in $O(V/B)$ transfers level-by-level bottom-up. The amortized cost of a

buffered delete is equal to the cost of inserting E' elements into a BRT, plus the cost of extracting an element from a BRT, that is, $O((\frac{E'}{B} + 1) \log_2 V)$ memory transfers. \square

DFS algorithm. As mentioned, the directed DFS numbering algorithm by Buchsbaum et al. [26] utilizes a number of data structures: a stack S containing vertices on the path from the root of the DFS tree to the current vertex, a priority queue $P(v)$ for each vertex v containing edges (v, w) connecting v with a possibly unvisited vertex w , as well as one buffered repository tree D containing edges (v, w) incident to a vertex w that has already been visited but where (v, w) is still present in $P(v)$. The key of an edge (v, w) in D and $P(v)$ is w . For the priority queues $P(v)$ we use our buffered priority tree.

Initially each $P(v)$ is constructed on the E_v edges of the form (v, w) incident to v , the source vertex is placed on the stack S , and the BRT D is empty. To compute a DFS numbering, the vertex u on the top of the stack is repeatedly considered. All edges in D of the form (u, w) are extracted and deleted from $P(u)$ using a buffered delete operation. If $P(u)$ is now empty, so that no minimal element (edge (u, v)) is returned, all neighbors of u have already been visited and u is popped off S . Otherwise vertex v is visited next; it is numbered and pushed on S , and all edges (w, v) incident to it are inserted in D (since v is now visited). In [26] it is shown that this algorithm correctly computes a DFS numbering.

To analyze the above algorithm, we first note that each vertex is considered on the top of S a number of times equal to one greater than its number of children in the DFS tree. Thus the total number of times we consider a vertex on top of S is $2V - 1$. When considering a vertex u we first perform a stack operation on S , an extract operation on D , and a buffered delete operation on $P(u)$. The stack operation requires $O(1/B)$ memory transfers and the extraction $O(\log_2 V)$ transfers, since the optimal paging strategy can keep the relevant $O(1)$ blocks of these structures in main memory at all times. Thus these costs add up to $O(V \log_2 V)$. The buffered delete operation requires $O((1 + \frac{E'}{B}) \log_2 V)$ memory transfers if E' is the number of deleted elements (edges). Each edge is deleted once, so overall the buffered deletes costs add up to $O((V + \frac{E}{B}) \log_2 V)$. Next insert the, say E'' , edges incident to v in D . This requires $O(1 + \frac{E''}{B} \log_2 V)$ memory transfers, or $O(V + \frac{E}{B} \log_2 V)$ transfers over the whole algorithm. Note that the E'' edges are not immediately deleted directly from the relevant $P(w)$'s since that could cost a memory transfer per edge. In addition, the initial construction of the buffered priority trees requires $O(\text{sort}(E))$ memory transfers. Thus overall the algorithm uses $O((V + \frac{E}{B}) \log_2 V + \text{sort}(E))$ memory transfers.

3.3.2 Breadth-First Search

The DFS algorithm described above can be modified to perform a breadth-first search simply by replacing the stack S with a queue. Queues, like stacks, can be implemented using a doubling array, and the optimal paging strategy can keep the two partial blocks in use by the *enqueue* and *dequeue* operations in memory, such that each queue operation requires $O(1/B)$ memory transfers amortized. Thus we also obtain an $O((V + \frac{E}{B}) \log_2 V)$ directed BFS numbering algorithm.

Our directed DFS and BFS algorithms can of course also be used on undirected graphs. For undirected graphs, improved $O(V + \text{sort}(E))$ and $O(\sqrt{\frac{V \cdot E}{B}} + \text{sort}(E))$ I/O model algorithms have been developed [40, 39]. The idea in the algorithm by Munagala and Ranade [40], which can immediately be made cache-oblivious, is to visit the vertices in “layers” of vertices of equal distance from the source vertex s . The algorithm utilizes that in an undirected graph any vertex adjacent to a vertex in layer i is either in layer $i - 1$, layer i , or layer $i + 1$. It maintains two sorted lists of vertices in the last two layers i and $i - 1$. To create a sorted list of vertices in layer $i + 1$, a list of possible layer $i + 1$ vertices is first produced by collecting all vertices with a neighbor in level

i (using a scan of the adjacency lists of layer i vertices). Then this list is sorted, and in a scan of the list and the (sorted) lists of vertices in level i and $i - 1$ all previously visited vertices are removed. Apart from the sorting steps, overall this algorithm uses $O(V + E/B)$ memory transfers to access the edge lists for all vertices, as well as $O(E/B)$ transfers to scan the lists. Since each vertex is included in a sort once for each of its incident edges, the a total cost of all sorting steps is $O(\text{sort}(E))$. Thus in total the algorithm uses $O(V + \text{sort}(E))$ memory transfers. Since it only uses scans and sorts, it is cache-oblivious without modification. Refer to [40] for full details.

Theorem 4 *The DFS or BFS numbering of a directed graph can be computed cache-obliviously in $O((V + \frac{E}{B}) \log_2 V + \text{sort}(E))$ memory transfers. The BFS numbering of an undirected graph can be computed cache-obliviously in $O(V + \text{sort}(E))$ memory transfers.*

3.4 Minimal Spanning Forest

In this section we consider algorithms for computing the minimal spanning forest (MSF) of an undirected weighted graph. Without loss of generality, we assume that all edge weights are distinct. In the I/O model, a sequence of algorithms have been developed for the problem [27, 1, 36, 13], culminating in an algorithm using $O(\text{sort}(E) \cdot \log_2 \log_2(\frac{VB}{E}))$ memory transfers developed by Arge et al. [13]. This algorithm consists of two phases. In the first phase an edge contraction algorithm inspired by PRAM algorithms [28, 30] is used to reduce the number of vertices to $O(E/B)$. In the second phase a modified version of Prim’s algorithm [42] is used to finish the MSF computation. Using our cache-oblivious priority queue we can relatively easily modify both of the phases to work cache-obliviously. However, since we cannot decide cache-obliviously when the first phase has reduced the number of vertices to $O(E/B)$, we are not able to combine the two phases as effectively as in the I/O model. Below we first describe how to make the algorithms used in the two phases cache-oblivious. Then we discuss their combination.

3.4.1 Phase 1

The basic edge contraction based MSF algorithm proceeds in stages [28, 27, 36]. In each stage the minimum weight edge incident to each vertex is selected and output as part of the MSF, and the vertices connected by the selected edges are contracted into super-vertices (that is, the connected components of the graph of selected edges are contracted). See e.g. [13] for a proof that the selected edges along with the edges in a MSF of the contracted graph constitute a MSF for the original graph.

In the following we sketch how we can perform a contraction stage on a graph G cache-obliviously in $O(\text{sort}(E))$ memory transfers as in [13]. We can easily select the minimum weight edges in $O(\text{sort}(E))$ memory transfers using a few scans and sorts. To perform the contraction, we select a *leader vertex* in each connected component of the graph G_s of selected edges, and replace every edge (u, v) in G with the edge $(\text{leader}(u), \text{leader}(v))$. To select the leaders, we utilize that the connected components of G_s are trees, except that one edge in each component (namely the minimal weight edge) appears twice [13]. In each component, we simply use one of the vertices incident to the edge appearing twice as leader. This way we can easily identify all the leaders in $O(\text{sort}(E))$ memory transfers using a few sorts and scans (by identifying all edges that appear twice). We can then use our cache-oblivious tree algorithms developed in Section 3.2 to distribute the identity of the leader to each vertex in each component in $O(\text{sort}(V))$ memory transfers: we add an edge between each leader in G_s and a pseudo root vertex s and perform a DFS numbering of the resulting tree starting in s . Since all vertices in the same connected component (tree) will have consecutive DFS numbers,

we can then mark each vertex with its leader using a few sorts and scans. Finally, after marking each vertex v with $leader(v)$, we can easily replace each edge (u, v) in G with $(leader(u), leader(v))$ in $O(\text{sort}(E))$ memory transfers using a few sort and scan steps on the vertices and edges.

Since each contraction stage reduces the number of vertices by a factor of two, and since a stage is performed in $O(\text{sort}(E))$ memory transfers, we can reduce the number of vertices to $V' = V/2^i$ in $O(\text{sort}(E) \cdot \log_2(V/V'))$ memory transfers by performing i stages after each other. Thus we obtain an $O(\text{sort}(E) \cdot \log_2 V)$ algorithm by continuing the contraction until we are left with no edges. In the I/O model, Arge et al. [13] showed how to improve this bound to $O(\text{sort}(E) \cdot \log_2 \log_2 V)$ by grouping the stages into “super-stages” and only work on a subset of the edges of G in each super-stage. The extra steps involved in their improvement are all sorting or scanning of the edges and vertices, and therefore the improvement is immediately cache-oblivious.

Lemma 6 *The minimal spanning forest of an undirected weighted graph can be computed cache-obliviously in $O(\text{sort}(E) \cdot \log_2 \log_2 V)$ memory transfers.*

3.4.2 Phase 2

Prim’s algorithm [42] grows a minimal spanning tree (MST) of a connect graph iteratively from a source vertex using a priority queue P on the vertices not already included in the MST. The key of a vertex v in P is equal to the weight of the minimal weight edge connecting v to the current MST. In each step of the algorithm a deletemin is used to obtain the next vertex u to add to the MST, and the keys of all neighbors of u in P are (possibly) updated. A standard implementation of this algorithm uses $\Omega(E)$ memory transfers, since an transfer is needed to obtain the current key of each neighbor vertex. In the I/O model, Arge et al. [13] showed how to modify the algorithm to use $O(V + \text{sort}(E))$ memory transfers by storing edges in the priority queue rather than vertices. Below we describe this algorithm in order to show that it can be implemented cache-obliviously.

Like Prim’s algorithm, the algorithm by Arge et al. [13] grows the MST iteratively. During the algorithm, we maintain a priority queue P containing (at least) all edges connecting vertices in the current MST with vertices not in the tree; P can also contain edges between two vertices in the MST. Initially it contains all edges incident to the source vertex. In each step of the algorithm we extract the minimum weight edge (u, v) from P : if v is already in the MST we discard the edge; otherwise we include v in the MST and insert all edges incident to v , except (v, u) , in the priority queue. We can efficiently determine if v is already in the MST, since if u and v are both already in the MST then (u, v) must be in the priority queue twice; thus if the next edge we extract from P is also (u, v) then v is already in the MST.

The correctness of the above algorithm follows immediately from the correctness of Prim’s algorithm. During the algorithm we access the edges in the adjacency list of each vertex v once (when v is included in the MST) for a total of $O(V + E/B)$ memory transfers. We also perform $O(E)$ priority queue operations, for a total of $O(\text{sort}(E))$ memory transfers. Thus the algorithm uses $O(V + \text{sort}(E))$ memory transfers. Using standard techniques, all of the above can easily be modified to compute a MSF for a unconnected graph rather than a MST for a connected graph. Thus we have obtained the following.

Lemma 7 *The minimal spanning forest of an undirected weighted graph can be computed cache-obliviously in $O(V + \text{sort}(E))$ memory transfers.*

3.4.3 Combined algorithm

In the I/O model, an $O(\text{sort}(E) \cdot \log_2 \log_2(\frac{VB}{E}))$ MSF algorithm can be obtained by running the phase 1 algorithm until the number of vertices have been reduced to $V' = E/B$ using $O((\text{sort}(E) \cdot \log_2 \log_2(\frac{VB}{E}))$ memory transfers, and then finishing the MSF in $O(V' + \text{sort}(E)) = O(\text{sort}(E))$ memory transfers using the phase 2 algorithm. As mentioned, we cannot combine the two phases as effectively in the cache-oblivious model. In general however, we can combine the two algorithms to obtain an $O(\text{sort}(E) \cdot \log_2 \log_2(V/V') + V')$ algorithm for any V' independent of B and M .

Theorem 5 *The minimal spanning forest of an undirected weighted graph can be computed cache-obliviously in $O(\text{sort}(E) \cdot \log_2 \log_2(V/V') + V')$ memory transfers for any V' independent of B and M .*

4 Conclusions

In this paper, we presented an optimal cache-oblivious priority queue and used it to develop efficient cache-oblivious algorithms for several graph problems. We believe the ideas utilized in the development of the priority queue and our graph algorithms will prove useful in the development of other cache-oblivious data structures.

Many important problems still remains open in the area of cache-oblivious algorithms and data structures. In the area of graph algorithms, for example, it remains open to develop a cache-oblivious MSF algorithm with complexity matching the best known cache-aware algorithm. Cache-oblivious shortest path algorithms also still have to be developed.

References

- [1] J. Abello, A. L. Buchsbaum, and J. R. Westbrook. A functional approach to external graph algorithms. *Algorithmica*, 32(3):437–458, 2002.
- [2] P. K. Agarwal, L. Arge, A. Danner, and B. Holland-Minkley. Cache-oblivious data structures for orthogonal range searching. In *Proc. ACM Symposium on Computational Geometry*, 2003.
- [3] A. Aggarwal, B. Alpern, A. K. Chandra, and M. Snir. A model for hierarchical memory. In *Proc. ACM Symposium on Theory of Computation*, pages 305–314, 1987.
- [4] A. Aggarwal and A. K. Chandra. Virtual memory algorithms. In *Proc. ACM Symposium on Theory of Computation*, pages 173–185, 1988.
- [5] A. Aggarwal, A. K. Chandra, and M. Snir. Hierarchical memory with block transfer. In *Proc. IEEE Symposium on Foundations of Computer Science*, pages 204–216, 1987.
- [6] A. Aggarwal and J. S. Vitter. The Input/Output complexity of sorting and related problems. *Communications of the ACM*, 31(9):1116–1127, 1988.
- [7] B. Alpern, L. Carter, E. Feig, and T. Selker. The uniform memory hierarchy model of computation. *Algorithmica*, 12(2-3), 1994.
- [8] R. J. Anderson and G. L. Miller. A simple randomized parallel algorithm for list-ranking. *Information Processing Letters*, 33:269–273, 1990.

- [9] L. Arge. The I/O-complexity of ordered binary-decision diagram manipulation. In *Proc. International Symposium on Algorithms and Computation, LNCS 1004*, pages 82–91, 1995. A complete version appear as BRICS Technical Report RS-96-29, University of Aarhus.
- [10] L. Arge. External memory data structures. In J. Abello, P. M. Pardalos, and M. G. C. Resende, editors, *Handbook of Massive Data Sets*, pages 313–358. Kluwer Academic Publishers, 2002.
- [11] L. Arge. The buffer tree: A technique for designing batched external data structures. *Algorithmica*, 2003. To appear.
- [12] L. Arge, M. Bender, E. Demaine, B. Holland-Minkley, and J. I. Munro. Cache-oblivious priority-queue and graph algorithms. In *Proc. ACM Symposium on Theory of Computation*, pages 268–276, 2002.
- [13] L. Arge, G. S. Brodal, and L. Toma. On external memory MST, SSSP and multi-way planar graph separation. In *Proc. Scandinavian Workshop on Algorithms Theory, LNCS 1851*, pages 433–447, 2000.
- [14] L. Arge, J. Chase, J. Vitter, and R. Wickremesinghe. Efficient sorting using registers and caches. *ACM Journal on Experimental Algorithmics*, 7(9), 2002.
- [15] R. Bayer and E. McCreight. Organization and maintenance of large ordered indexes. *Acta Informatica*, 1:173–189, 1972.
- [16] M. Bender, R. Cole, E. Demaine, and M. Farach-Colton. Scanning and traversing: Maintaining data for traversals in memory hierarchy. In *Proc. European Symposium on Algorithms*, pages 152–164, 2002.
- [17] M. A. Bender, R. Cole, and R. Raman. Exponential structures for cache-oblivious algorithms. In *Proc. International Colloquium on Automata, Languages, and Programming*, pages 195–207, 2002.
- [18] M. A. Bender, E. D. Demaine, and M. Farach-Colton. Cache-oblivious B-trees. In *Proc. IEEE Symposium on Foundations of Computer Science*, pages 339–409, 2000.
- [19] M. A. Bender, Z. Duan, J. Iacono, and J. Wu. A locality-preserving cache-oblivious dynamic dictionary. In *Proc. ACM-SIAM Symposium on Discrete Algorithms*, pages 29–38, 2002.
- [20] R. D. Blumofe, M. Frigo, C. F. Joerg, C. E. Leiserson, and K. H. Randall. An analysis of dag-consistent distributed shared-memory algorithms. In *Proc. ACM Symposium on Parallel Algorithms and Architectures*, pages 297–308, 1996.
- [21] G. S. Brodal and R. Fagerberg. Cache oblivious distribution sweeping. In *Proc. International Colloquium on Automata, Languages, and Programming*, pages 426–438, 2002.
- [22] G. S. Brodal and R. Fagerberg. Funnel heap - a cache oblivious priority queue. In *Proc. International Symposium on Algorithms and Computation, LNCS 2518*, pages 219–228, 2002.
- [23] G. S. Brodal and R. Fagerberg. On the limits of cache-obliviousness. In *Proc. ACM Symposium on Theory of Computation*, 2003.
- [24] G. S. Brodal, R. Fagerberg, and R. Jacob. Cache oblivious search trees via binary trees of small height. In *Proc. ACM-SIAM Symposium on Discrete Algorithms*, pages 39–48, 2002.

- [25] G. S. Brodal and J. Katajainen. Worst-case efficient external-memory priority queues. In *Proc. Scandinavian Workshop on Algorithms Theory, LNCS 1432*, pages 107–118, 1998.
- [26] A. L. Buchsbaum, M. Goldwasser, S. Venkatasubramanian, and J. R. Westbrook. On external memory graph traversal. In *Proc. ACM-SIAM Symposium on Discrete Algorithms*, pages 859–860, 2000.
- [27] Y.-J. Chiang, M. T. Goodrich, E. F. Grove, R. Tamassia, D. E. Vengroff, and J. S. Vitter. External-memory graph algorithms. In *Proc. ACM-SIAM Symposium on Discrete Algorithms*, pages 139–149, 1995.
- [28] F. Chin, J. Lam, and I. Chen. Efficient parallel algorithms for some graph problems. *Communications of ACM*, 25:659–665, 1982.
- [29] R. Cole and U. Vishkin. Deterministic coin tossing with applications to optimal list-ranking. *Information and Control*, 70(1):32–53, 1986.
- [30] R. Cole and U. Vishkin. Approximate parallel scheduling. II. Applications to logarithmic-time optimal parallel graph algorithms. *Information and Computation*, 92(1):1–47, 1991.
- [31] D. Comer. The ubiquitous B-tree. *ACM Computing Surveys*, 11(2):121–137, 1979.
- [32] R. Fadel, K. V. Jakobsen, J. Katajainen, and J. Teuhola. Heaps and heapsort on secondary storage. *Theoretical Computer Science*, 220(2):345–362, 1999.
- [33] M. Frigo, C. E. Leiserson, H. Prokop, and S. Ramachandran. Cache-oblivious algorithms. In *Proc. IEEE Symposium on Foundations of Computer Science*, pages 285–298, 1999.
- [34] S. Huddleston and K. Mehlhorn. A new data structure for representing sorted lists. *Acta Informatica*, 17:157–184, 1982.
- [35] D. E. Knuth. *Sorting and Searching*, volume 3 of *The Art of Computer Programming*. Addison-Wesley, Reading MA, second edition, 1998.
- [36] V. Kumar and E. Schwabe. Improved algorithms and data structures for solving graph problems in external memory. In *Proc. IEEE Symp. on Parallel and Distributed Processing*, pages 169–177, 1996.
- [37] R. Ladner, J. Fix, and A. LaMarca. Cache performance analysis of traversals and random accesses. In *Proc. ACM-SIAM Symposium on Discrete Algorithms*, pages 613–622, 1999.
- [38] A. LaMarca and R. E. Ladner. The influence of caches on the performance of heaps. *Journal of Experimental Algorithmics*, 1(4), 1996.
- [39] U. Meyer. External memory bfs on undirected graphs with bounded degree. In *Proc. ACM-SIAM Symposium on Discrete Algorithms*, pages 87–88, 2001.
- [40] K. Munagala and A. Ranade. I/O-complexity of graph algorithms. In *Proc. ACM-SIAM Symposium on Discrete Algorithms*, pages 687–694, 1999.
- [41] M. H. Overmars. *The Design of Dynamic Data Structures*. Springer-Verlag, LNCS 156, 1983.
- [42] R. C. Prim. Shortest connection networks and some generalizations. *Bell Syst. Tech. J.*, 36:1389–1401, 1957.

- [43] N. Rahman, R. Cole, and R. Raman. Optimized predecessor data structures for internal memory. In *Proc. Workshop on Algorithm Engineering, LNCS 2141*, pages 67–78, 2001.
- [44] N. Rahman and R. Raman. Analysing cache effects in distribution sorting. In *Proc. Workshop on Algorithm Engineering*, 1999.
- [45] J. H. Reif, editor. *Synthesis of Parallel Algorithms*. Morgan Kaufmann, 1993.
- [46] J. H. Reif, editor. *Synthesis of Parallel Algorithms*, chapter 2, pages 61–114. Morgan Kaufmann, 1993.
- [47] P. Sanders. Fast priority queues for cached memory. In *Proc. Workshop on Algorithm Engineering and Experimentation, LNCS 1619*, pages 312–327, 1999.
- [48] J. E. Savage. Extending the Hong-Kung model to memory hierarchies. In *Proc. Annual International Conference on Computing and Combinatorics, LNCS 959*, pages 270–281, 1995.
- [49] S. Sen, S. Chatterjee, and N. Dumir. Towards a theory of cache-efficient algorithms. *Journal of the ACM*, 49(6):828–858, 2002.
- [50] D. D. Sleator and R. E. Tarjan. Amortized efficiency of list update and paging rules. *Commun. ACM*, 28:202–208, 1985.
- [51] R. E. Tarjan and U. Vishkin. An efficient parallel biconnectivity algorithm. *SIAM Journal on Computing*, 14(4):862–874, 1985.
- [52] S. Toledo. Locality of reference in *LU* decomposition with partial pivoting. *SIAM Journal on Matrix Analysis and Applications*, 18(4):1065–1081, 1997.
- [53] U. Vishkin. On efficient parallel strong orientation. *Information Processing Letters*, 20:235–240, 1985.
- [54] J. S. Vitter. External memory algorithms and data structures: Dealing with MASSIVE data. *ACM Computing Surveys*, 33(2):209–271, 2001.
- [55] N. Zeh. *I/O-Efficient Algorithms for Shortest Path Related Problems*. PhD thesis, School of Computer Science, Carleton University, 2002.