

# External Memory Priority Queues with Decrease-Key and Applications to Graph Algorithms

John Iacono 

Department of Computer Science, Université Libre de Bruxelles, Belgium  
<http://johniacono.com/>  
ulb@johniacono.com

Riko Jacob

Computer Science Department, IT University of Copenhagen, Denmark  
<http://www.itu.dk/people/rikj/>  
rikj@itu.dk

Konstantinos Tsakalidis 

Department of Computer Science, University of Liverpool, United Kingdom  
<https://cgi.csc.liv.ac.uk/~tsakalid/>  
K.Tsakalidis@liverpool.ac.uk

---

## Abstract

We present priority queues in the external memory model with block size  $B$  and main memory size  $M$  that support on  $N$  elements, operation UPDATE (a combination of operations INSERT and DECREASEKEY) in  $O\left(\frac{1}{B} \log_{\frac{M}{B}} \frac{N}{B}\right)$  amortized I/Os and operations EXTRACTMIN and DELETE in  $O\left(\lceil \frac{M^\varepsilon}{B} \log_{\frac{M}{B}} \frac{N}{B} \rceil \log_{\frac{M}{B}} \frac{N}{B}\right)$  amortized I/Os, for any real  $\varepsilon \in (0, 1)$ , using  $O\left(\frac{N}{B} \log_{\frac{M}{B}} \frac{N}{B}\right)$  blocks. Previous I/O-efficient priority queues either support these operations in  $O\left(\frac{1}{B} \log_2 \frac{N}{B}\right)$  amortized I/Os [Kumar and Schwabe, SPDP '96] or support only operations INSERT, DELETE and EXTRACTMIN in optimal  $O\left(\frac{1}{B} \log_{\frac{M}{B}} \frac{N}{B}\right)$  amortized I/Os, however without supporting DECREASEKEY [Fadel et al., TCS '99].

We also present buffered repository trees that support on a multi-set of  $N$  elements, operation INSERT in  $O\left(\frac{1}{B} \log_{\frac{M}{B}} \frac{N}{B}\right)$  I/Os and operation EXTRACT on  $K$  extracted elements in  $O\left(M^\varepsilon \log_{\frac{M}{B}} \frac{N}{B} + K/B\right)$  amortized I/Os, using  $O\left(\frac{N}{B}\right)$  blocks. Previous results achieve  $O\left(\frac{1}{B} \log_2 \frac{N}{B}\right)$  I/Os and  $O\left(\log_2 \frac{N}{B} + \frac{K}{B}\right)$  I/Os, respectively [Buchsbaum et al., SODA '00].

Our results imply improved  $O\left(\frac{E}{B} \log_{\frac{M}{B}} \frac{E}{B}\right)$  I/Os for single-source shortest paths, depth-first search and breadth-first search algorithms on massive directed dense graphs  $(V, E)$  with  $E = \Omega(V^{1+\varepsilon})$ ,  $\varepsilon > 0$  and  $V = \Omega(M)$ , which is equal to the I/O-optimal bound for sorting  $E$  values in external memory.

**2012 ACM Subject Classification** Theory of computation  $\rightarrow$  Graph algorithms analysis; Theory of computation  $\rightarrow$  Shortest paths; Theory of computation  $\rightarrow$  Data structures design and analysis; Hardware  $\rightarrow$  External storage

**Keywords and phrases** priority queues, external memory, graph algorithms, shortest paths, depth-first search, breadth-first search

**Digital Object Identifier** 10.4230/LIPIcs.ESA.2019.60

**Related Version** A full version of the paper is available at <https://arxiv.org/abs/1903.03147>.

**Funding** *John Iacono*: Supported by Fonds de la Recherche Scientifique-FNRS under Grant no MISU F 6001 1. Partially supported by NSF grants CCF-1319648 and CCF-1533564.

*Konstantinos Tsakalidis*: Partially supported by NSF grants CCF-1319648 and CCF-1533564.



© John Iacono, Riko Jacob, and Konstantinos Tsakalidis;  
licensed under Creative Commons License CC-BY

27th Annual European Symposium on Algorithms (ESA 2019).

Editors: Michael A. Bender, Ola Svensson, and Grzegorz Herman; Article No. 60; pp. 60:1–60:14

Leibniz International Proceedings in Informatics



Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

## 1 Introduction

Priority queues are fundamental data structures with numerous applications across computer science, most prominently in the design of efficient graph algorithms. They support the following operations on a set of  $N$  stored *elements* of the type  $(key, priority)$ , where “key” serves as an identifier and “priority” is a value from a total order:

- INSERT(element  $e$ ): Insert element  $e$  to the priority queue.
- DELETE(key  $k$ ): Remove all elements with key  $k$  from the priority queue.
- element  $e = \text{EXTRACTMIN}()$ : Remove and return the element  $e$  in the priority queue with the smallest priority.
- DECREASEKEY(element  $(k, p)$ ): Given that an element with key  $k$  and priority  $p'$  is stored in the priority queue, if priority  $p < p'$ , replace the element’s priority  $p'$  with  $p$ .

Operation UPDATE(element  $(k, p)$ ) is a combination of operations INSERT and DECREASEKEY, such that if the priority queue does not contain any element with key  $k$ , INSERT( $(k, p)$ ) is executed, otherwise DECREASEKEY( $(k, p)$ ) is executed.

We study the problem of designing priority queues that support all these operations in external memory. In the *external memory model* (also known as the *I/O model*) [1] the amount of input data is assumed to be much larger than the main memory size  $M$ . Thus, the data is stored in an external memory device (i.e. disk) that is divided into consecutive blocks of size  $B$  elements. Time complexity is measured in terms of *I/O operations* (or *I/Os*), namely block transfers from external to main memory and vice versa, while computation in main memory is considered to be free. Space complexity is measured in the number of blocks occupied by the input data in external memory. Algorithms and data structures in this model are considered *cache-aware*, since they are parameterized in terms of  $M$  and  $B$ . In contrast, *cache-oblivious* algorithms and data structures [11] are oblivious to both these values, which allows them to be efficient along all levels of a memory hierarchy. I/O-optimally scanning and sorting  $x$  consecutive elements in an array are commonly denoted to take  $\text{Scan}(x) = O\left(\frac{x}{B}\right)$  I/Os and  $\text{Sort}(x) = O\left(\frac{x}{B} \log_{\frac{M}{B}} \frac{x}{B}\right)$  I/Os, respectively [1, 11].

Priority queues are a basic component in several fundamental graph algorithms, including:

- The *single-source shortest paths (SSSP)* algorithm on directed graphs with positively weighted edges, which computes the minimum edge-weight paths from a given source node to all other nodes in the graph.
- The *depth-first search (DFS)* and *breadth-first search (BFS)* algorithms on directed unweighted graphs, which number all nodes of the graph according to a depth-first or a breadth-first exploration traversal starting from a given source node, respectively.

Another necessary component for these algorithms are I/O-efficient *buffered repository trees (BRTs)* [6, 2, 7]. They are used by external memory graph algorithms in order to confirm that a given node has already been visited by the algorithm. This avoids expensive random-access I/Os incurred by internal memory methods. In particular, BRTs support the following operations on a stored multi-set of  $N$  (key, value) elements, where “key” serves as an identifier and “value” is from a total order:

- INSERT(element  $e$ ): Insert element  $e$  to the BRT.
- element  $e_i = \text{EXTRACT}(\text{key } k)$ : Remove and return all  $K$  elements  $e_i$  (for  $i \in [1, K]$ ) in the BRT with key  $k$ .

■ **Table 1** Asymptotic amortized I/O-bounds of cache-aware and cache-oblivious priority queue operations (respectively, above and below the horizontal line) on  $N$  elements and real  $\varepsilon \in (0, 1)$ . \*Expected I/Os.

	INSERT	DELETE	EXTRACTMIN	DECREASEKEY
[10]	$\frac{1}{B} \log_{\frac{M}{B}} \frac{N}{B}$	$\frac{1}{B} \log_{\frac{M}{B}} \frac{N}{B}$	$\frac{1}{B} \log_{\frac{M}{B}} \frac{N}{B}$	–
[14]	$\frac{1}{B} \log_2 \frac{N}{B}$	$\frac{1}{B} \log_2 \frac{N}{B}$	$\frac{1}{B} \log_2 \frac{N}{B}$	$\frac{1}{B} \log_2 \frac{N}{B}$
[13]*	$\frac{1}{B} \log_{\log N} \frac{N}{B}$	$\frac{1}{B} \log_{\log N} \frac{N}{B}$	$\frac{1}{B} \log_{\log N} \frac{N}{B}$	$\frac{1}{B} \log_{\log N} \frac{N}{B}$
New	$\frac{1}{B} \log_{\frac{M}{B}} \frac{N}{B}$	$\lceil \frac{M^\varepsilon}{B} \log_{\frac{M}{B}} \frac{N}{B} \rceil \log_{\frac{M}{B}} \frac{N}{B}$	$\lceil \frac{M^\varepsilon}{B} \log_{\frac{M}{B}} \frac{N}{B} \rceil \log_{\frac{M}{B}} \frac{N}{B}$	$\frac{1}{B} \log_{\frac{M}{B}} \frac{N}{B}$
[2]	$\frac{1}{B} \log_{\frac{M}{B}} N$	$\frac{1}{B} \log_{\frac{M}{B}} N$	$\frac{1}{B} \log_{\frac{M}{B}} N$	–
[5, 7]	$\frac{1}{B} \log_2 N$	$\frac{1}{B} \log_2 N$	$\frac{1}{B} \log_2 N$	$\frac{1}{B} \log_2 N$

### 1.1 Previous work

Designing efficient external memory priority queues able to support operation DECREASEKEY (or at least operation UPDATE) has been a long-standing open problem [14, 10, 16, 9, 7, 13]. I/O-efficient adaptations of the standard heap data structure [10] or other sorting-based approaches [16], despite achieving optimal base- $(M/B)$  logarithmic amortized I/O-complexity, fail to support operation DECREASEKEY. (Nevertheless, we use these priority queues as subroutines in our structure.) Adaptations of the tournament tree data structure support all operations, albeit in not so efficient base-2 logarithmic amortized I/Os [14, 7]. Indeed, in the recent work of Eenberg, Larsen and Yu [9] it is shown that for a sequence of  $N$  operations, any external-memory priority queue supporting DECREASEKEY must spend  $\max\{\text{INSERT}, \text{DELETE}, \text{EXTRACTMIN}, \text{DECREASEKEY}\} = \Omega\left(\frac{1}{B} \log_{\log N} B\right)$  amortized I/Os. Randomized priority queues with matching complexity were recently presented by Jiang and Larsen [13].

The BRTs introduced by Buchsbaum et al. [6, Lemma 2.1] and their cache-oblivious counterparts [2] support INSERT in  $O\left(\frac{1}{B} \log_2 N\right)$  amortized I/Os and EXTRACT on  $K$  extracted elements in  $O(\log_2 N + K/B)$  amortized I/Os on a multi-set of  $N$  stored elements.

### 1.2 Our contributions

We present I/O-efficient priority queues that support on  $N$  stored elements, operation UPDATE in optimal  $O\left(\frac{1}{B} \log_{\frac{M}{B}} \frac{N}{B}\right)$  amortized I/Os and operations EXTRACTMIN and DELETE in  $O\left(\lceil \frac{M^\varepsilon}{B} \log_{\frac{M}{B}} \frac{N}{B} \rceil \log_{\frac{M}{B}} \frac{N}{B}\right)$ , for any real  $\varepsilon \in (0, 1)$ . Our priority queues are the first to support operation UPDATE (and thus DECREASEKEY) in a cache-aware setting in optimal I/Os, while also I/O-optimally supporting operation INSERT. These bounds improve upon previous priority queues supporting DECREASEKEY [14], albeit at the expense of suboptimal I/O-efficiency for EXTRACTMIN and DELETE (respecting the lower bound of [9] for  $M = \Omega(B \log_2 N)$ ). See Table 1 for a comparison with previous cache-aware and cache-oblivious I/O-efficient priority queues.

We also present I/O-efficient BRTs that support on a multi-set of  $N$  elements, operation INSERT in  $O\left(\frac{1}{B} \log_{\frac{M}{B}} \frac{N}{B}\right)$  amortized I/Os and operation EXTRACT on  $K$  extracted elements in  $O\left(M^\varepsilon \log_{\frac{M}{B}} \frac{N}{B} + K/B\right)$  amortized I/Os. Previous cache-aware bounds were  $O\left(\frac{1}{B} \log_2 \frac{N}{B}\right)$  and  $O\left(\log_2 \frac{N}{B} + \frac{K}{B}\right)$ , respectively [6]. Combined with our priority queues, for external memory SSSP, DFS and BFS algorithms on graphs with  $V$  nodes and  $E$  directed edges, we achieve  $O\left(V \frac{M^{\frac{\alpha}{1+\alpha}}}{B} \log_{\frac{M}{B}}^2 \frac{E}{B} + V \log_{\frac{M}{B}} \frac{E}{B} + \frac{E}{B} \log_{\frac{M}{B}} \frac{E}{B}\right)$  I/Os. This compares to previous

$O\left(\left(V + \frac{E}{B}\right) \log_2 E\right)$  I/Os for directed SSSP [14, 15, 7] and  $O\left(\left(V + \frac{E}{B}\right) \log_2 \frac{V}{B} + \frac{E}{B} \log_{\frac{M}{B}} \frac{E}{B}\right)$  I/Os for directed DFS and BFS [6, 2]. Our bounds are I/O-optimal for dense graphs with  $E = \Omega(V^{1+\epsilon})$  and  $V = \Omega(M)$ .

### 1.3 Our approach

The main component of our priority queues is the *x-treap*, a recursive structure inspired by similar cache-oblivious *x-box* [4] and cache-aware hashing data structures [12] that solve the dynamic dictionary problem in external memory (respectively, under predecessor and membership queries on a dynamic set of keys). To solve the priority queue problem, we adapt this recursive scheme to also handle priorities, inspired by the cache-oblivious priority queues of Brodal et al. [5] that support UPDATE, yet in suboptimal I/Os. Here we discuss informally these ideas, the rationale for combining them, and a back-of-the-envelope intuitive, but incomplete analysis. It is hoped that this will provide the intuition to more easily follow the full details in the sequel.

The idea behind the cache-oblivious priority queues of Brodal et al. [5] is simple. The structure has a logarithmic number of levels, where level  $i$  has two arrays, or buffers, of size roughly  $2^i$ . These buffers are called the *front* and *rear* buffers. They contain key-priority pairs or a key-delete message (described later). The idea is that the front buffers are sorted, with everything in the  $i$ -th front buffer having smaller priorities than everything in the  $(i + 1)$ -th front buffer. The items in the rear buffers do not have this rigorous ordering, but instead must be larger than the items in the rear buffer at the smaller levels. When an UPDATE operation occurs, the key-priority pair gets placed in the first rear buffer; when an EXTRACTMIN operation occurs, the key-priority pair with the smallest priority is removed from the first front buffer. Every time a level- $i$  buffer gets too full or empty relative to its target size of  $2^i$ , this is fixed by moving things up or down as needed, and moving things from the rear to front buffer if that respects the ordering of items in the front buffer. This resolution of problems is done efficiently using a scan of the affected and neighbouring levels. Thus, looking in a simplified manner at the lifetime of an UPDATED item, it will be inserted in the smallest rear buffer, be pushed down to larger rear buffers as they overflow, be moved from a rear buffer to a front buffer once it has gone down to a level where its priority is compatible with those in the corresponding front buffer, then moves up from the front buffer to smaller front buffers as they underflow, and is finally removed from the smallest front buffer during an EXTRACTMIN. Thus, during its lifetime, it could be moved from one level to another a total of  $O\left(\log_2 \frac{N}{B}\right)$  times at an I/O-cost of  $O\left(\frac{1}{B}\right)$  per level, for a total cost of  $O\left(\frac{1}{B} \log_2 \frac{N}{B}\right)$  I/Os. One detail is that when an item moves from a rear to a front buffer, we want to make sure that no items in larger levels with the same key and larger priority are ever removed. This is done through special delete messages, which stay in the rear buffers and percolate down, removing any key-priority pairs with the given key that they encounter in their buffer or the corresponding front buffer.

The problem with this approach is that the base-2 logarithm seems unavoidable, with the simple idea of a geometrically increasing buffer size. So here instead we use the more complicated recursion introduced with the cache-oblivious *x-box* [4] structure and also used in the cache-aware hashing data structures [12]. In its simplest form, used for a dictionary, an *x-box* has three buffers: top, middle and bottom (respectively of approximate size  $x$ ,  $x^{1.5}$  and  $x^2$ ), as well as  $\sqrt{x}$  recursive *upper-level*  $\sqrt{x}$ -boxes (ordered logically between the top and middle buffers) and  $x$  recursive *lower-level*  $\sqrt{x}$ -boxes (ordered logically between the middle and bottom buffers). Data in each buffer is sorted, and all keys in a given recursive buffer are

smaller than all keys in subsequent recursive buffers in the same level (upper or lower). There is no enforced order among keys in different buffers or in a recursive upper- or lower-level  $\sqrt{x}$ -box. The key feature of this construction is that the top/middle/bottom buffers have the same size as the neighbouring recursive buffers: the top buffer has size  $x$ , the top buffers of the upper-level recursive  $\sqrt{x}$ -boxes have total size  $x$ ; the middle buffer, sum of the bottom buffers of the upper-level, and sum of the top buffers of the lower-level recursive structures all have size  $x^{1.5}$ ; the sum of the bottom buffers of the lower-level recursive structures and the bottom buffer both have size  $x^2$ . Therefore, when for example a top buffer overflows, it can be fixed by moving excess items to the top buffers of the top recursive substructures. In a simplified view with only insertions, as buffers overflow, an item over its lifetime will percolate from the top buffer to the upper-level substructures, to the middle buffer, to the lower-level substructures, and to the bottom buffer, with each overflow handled only using scans. Assuming a base case of size  $M$ , there will be  $O(\log_M N)$  times that an item will move from one buffer to another and an equal number of times that an item will pass through a base case. One major advantage of this recursive approach, is that an item will pass through a small base case not just once at the structure's top, as before, but many times.

We combine these ideas to form the  $x$ -treap, described at a high level as follows: Everywhere an  $x$ -box has a buffer, we replace it with front and rear buffers storing key-priority pairs. The order used by  $x$ -box is imposed on the keys, not the priorities. The order imposed on priorities in the Brodal et al. structure are carried over and imposed on the priorities in different levels of the  $x$ -treap; this is aided by the fact that the buffers in the  $x$ -treap form a DAG, thus the buffers where items with a given key can appear, form a natural total order. Hence, this forms a treap-like arrangement where we use the keys for order in one dimension and priorities for order in the other. We use a separate trivial base case structure which is invoked at a size smaller than the memory size; it stores items in no particular order and thus supports fast insertion of items when a neighbouring buffer adds them ( $O(\frac{1}{B})$ ), but slow ( $O(M^\epsilon)$  amortized) removal of items with small priorities to fix the underflow of a front buffer above. Thus, again considering the typical hypothetical lifetime of an item, it will be inserted at the top in the rear buffer, percolate down  $O\left(\log_{\frac{M}{B}} \frac{N}{B}\right)$  levels and base cases at a cost of  $O(\frac{1}{B})$  amortized each, move over to a front buffer, then percolate up  $O\left(\log_{\frac{M}{B}} \frac{N}{B}\right)$  levels at a cost of  $O(\frac{M^\epsilon}{B})$  amortized each. Thus, the total amortized cost for an item that is eventually removed by an EXTRACTMIN is  $O\left(\frac{M^\epsilon}{B} \log_{\frac{M}{B}} \frac{N}{B}\right)$ .

However, we want the amortized cost for an item that is inserted via UPDATE to be much faster than this, i.e.  $O\left(\frac{1}{B} \log_{\frac{M}{B}} \frac{N}{B}\right)$ . This requires additional observations and tricks. The first is that, unlike Brodal et al., we do not use delete-type messages that percolate down to eliminate items with larger than minimum priority in order to prevent their removal from EXTRACTMIN. Instead, we adopt a much simpler approach, and use a hash table to keep track of all keys that have been removed by an EXTRACTMIN, and when an EXTRACTMIN returns a key that has been seen before, it is discarded and EXTRACTMIN is repeated. The second trick is to simply ensure that each buffer has at most one item with each key (and removes key-priority pairs other than the one with the minimum priority among those with the same key in the buffer). This has the effect that if there are a total of  $u$  updates performed on a key before it is removed by an EXTRACTMIN, the total cost will involve up to  $O\left(u \log_{\frac{M}{B}} \frac{N}{B}\right)$  percolations down at a cost of  $O(\frac{1}{B})$ , but only  $O\left(\log_{\frac{M}{B}}^2 \frac{N}{B}\right)$  percolations up at a cost of  $O(\frac{M^\epsilon}{B})$  amortized each. After the EXTRACTMIN, some items may still remain in the structure and will be discarded when removed by EXTRACTMIN however, due to the no-duplicates-per-level property there will only be  $O\left(\log_{\frac{M}{B}} N\right)$  such items (called *ghosts*) which

will incur a cost of at most  $O\left(\lceil \frac{M^\epsilon}{B} \log_{\frac{M}{B}} \frac{N}{B} \rceil\right)$  amortized each, where the ceiling accounts for accessing the hash table. Thus the total amortized cost for the lifetime of the  $u$  UPDATES and one EXTRACTMIN involving a single key is  $O\left(\frac{u}{B} \log_{\frac{M}{B}} \frac{N}{B} + \lceil \frac{M^\epsilon}{B} \log_{\frac{M}{B}} \frac{N}{B} \rceil \log_{\frac{M}{B}} \frac{N}{B}\right)$ . This cost can be apportioned in the amortized sense by having the EXTRACTMIN cost  $O\left(\lceil \frac{M^\epsilon}{B} \log_{\frac{M}{B}} \frac{N}{B} \rceil \log_{\frac{M}{B}} \frac{N}{B}\right)$  amortized and the updates cost  $O\left(\frac{1}{B} \log_{\frac{M}{B}} \frac{N}{B}\right)$  amortized each, assuming that the treap finishes in an empty state and no item can be UPDATED after it has been EXTRACTMIN'd.

The details that implement these rough ideas consume the rest of the paper. One complication that eludes the above discussion is that items don't just percolate down and then up; they could move up and down repeatedly and this can be handled through an appropriate potential function. The various layers of complexity needed for the  $x$ -treap recursion combined with the front/rear buffer idea, various types of over/underflows of buffers, a special base case, having the middle and bottom buffers be of size  $x^{1+\frac{\alpha}{2}}$  and  $x^{1+\alpha}$  for a suitable parameter  $\alpha$  rather than  $x^{1.5}$  and  $x^2$  as described above, and a duplicate-catching hash table, result in a complex structure with an involved potential analysis, but that follows naturally from the above high-level description.

## 2 $x$ -Treap

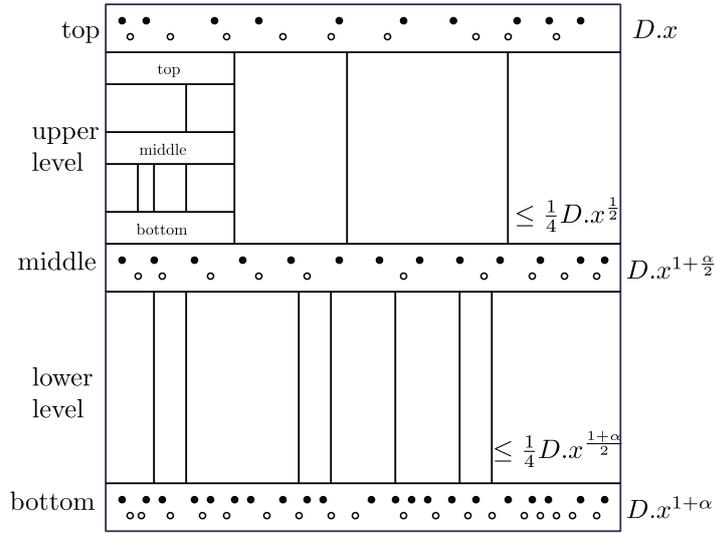
Given real parameter  $\alpha \in (0, 1]$  and *key range*  $[k_{\min}, k_{\max}) \subseteq \mathbb{R}$ , an  $x$ -treap  $D$  stores a set of at most  $2(D.x)^{1+\alpha}$  *elements*  $(*, k, p)$  associated with a *key*  $k \in [D.k_{\min}, D.k_{\max})$  and a *priority*  $p$  from a totally ordered set.  $D$  represents a set  $D.rep$  of pairs (key, priority), such that a particular key  $k$  contained in  $D$  is represented to have the smallest priority  $p$  of any element with key  $k$  stored in  $D$ , unless an element with key  $k$  and a smaller priority has been removed from the structure. In particular, we call the key and priority *represented*, when the pair (key, priority)  $\in D.rep$ . A *representative* element contains a represented key and its represented priority. More formally, we define:

$$D.rep := \bigcup_{\{k|(k,p) \in D\}} \left\{ \left( k, \min_p (k, p) \in D \right) \right\}$$

The proposed representation scheme works under the assumption that a key that is not represented by the structure anymore, cannot become represented again. In other words, a key returned by operation EXTRACTMIN cannot be INSERTED to the structure again.

The following *interface operations* are supported (See full version for their correctness):

- BATCHED-INSERT( $D, e_1, e_2, \dots, e_b$ ): For constant  $c \in (0, \frac{1}{3}]$ , insert  $b \leq c \cdot D.x$  elements  $e_1, e_2, \dots, e_b$  to  $D$ , given they are key-sorted with keys  $e_i.k \in [D.k_{\min}, D.k_{\max})$ ,  $i \in [1, b]$ . BATCHED-INSERT adds the pairs  $(e_i.k, e_i.p)$  to  $D.rep$  with key  $e_i.k$  that is not contained in  $D$  already. BATCHED-INSERT decreases the priority of a represented key  $e_i.k$  to  $e_i.p$ , if its represented priority is larger than  $e_i.p$  before the operation. More formally, let  $X_{new}$  contain the inserted pairs  $(e_i.k, e_i.p)$  with  $e_i.k \notin D.rep$ . Let  $X_{old}$  contain the pairs in  $D.rep$  with an inserted key, but with larger priority than the inserted one, and let  $X_{dec}$  contain these inserted pairs. After BATCHED-INSERT, a new  $x$ -treap  $D'$  is created where  $D'.rep = D.rep \cup X_{new} \cup X_{dec} \setminus X_{old}$ .
- BATCHED-EXTRACTMIN( $D$ ): For constant  $c \in (0, \frac{1}{4}]$ , remove and return the at most  $c \cdot D.x$  elements  $(k, p)$  with the smallest priorities in  $D$ . BATCHED-EXTRACTMIN removes the pairs  $X_{\min}$  from  $D.rep$  with the at most  $c \cdot D.x$  smallest priorities. Let  $X_{key}$  contain the pairs in  $D$  with keys in  $X_{\min}$ . After BATCHED-EXTRACTMIN, a new  $x$ -treap  $D'$  is created where  $D'.rep = D.rep \setminus X_{\min} \setminus X_{key}$ .



■ **Figure 1** Overview of an  $x$ -treap  $D$  on “key”  $\times$  “partial order” space. Black/white dots represent elements in the front/rear buffers, respectively. All buffers are resolved. Buffer sizes and maximum number of subtrees in a level are shown on the right-hand side.

► **Theorem 1.** *An  $x$ -treap  $D$  supports BATCHED-EXTRACTMIN in  $O\left(M^{\frac{\alpha}{1+\alpha}} \frac{1+\alpha}{B} \log_M D.x\right)$  amortized I/Os per element and BATCHED-INSERT in  $O\left(\frac{1+\alpha}{B} \log_M D.x\right)$  amortized I/Os per element, using  $O\left(\frac{(D.x)^{1+\alpha}}{B} \log_M D.x\right)$  blocks, for any real  $\alpha \in (0, 1]$ .*

The structure is recursive. The base case is described separately in Subsection 2.3. The base case structure is used when  $D.x \leq c' M^{\frac{1}{1+\alpha}}$  (for an appropriately chosen constant  $c' > 0$ ). Thus assuming  $D.x > c' M^{\frac{1}{1+\alpha}}$ , we define an  $x$ -treap to contain three *buffers* (which are arrays that store elements) and many  $\sqrt{x}$ -treaps (called *subtreaps*). Specifically, the *top*, *middle* and *bottom* buffers have *sizes*  $D.x$ ,  $(D.x)^{1+\frac{\alpha}{2}}$  and  $(D.x)^{1+\alpha}$ , respectively. Each buffer is divided in the middle into a *front* and a *rear buffer*. The subtrees are divided into the *upper* and the *lower level* that contain at most  $\frac{1}{4} (D.x)^{\frac{1}{2}}$  and  $\frac{1}{4} (D.x)^{\frac{1+\alpha}{2}}$  subtrees, respectively. Let  $|b|$  denote the *size* of a buffer  $b$ . We define the *capacity* of an  $x$ -treap  $D$  to be the maximum number of elements it can contain, which is  $D.x + \frac{5}{4} (D.x)^{1+\frac{\alpha}{2}} + \frac{5}{4} (D.x)^{1+\alpha} < 2 (D.x)^{1+\alpha}$ .

We define a partial order ( $\preceq$ ) using the terminology “above/below” among the buffers of an  $x$ -treap and all of the buffers in recursive subtrees or base case structures. In this order we have top buffer  $\preceq$  upper level recursive subtrees  $\preceq$  middle buffer  $\preceq$  lower level recursive subtrees  $\preceq$  bottom buffer.

Along with all buffers of  $D$ , we also store several additional pieces of bookkeeping information: a counter with the total number of elements stored in  $D$  and an index indicating which subtree is stored in which space in memory.

## 2.1 Invariants

An  $x$ -treap  $D$  maintains the following invariants with respect to every one of its top/middle/bottom buffers  $b$ . The invariants hold after the execution of each interface operation, but may be violated during the execution. They allow changes to  $D$  that do not change  $D.rep$ .

1. The front and rear buffers of  $b$  store elements sorted by key and left-justified.
2. The front buffer’s elements’ priorities are smaller than the rear buffer’s elements’ priorities.

3. The front buffer's elements' priorities are smaller than all elements' priorities in buffers below  $b$  in  $D$ .
4. For a top or middle buffer  $b$  with key range  $[b.k_{\min}, b.k_{\max})$ , the  $r$  upper or lower subtrees  $D_i, i \in \{1, r\}$ , respectively, have distinct key ranges  $[D_i.k_{\min}, D_i.k_{\max})$ , such that  $b.k_{\min} = D_1.k_{\min} < D_1.k_{\max} = D_2.k_{\min} < \dots < D_r.k_{\max} = b.k_{\max}$ .
5. If the middle or bottom buffer  $b$  is not empty, then at least one upper or lower subtree is not empty, respectively.

## 2.2 Auxiliary operations

The operations BATCHED-INSERT and BATCHED-EXTRACTMIN make use of the following *auxiliary operations* (See full version for their implementation and correctness):

- Operation RESOLVE( $D, b$ ). We say that a buffer  $b$  is *resolved*, when the front and rear buffers contain elements with pairs (key,priority)  $(k, p)$ , such that  $k$  is a represented key, and when the front buffer contains those elements with smallest priorities in the buffer. To resolve  $b$ , operation RESOLVE assigns to the elements with represented keys, the key's minimum priority stored in  $b$ . Also, it removes any elements with non-represented keys from  $b$ . RESOLVE restores Invariant 2 in  $b$ , when it is temporarily violated by other (interface or auxiliary) operations that call it.
- Operation INITIALIZE( $D, e_1, e_2, \dots, e_b$ ) distributes to a new  $x$ -treap  $D$ ,  $\frac{1}{4}(D.x) \leq b \leq \frac{1}{2}(D.x)^{1+\alpha}$  elements  $e_i, i \in [1, b]$  from a temporary array (divided in the middle into a front and a rear array, respecting Invariants 1 and 2).
- Operation FLUSH-UP( $D$ ) ensures that the front top buffer of  $D$  contains at least  $\frac{1}{4}D.x$  elements (unless all buffers of  $D$  contains less elements altogether, in which case FLUSH-UP moves them all to the top front buffer of  $D$ ). By Invariants 2 and 3, these are the elements in  $D$  with smallest priority.
- Operation FLUSH-DOWN( $D$ ) is called by BATCHED-INSERT on an  $x$ -treap  $D$  whose bottom buffer contains between  $\frac{1}{2}(D.x)^{1+\alpha}$  and  $(D.x)^{1+\alpha}$  elements. It moves to a new temporary array, at least  $\frac{1}{6}(D.x)^{1+\alpha}$  and at most  $\frac{2}{3}(D.x)^{1+\alpha}$  elements from the bottom buffer of  $D$ . It ensures that the largest priority elements are removed from  $D$ .
- Operation SPLIT( $D$ ) is called by BATCHED-INSERT on an  $x$ -treap  $D$  that contains between  $\frac{1}{2}(D.x)^{1+\alpha}$  and  $(D.x)^{1+\alpha}$  elements. It moves to a new temporary (front and rear) array, the at most  $\frac{1}{3}(D.x)^{1+\alpha}$  elements with largest keys in  $D$ .

## 2.3 Base case

The  $x$ -treap is a recursive structure. When the  $x$ -treap stores few enough elements so that it can be stored in internal memory, we use simple arrays to support the interface operations and operation FLUSH-UP.

► **Lemma 2.** *An  $O\left(M^{\frac{1}{1+\alpha}}\right)$ -treap fits in internal memory and supports operation BATCHED-INSERT in  $O(1/B)$  amortized I/Os per element and operations BATCHED-EXTRACTMIN and FLUSH-UP in  $\text{Scan}\left(M^{\frac{\alpha}{1+\alpha}}\right)$  amortized I/Os per element.*

**Proof.** For a universal positive constant  $c_0$  and a constant parameter  $c' < c_0^{\frac{1}{\alpha}+1}$ , we allocate an array of size  $\left(c' \left(M^{\frac{1}{1+\alpha}}\right)\right)^{\frac{\alpha}{1+\alpha}} \leq c_0 M$  and divide it in the middle into a front and a rear buffer that store elements and maintain only Invariants 1 and 2. To implement BATCHED-INSERT on at most  $\frac{c'}{2}M^{\frac{1}{1+\alpha}}$  elements, we simply add them to the rear buffer and

update the counter. This costs  $O\left(\frac{M^{\frac{1}{1+\alpha}}}{B}/\frac{1}{2}M^{\frac{1}{1+\alpha}}\right) = O\left(\frac{1}{B}\right)$  I/Os amortized per added element, since we only scan the part of the rear buffer where the elements are being added to. To implement BATCHED-EXTRACTMIN on at most  $\frac{c'}{2}M^{\frac{1}{1+\alpha}}$  extracted elements, we RESOLVE the array (as implemented for Theorem 1), remove and return all elements in the front buffer, and update the counter. By Lemma 5 (proven later in Subsection 2.5) this costs  $O\left(\frac{M}{B}/\frac{1}{2}M^{\frac{1}{1+\alpha}}\right) = O\left(\frac{M^{\frac{\alpha}{1+\alpha}}}{B}\right)$  I/Os amortized per extracted element. FLUSH-UP is implemented like BATCHED-EXTRACTMIN with the difference that the returned elements are not removed from the array. ◀

## 2.4 Interface operations

(See full version for the correctness of the interface operations.)

### 2.4.1 Inserting elements to an $x$ -treap

Interface operation BATCHED-INSERT on an  $x$ -treap  $D$  is implemented by means of the recursive subroutine BATCHED-INSERT( $D, e_1, \dots, e_{c \cdot |b|}, b$ ) that also takes as argument a top or middle buffer  $b$  of  $D$  and inserts  $c \cdot |b|$  elements  $e_1, \dots, e_{c \cdot |b|}$  (contained in a temporary array) inside and below  $b$  in  $D$ , for constant  $c \in (0, \frac{1}{3}]$ . For a bottom buffer  $b$ , a non-recursive subroutine BATCHED-INSERT( $D, b$ ) simply executes Step 1 below and discards the temporary array. BATCHED-INSERT( $D, e_1, \dots, e_{c \cdot |b|}, b$ ) is implemented as following:

1. If  $D.x > c' M^{\frac{1}{1+\alpha}} + c|b|$ :
  - 1.1. 2-way merge into the temporary array, the elements in the temporary array and in the rear buffer of  $b$  (by simultaneous scans in increasing key-order). RESOLVE  $b$  considering the temporary array as the rear buffer of  $b$ .
  - 1.2. Implicitly partition the front buffer of  $b$  and the temporary array by the key ranges of the subtrees immediately below  $b$ . Consider the subtrees in increasing key-order by reading the index of  $D$ . For every key range (associated with subtree  $D'$ ) that contains at least  $\frac{1}{3}(D.x)^{\frac{1}{2}}$  elements in either the front buffer of  $b$  or the temporary array: While the key range in the front buffer of  $b$  and in the temporary array contains at most  $\frac{2}{3}(D.x)^{\frac{1}{2}}$  elements, do:
    - 1.2.1. Find the  $\left(\frac{2}{3}(D.x)^{\frac{1}{2}}\right)$ -th smallest priority within the key range in the front buffer of  $b$  and in the temporary array (by an external memory order-statistics algorithm [3]) and move the elements in the key range with larger priority to a new auxiliary array (by simultaneous scans in increasing key-order).
    - 1.2.2. If the counter of  $D'$  plus the auxiliary array's size does not exceed the capacity of  $D'$ : BATCHED-INSERT the elements in the auxiliary array to the top buffer of  $D'$ . Discard the auxiliary array.
    - 1.2.3. Else, if there are fewer than the maximum allowed number of subtrees in the level immediately below  $b$ : SPLIT  $D'$ . Let  $k$  be the smallest key in the array returned by SPLIT (determined by a constant number of random accesses to the leftmost elements in the returned front/rear array). Move the elements in the auxiliary array with key smaller than  $k$  to a new temporary array (by a scan), BATCHED-INSERT these elements to  $D'$  and discard this temporary array. 2-way merge the remaining elements in the auxiliary array into the returned rear array and discard the auxiliary array. INITIALIZE a new subtree with the elements in the returned array. Discard the returned array.

1.2.4. Else, FLUSH-DOWN all subtrees immediately below  $b$ , which writes them to many returned arrays. 2-way merge into a new temporary array, all elements in  $b$  and in all returned arrays (by simultaneous scans in increasing key-order). (When the scan on a subtree's temporary array is over, determine the subtree with the key-next elements in the level by reading the index of  $D$ .) BATCHED-INSERT the elements in the new temporary array to the buffer  $b'$  immediately below  $b$ . Discard the new temporary array and all returned arrays.

1.3. Discard the temporary array and update the counter of  $D$ .

1.4. Else if  $D.x \leq c'M^{\frac{1}{1+\alpha}} + c|b|$ : BATCHED-INSERT the elements to the base case structure.

## 2.4.2 Extracting minimum-priority elements from an $x$ -treap

Interface operation BATCHED-EXTRACTMIN on an  $x$ -treap  $D$  is implemented as following:

1. If  $D.x > c'M^{\frac{1}{1+\alpha}}$ :
  - 1.1 If the front top buffer contains less than  $\frac{1}{4}D.x$  elements: FLUSH-UP the top buffer.
  - 1.2 Remove and return all the elements  $(e_i.k, e_i.p)$  from the front top buffer.
  - 1.3 Update the counter of  $D$ .
2. Else if  $D.x \leq c'M^{\frac{1}{1+\alpha}}$ : BATCHED-EXTRACTMIN the base case structure.

## 2.5 Analysis

(See full version for the proofs of Lemmata 3, 4, 5, 6 and 7, respectively.)

► **Lemma 3.** *An  $x$ -treap  $D$  has  $O\left(\log_{\frac{M}{B}} D.x\right)$  levels and occupies  $O\left((D.x)^{1+\alpha} \log_{\frac{M}{B}} D.x\right)$  blocks of space.*

► **Lemma 4.** *By the tall-cache assumption, scanning the buffers of an  $x$ -treap  $D$  and randomly accessing  $O\left((D.x)^{\frac{1+\alpha}{2}}\right)$  subtrees takes  $\text{Scan}\left((D.x)^{1+\alpha}\right)$  I/Os, for any real  $\alpha \in (0, 1]$ .*

A buffer  $b_i$  at level  $i \leq h = O\left(\log_{\frac{M}{B}} D.x\right)$  with current number of elements in the front and rear buffers  $b_f, b_r$ , respectively, has potential  $\Phi(b_i) = \Phi_f(b_i) + \Phi_r(b_i)$ , such that (for constants  $\varepsilon := \frac{\alpha}{1+\alpha}$  and  $c_0 \geq 1$ ):

$$\begin{aligned} \blacksquare \quad \Phi_f(b_i) &= \begin{cases} 0, & \text{if } \frac{1}{4}|b_i| \leq b_f \leq \frac{1}{3}|b_i|, \\ \frac{c_0}{B} M^\varepsilon \cdot \left(\frac{|b_i|}{4} - b_f\right) \cdot (h - i), & \text{if } b_f < \frac{1}{4}|b_i|, \\ \frac{c_0}{B} \cdot \left(b_f - \frac{|b_i|}{3}\right) \cdot (h - i), & \text{if } b_f > \frac{1}{3}|b_i|, \end{cases} \\ \blacksquare \quad \Phi_r(b_i) &= \begin{cases} 2\frac{c_0}{B} \cdot \left(b_r - \frac{|b_i|}{2}\right) \cdot (h - i), & \text{if } b_r > 0. \end{cases} \end{aligned}$$

In general, a particular element will be added to a rear buffer and will be moved down the levels of the structure over rear buffers by operation FLUSH-DOWN. A RESOLVE operation will move the element from the rear to the front buffer, if it is a representative element. From this point, it will be moved up the levels over front buffers by operation FLUSH-UP. If it is not representative, it will either get discarded by RESOLVE (when there is an element with the same key and with smaller priority in the same buffer) or it will keep going down the structure. Since RESOLVE leaves only one element per key at the level it operates,  $O\left(\log_{\frac{M}{B}} D.x\right)$  elements with the same key (i.e. at most one per level) will remain in the structure after the extraction of the representative element for this key.

The  $M^\varepsilon$ -factor accounts for the extra cost of FLUSH-UP and BATCHED-EXTRACTMIN, the  $(h - i)$ -factor allows for moving elements up or down a level by FLUSH-UP and FLUSH-DOWN and the 2-factor accounts for moving elements from the rear to the front buffer.

► **Lemma 5.** RESOLVE on a buffer  $b_i$  takes  $\text{Scan}(|b_i|) + O(1)$  amortized I/Os.

► **Lemma 6.** BATCHED-INSERT on an  $x$ -treap  $D$  takes  $O\left(\frac{1+\alpha}{B} \log_{\frac{M}{B}} D.x\right)$  amortized I/Os per element, for any real  $\alpha \in (0, 1]$ .

► **Lemma 7.** BATCHED-EXTRACTMIN on an  $x$ -treap  $D$  takes  $O\left(M^{\frac{\alpha}{1+\alpha}} \frac{1+\alpha}{B} \log_{\frac{M}{B}} D.x\right)$  amortized I/Os per element, for any real  $\alpha \in (0, 1]$ .

### 3 Priority queues

Priority queues support operations UPDATE and EXTRACTMIN that are defined similarly to BATCHED-INSERT and BATCHED-EXTRACTMIN, respectively, but on a *single* element.

To support these operations, we compose a priority queue out of its batched counterpart in Theorem 1. The data structure on  $N$  elements consists of  $1 + \log_{1+\alpha} \log_2 N$   $x$ -treaps of doubly increasing size with parameter  $\alpha$  being set the same in all of them. Specifically, for  $i \in \{0, \log_{1+\alpha} \log_2 N\}$ , the  $i$ -th  $x$ -treap  $D_i$  has  $D_i.x = 2^{(1+\alpha)^i}$ . We store all keys returned by EXTRACTMIN in a hash table  $X$  [12, 8].

For  $i \in \{0, \log_{1+\alpha} \log_2 N - 1\}$ , we define the top buffer of  $D_i$  to be “below” the bottom buffer of  $D_{i-1}$  and the bottom buffer of  $D_i$  to be “above” the top buffer of  $D_{i+1}$ . We define the set of represented pairs (key, priority)  $rep = \bigcup_{i=0}^{\log_{1+\alpha} \log_2 N} D_i.rep \setminus \{(k, p) \mid k \in X\}$  and call *represented* the keys and priorities in  $rep$ . We maintain the invariant that the maximum represented priority in  $D_i.rep$  is smaller than the smallest represented priority below it.

To implement UPDATE on a pair (key, priority)  $\in rep$ , we BATCHED-INSERT the corresponding element to  $D_0$ .  $D_0$  handles single-element batches, since for  $i = 0 \Rightarrow x = \Theta(1)$ . When  $D_i$  reaches capacity (i.e. contains  $(D_i.x)^{1+\alpha}$  elements), we call FLUSH-DOWN on it, BATCHED-INSERT the elements in the returned temporary array to  $D_{i+1}$  and discard the array. This process terminates at the first  $x$ -treap that can accommodate these elements without reaching capacity.

To implement EXTRACTMIN, we call BATCHED-EXTRACTMIN to the first  $x$ -treap  $D_i$  with a positive counter, add the extracted elements to the (empty) bottom front buffer of  $D_{i-1}$  and repeat this process on  $D_{i-1}$ , until  $D_0$  returns at least one element. If the returned key does not belong to  $X$ , we insert it. Else, we discard the element and repeat EXTRACTMIN.

To implement DELETE of a key, we add the key to  $X$ .

(See full version for the proof of Theorem 8.)

► **Theorem 8.** There exist priority queues on  $N$  elements that support operation UPDATE in  $O\left(\frac{1}{B} \log_{\frac{M}{B}} \frac{N}{B}\right)$  amortized I/Os per element, operations EXTRACTMIN and DELETE in amortized  $O\left(\lceil \frac{M^{1+\alpha}}{B} \log_{\frac{M}{B}} \frac{N}{B} \rceil \log_{\frac{M}{B}} \frac{N}{B}\right)$  I/Os per element, using  $O\left(\frac{N}{B} \log_{\frac{M}{B}} \frac{N}{B}\right)$  blocks, for any real  $\alpha \in (0, 1]$ .

### 4 Buffered repository trees

A *buffered repository tree (BRT)* [6, 2, 7] stores a multi-set of at most  $N$  elements, each associated with a *key* in the range  $[1 \dots k_{\max}]$ . It supports the operations INSERT and EXTRACT that, respectively, insert a new element to the structure and remove and report all

elements in the structure with a given key. To implement a BRT, we make use of the  $x$ -box [4]. Given positive real  $\alpha \leq 1$  and key range  $[k_{\min}, k_{\max}] \subseteq \mathfrak{R}$ , an  $x$ -box  $D$  stores a set of at most  $\frac{1}{2}(D.x)^{1+\alpha}$  elements associated with a key  $k \in [D.k_{\min}, D.k_{\max}]$ . An  $x$ -box supports the following operations:

- **BATCHED-INSERT**( $D, e_1, e_2, \dots, e_b$ ): For constant  $c \in (0, \frac{1}{2}]$ , insert  $b \leq c \cdot D.x$  elements  $e_1, e_2, \dots, e_b$  to  $D$ , given they are key-sorted with keys  $e_i.k \in [D.k_{\min}, D.k_{\max}]$ ,  $i \in [1, b]$ .
- **SEARCH**( $D, \kappa$ ): Return pointers to all elements in  $D$  with key  $\kappa$ , given they exist in  $D$  and  $\kappa \in [D.k_{\min}, D.k_{\max}]$ .

To implement operation **EXTRACT**( $D, \kappa$ ) that extracts all elements with key  $\kappa$  from an  $x$ -box  $D$ , we **SEARCH**( $D, \kappa$ ) and remove from  $D$  all returned pointed elements.

The BRT on  $N$  elements consists of  $1 + \log_{1+\alpha} \log_2 N$   $x$ -boxes of doubly increasing size with parameter  $\alpha$  being set the same in all of them. We obtain the stated bounds by modifying the proof of the  $x$ -box [4, Theorem 5.1] to account for Lemmata 9 and 10.

► **Lemma 9.** For  $D.x = \Omega\left(M^{\frac{1}{1+\alpha}}\right)$ , an  $x$ -box supports operation **BATCHED-INSERT** in amortized  $O\left(\frac{1+\alpha}{B} \log_{\frac{M}{B}} \frac{D.x}{B}\right)$  I/Os and operation **EXTRACT** on  $K$  extracted elements in amortized  $O\left((1+\alpha) \log_{\frac{M}{B}} \frac{D.x}{B} + \frac{K}{B}\right)$  I/Os, using  $O\left(\frac{(D.x)^{1+\alpha}}{B}\right)$  blocks of space.

**Proof.** Regarding **BATCHED-INSERT** on a cache-aware  $x$ -box, we obtain  $O\left(\frac{1+\alpha}{B} \log_{\frac{M}{B}} \frac{D.x}{B}\right)$  amortized I/Os by modifying the proof of **BATCHED-INSERT** [4, Theorem 4.1] according to the proof of Lemma 6. Specifically, every element is charged  $O(1/B)$  amortized I/Os, instead of  $O\left(1/B^{\frac{1}{1+\alpha}}\right)$ , and the recursion stops when  $D.x = O\left(M^{\frac{1}{1+\alpha}}\right)$ , instead of  $D.x = O\left(B^{\frac{1}{1+\alpha}}\right)$ .

Regarding **SEARCHING** for the first occurrence of a key in a cache-aware  $x$ -box, we obtain  $O\left(\log_{\frac{M}{B}} \frac{D.x}{B}\right)$  amortized I/Os by modifying the proof of **SEARCH** [4, Lemma 4.1], such that the recursion stops when  $D.x = O\left(M^{\frac{1}{1+\alpha}}\right)$ , instead of  $D.x = O\left(B^{\frac{1}{1+\alpha}}\right)$ . To **EXTRACT** all  $K$  occurrences of the searched key, we access them by scanning the  $x$ -box and by following fractional cascading pointers, which incurs an extra  $O(K/B)$  I/Os. ◀

► **Lemma 10.** An  $O\left(M^{\frac{1}{1+\alpha}}\right)$ -box fits in internal memory and supports operations **BATCHED-INSERT** in  $O(1/B)$  amortized I/Os per element and operation **EXTRACT** on  $K$  extracted elements in  $O\left(M^{\frac{\alpha}{1+\alpha}}\right)$  amortized I/Os per element.

**Proof.** We allocate an array of size  $O(M)$  and implement **BATCHED-INSERT** by simply appending the inserted element to the array and **EXTRACT** by scanning the array and removing and returning all occurrences of the searched key. ◀

► **Theorem 11.** There exist buffered priority trees on a multi-set of  $N$  elements and of  $K$  extracted elements that support operations **INSERT** and **EXTRACT** in amortized  $O\left(\frac{1}{B} \log_{\frac{M}{B}} \frac{N}{B}\right)$  and  $O\left(\frac{M^{\frac{\alpha}{1+\alpha}}}{B} \log_{\frac{M}{B}} \frac{N}{B} + \frac{K}{B}\right)$  I/Os per element, using  $O\left(\frac{N}{B}\right)$  blocks, for any real  $\alpha \in (0, 1]$ .

## 5 Graph algorithms

► **Theorem 12.** Single source shortest paths on a directed graph with  $V$  nodes and  $E$  edges can be computed in  $O\left(V \frac{M^{\frac{1+\alpha}{1+\alpha}}}{B} \log_{\frac{M}{B}}^2 \frac{E}{B} + V \log_{\frac{M}{B}} \frac{E}{B} + \frac{E}{B} \log_{\frac{M}{B}} \frac{E}{B}\right)$  I/Os, for any real  $\alpha \in (0, 1]$ .

**Proof.** The algorithm of Vitter [15] (described in detail in [7, Lemma 4.1] for the cache-oblivious model) makes use of a priority queue that supports the UPDATE operation and of a BRT on  $O(E)$  elements. Specifically, it makes  $V$  calls to EXTRACTMIN and  $E$  calls to UPDATE on the priority queue and  $V$  calls to EXTRACT and  $E$  calls to INSERT on the BRT. Hence, we obtain the stated bounds, by using Theorems 8 and 11. ◀

► **Theorem 13.** *Depth-first search and breadth-first search numbers can be assigned to a directed graph with  $V$  nodes and  $E$  edges in  $O\left(V \frac{M^{1+\alpha}}{B} \log^2 \frac{E}{M} \frac{E}{B} + V \log \frac{E}{M} \frac{E}{B} + \frac{E}{B} \log \frac{E}{M} \frac{E}{B}\right)$  I/Os, for any real  $\alpha \in (0, 1]$ .*

**Proof.** The algorithm of Buchsbaum et al. [6] makes use of a priority queue and of a BRT on  $O(E)$  elements. Specifically, it makes  $2V$  calls to EXTRACTMIN and  $E$  calls to INSERT on the priority queue and  $2V$  calls to EXTRACT and  $E$  calls to INSERT on the BRT [6, Theorem 3.1]. Hence, we obtain the stated bounds, by using Theorems 8 and 11. ◀

---

## References

- 1 Alok Aggarwal and S. Vitter, Jeffrey. The Input/Output Complexity of Sorting and Related Problems. *Commun. ACM*, 31(9):1116–1127, September 1988. doi:10.1145/48529.48535.
- 2 Lars Arge, Michael A. Bender, Erik D. Demaine, Bryan Holland-Minkley, and J. Ian Munro. An Optimal Cache-Oblivious Priority Queue and Its Application to Graph Algorithms. *SIAM Journal on Computing*, 36(6):1672–1695, 2007. doi:10.1137/S0097539703428324.
- 3 Manuel Blum, Robert W. Floyd, Vaughan Pratt, Ronald L. Rivest, and Robert E. Tarjan. Time Bounds for Selection. *J. Comput. Syst. Sci.*, 7(4):448–461, August 1973. doi:10.1016/S0022-0000(73)80033-9.
- 4 Gerth Stølting Brodal, Erik D. Demaine, Jeremy T. Fineman, John Iacono, Stefan Langerman, and J. Ian Munro. Cache-oblivious Dynamic Dictionaries with Update/Query Tradeoffs. In *Proceedings of the Twenty-first Annual ACM-SIAM Symposium on Discrete Algorithms, SODA '10*, pages 1448–1456, Philadelphia, PA, USA, 2010. Society for Industrial and Applied Mathematics. URL: <http://dl.acm.org/citation.cfm?id=1873601.1873718>.
- 5 Gerth Stølting Brodal, Rolf Fagerberg, Ulrich Meyer, and Norbert Zeh. Cache-Oblivious Data Structures and Algorithms for Undirected Breadth-First Search and Shortest Paths. In Torben Hagerup and Jyrki Katajainen, editors, *Algorithm Theory - SWAT 2004*, pages 480–492, Berlin, Heidelberg, 2004. Springer Berlin Heidelberg.
- 6 Adam L. Buchsbaum, Michael Goldwasser, Suresh Venkatasubramanian, and Jeffery R. Westbrook. On External Memory Graph Traversal. In *Proceedings of the Eleventh Annual ACM-SIAM Symposium on Discrete Algorithms, SODA '00*, pages 859–860, Philadelphia, PA, USA, 2000. Society for Industrial and Applied Mathematics. URL: <http://dl.acm.org/citation.cfm?id=338219.338650>.
- 7 Rezaul A. Chowdhury and Vijaya Ramachandran. Cache-Oblivious Buffer Heap and Cache-Efficient Computation of Shortest Paths in Graphs. *ACM Trans. Algorithms*, 14(1):1:1–1:33, January 2018. doi:10.1145/3147172.
- 8 Alex Conway, Martín Farach-Colton, and Philip Shilane. Optimal Hashing in External Memory. In Ioannis Chatzigiannakis, Christos Kaklamanis, Dániel Marx, and Donald Sannella, editors, *45th International Colloquium on Automata, Languages, and Programming (ICALP 2018)*, volume 107 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 39:1–39:14, Dagstuhl, Germany, 2018. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik. doi:10.4230/LIPIcs.ICALP.2018.39.
- 9 Kasper Eenberg, Kasper Green Larsen, and Huacheng Yu. DecreaseKeys Are Expensive for External Memory Priority Queues. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2017*, pages 1081–1093, New York, NY, USA, 2017. ACM. doi:10.1145/3055399.3055437.

- 10 R. Fadel, K. V. Jakobsen, J. Katajainen, and J. Teuhola. Heaps and Heapsort on Secondary Storage. *Theor. Comput. Sci.*, 220(2):345–362, June 1999. doi:10.1016/S0304-3975(99)00006-7.
- 11 Matteo Frigo, Charles E. Leiserson, Harald Prokop, and Sridhar Ramachandran. Cache-Oblivious Algorithms. In *Proceedings of the 40th Annual Symposium on Foundations of Computer Science, FOCS '99*, pages 285–, Washington, DC, USA, 1999. IEEE Computer Society. URL: <http://dl.acm.org/citation.cfm?id=795665.796479>.
- 12 John Iacono and Mihai Pătraşcu. Using Hashing to Solve the Dictionary Problem. In *Proceedings of the Twenty-third Annual ACM-SIAM Symposium on Discrete Algorithms, SODA '12*, pages 570–582, Philadelphia, PA, USA, 2012. Society for Industrial and Applied Mathematics. URL: <http://dl.acm.org/citation.cfm?id=2095116.2095164>.
- 13 Shunhua Jiang and Kasper Green Larsen. A Faster External Memory Priority Queue with DecreaseKeys. In Timothy M. Chan, editor, *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2019, San Diego, California, USA, January 6-9, 2019*, pages 1331–1343. SIAM, 2019. doi:10.1137/1.9781611975482.81.
- 14 Vijay Kumar and Eric J. Schwabe. Improved Algorithms and Data Structures for Solving Graph Problems in External Memory. In *Proceedings of the 8th IEEE Symposium on Parallel and Distributed Processing (SPDP '96)*, SPDP '96, pages 169–, Washington, DC, USA, 1996. IEEE Computer Society. URL: <http://dl.acm.org/citation.cfm?id=829517.830723>.
- 15 Jeffrey Scott Vitter. External Memory Algorithms and Data Structures: Dealing with Massive Data. *ACM Comput. Surv.*, 33(2):209–271, June 2001. doi:10.1145/384192.384193.
- 16 Zhewei Wei and Ke Yi. Equivalence between Priority Queues and Sorting in External Memory. In Andreas S. Schulz and Dorothea Wagner, editors, *Algorithms - ESA 2014*, pages 830–841, Berlin, Heidelberg, 2014. Springer Berlin Heidelberg.