

Basic Research in Computer Science

BRICS RS-04-2 Brodal et al.: Cache-Oblivious Undirected Breadth-First Search and Shortest Paths

**Cache-Oblivious  
Data Structures and Algorithms for  
Undirected Breadth-First Search  
and Shortest Paths**

**Gerth Stølting Brodal  
Rolf Fagerberg  
Ulrich Meyer  
Norbert Zeh**

**BRICS Report Series**

**ISSN 0909-0878**

**RS-04-2**

**February 2004**

**Copyright © 2004, Gerth Stølting Brodal & Rolf Fagerberg & Ulrich Meyer & Norbert Zeh.  
BRICS, Department of Computer Science  
University of Aarhus. All rights reserved.**

**Reproduction of all or part of this work  
is permitted for educational or research use  
on condition that this copyright notice is  
included in any copy.**

**See back inner page for a list of recent BRICS Report Series publications.  
Copies may be obtained by contacting:**

**BRICS  
Department of Computer Science  
University of Aarhus  
Ny Munkegade, building 540  
DK-8000 Aarhus C  
Denmark  
Telephone: +45 8942 3360  
Telefax: +45 8942 3255  
Internet: BRICS@brics.dk**

**BRICS publications are in general accessible through the World Wide  
Web and anonymous FTP through these URLs:**

`http://www.brics.dk`  
`ftp://ftp.brics.dk`  
**This document in subdirectory RS/04/2/**

# Cache-Oblivious Data Structures and Algorithms for Undirected Breadth-First Search and Shortest Paths

Gerth Stølting Brodal<sup>\*,†,‡</sup>    Rolf Fagerberg<sup>\*,‡</sup>    Ulrich Meyer<sup>§,‡</sup>  
Norbert Zeh<sup>¶</sup>

February 11, 2004

## Abstract

We present improved cache-oblivious data structures and algorithms for breadth-first search (BFS) on undirected graphs and the single-source shortest path (SSSP) problem on undirected graphs with non-negative edge weights. For the SSSP problem, our result closes the performance gap between the currently best *cache-aware* algorithm and the *cache-oblivious* counterpart. Our cache-oblivious SSSP-algorithm takes nearly full advantage of block transfers for *dense* graphs. The algorithm relies on a new data structure, called *bucket heap*, which is the first cache-oblivious priority queue to efficiently support a weak DECREASEKEY operation. For the BFS problem, we reduce the number of I/Os for *sparse* graphs by a factor of nearly  $\sqrt{B}$ , where  $B$  is the cache-block size, nearly closing the performance gap between the currently best *cache-aware* and *cache-oblivious* algorithms.

---

\*BRICS (Basic Research in Computer Science, [www.brics.dk](http://www.brics.dk), funded by the Danish National Research Foundation), Department of Computer Science, University of Aarhus, DK-8000 Århus C, Denmark. E-mail: {gerth,rolf}@brics.dk.

†Supported by the Carlsberg Foundation (contract number ANS-0257/20).

‡Partially supported by the Future and Emerging Technologies programme of the EU under contract number IST-1999-14186 (ALCOM-FT).

§Max-Planck-Institut für Informatik, Stuhlsatzhausenweg 85, 66123 Saarbrücken, Germany. E-mail: [uli@uli-meyer.de](mailto:uli@uli-meyer.de). Partially supported by DFG grant SA 933/1-1.

¶Faculty of Computer Science, Dalhousie University, 6050 University Ave, Halifax, NS B3H 1W5, Canada. E-mail: [nzeh@cs.dal.ca](mailto:nzeh@cs.dal.ca).

# 1 Introduction

*Breadth-first search* (BFS) and the *single-source shortest path* (SSSP) problem are fundamental combinatorial optimization problems with numerous applications. SSSP is defined as follows: Let  $G = (V, E)$  be a graph with  $V$  vertices and  $E$  edges,<sup>1</sup> let  $s$  be a distinguished vertex of  $G$ , and let  $\omega$  be an assignment of non-negative real *weights* to the edges of  $G$ . The weight of a path is the sum of the weights of its edges. We want to find for every vertex  $v$  that is reachable from  $s$ , the weight  $\text{dist}(s, v)$  of a minimum-weight (“shortest”) path from  $s$  to  $v$ . BFS can be seen as the unweighted version of SSSP.

Both problems are well understood in the RAM model, where the cost of a memory access is assumed to be independent of the accessed memory location. However, modern computers contain a hierarchy of memory levels with each level acting as a cache for the next. Typical components of the memory hierarchy are registers, level-1 cache, level-2 cache, level-3 cache, main memory, and disk. The time for accessing a level increases for each new level (most dramatically when going from main memory to disk), making the cost of a memory access highly dependent on what is the currently lowest memory level that contains the accessed element. This is not accounted for in the RAM model, and current BFS and SSSP-algorithms, when run in memory hierarchies, have turned out to be notoriously inefficient for sparse input graphs because their memory access patterns are highly unstructured. The purpose of this paper is to provide improved data structures and algorithms for BFS and SSSP under the currently most powerful model for multi-level memory hierarchies.

**I/O model.** The most widely used model for the design of cache-aware algorithms is the *I/O-model* of Aggarwal and Vitter [2]. This model assumes a memory hierarchy consisting of two levels, the lower level having size  $M$  and the transfer between the two levels taking place in blocks of  $B$  consecutive data items. The cost of the computation is the number of blocks transferred (I/Os). The parameters  $M$  and  $B$  are assumed to be known to the algorithm. The strength of the I/O-model is its simplicity, while it still adequately models the situation where the I/Os between two levels of the memory hierarchy dominate the running time of the algorithm; this is often the case when the size of the data significantly exceeds the size of main memory. A comprehensive list of results for the I/O-model have been obtained—see the surveys [3, 17, 20] and the references therein. One of the fundamental facts is that, in the I/O-model, comparison-based sorting of  $N$  elements takes  $\Theta(\text{Sort}(N))$  I/Os in the worst case, where  $\text{Sort}(N) = \frac{N}{B} \log_{M/B} \frac{N}{B}$ .

---

<sup>1</sup>For convenience, we use the name of a set to denote both the set and its cardinality. Furthermore, we will assume that  $E = \Omega(V)$ , in order to simplify notation.

**Cache-oblivious model.** Recently, the concept of *cache-oblivious* algorithms was introduced by Frigo et al. [14]. In essence, algorithms are formulated in the RAM model (i.e.,  $M$  and  $B$  are not used in the description of the algorithms); but they are analyzed in the I/O-model for arbitrary block size  $B$  and memory size  $M$ . I/Os are assumed to be performed automatically by an optimal offline cache-replacement strategy. Since the analysis holds for any block and memory sizes, it holds for *all* levels of a multi-level memory hierarchy (see [14] for details). Thus, the cache-oblivious model elegantly combines the simplicity of the I/O-model with a coverage of the entire memory hierarchy. An additional benefit is that the designed algorithms are independent of the characteristics of the memory hierarchy and are therefore portable. Together with the model, Frigo et al. also introduced optimal cache-oblivious algorithms for matrix multiplication, matrix transposition, FFT, and sorting [14]. The cache-oblivious sorting bound matches that for the I/O-model:  $\mathcal{O}(\text{Sort}(N))$  I/Os. After the publication of [14], a number of results for the cache-oblivious model have appeared; see [11, 17] for recent surveys.

Some results in the cache-oblivious model, in particular those concerning sorting and algorithms and data structures that can be used to sort, such as priority queues, are proved under the assumption that  $M \geq B^2$ . This is also known as the *tall cache assumption*. In particular, this assumption is made in the *Funnelsort* algorithm of Frigo et al. [14]. A variant termed *Lazy-Funnelsort* [6] works under the weaker tall cache assumption that  $M \geq B^{1+\varepsilon}$ , for any fixed  $\varepsilon > 0$ . The cost is a slow-down by a factor of  $1/\varepsilon$  compared to the optimal sorting bound  $\Theta(\text{Sort}(N))$  for the case when  $M \gg B^{1+\varepsilon}$ . Recently, it has been shown [8] that a tall cache assumption is necessary for cache-oblivious comparison-based sorting algorithms, in the sense that the trade-off attained by Lazy-Funnelsort between the strength of the tall cache assumption and the cost for the case when  $M \gg B^{1+\varepsilon}$  is best possible.

## Previous and Related Work for BFS/SSSP in Memory Hierarchies

Graph algorithms for the I/O-model have received considerable attention in recent years; most efficient cache-aware algorithms do have a cache-oblivious counterpart that achieves the same performance; see Table 1. Despite these efforts, only little progress has been made on the solution of the SSSP-problem with general non-negative edge weights using either cache-aware or cache-oblivious algorithms: The best known lower bound is  $\Omega(\text{Sort}(E))$  I/Os, which can be obtained through a reduction from list-ranking; but the currently best algorithm, due to Kumar and Schwabe [15], performs  $\mathcal{O}(V + (E/B) \cdot \log(E/B))$  I/Os on undirected graphs.<sup>2</sup> For  $E = \mathcal{O}(V)$ , this is hardly better than naïvely

---

<sup>2</sup> $\log(x)$  denotes the binary logarithm of  $x$ .

Problem	Best cache-oblivious result
List ranking	$\mathcal{O}(\text{Sort}(V))$ [4]
Euler Tour	$\mathcal{O}(\text{Sort}(V))$ [4]
ST, Minimum ST	$\mathcal{O}(\text{Sort}(E) \cdot \log \log V)$ [4] $\mathcal{O}(\text{Sort}(E))$ (randomized) [1]
Undir. BFS	$\mathcal{O}(V + \text{Sort}(E))$ [19] $\mathcal{O}(\text{ST}(E) + \text{Sort}(E) + \frac{E}{B} \cdot \log V + \sqrt{V \cdot E/B})$ <b>New</b> $\mathcal{O}(\text{ST}(E) + \text{Sort}(E) + \frac{E}{B} \cdot \frac{1}{\epsilon} \cdot \log \log V + \sqrt{V \cdot E/B} \cdot \sqrt{V \cdot B/E^c})$ <b>New</b>
Dir. BFS & DFS	$\mathcal{O}((V + E/B) \cdot \log V + \text{Sort}(E))$ [4]
Undir. SSSP	$\mathcal{O}(V + (E/B) \cdot \log(E/B))$ <b>New</b>

Problem	Best cache-aware result
List ranking	$\mathcal{O}(\text{Sort}(V))$ [10]
Euler Tour	$\mathcal{O}(\text{Sort}(V))$ [10]
ST, Minimum ST	$\mathcal{O}(\text{Sort}(E) \cdot \log \log(V \cdot B/E))$ [5, 19] $\mathcal{O}(\text{Sort}(E))$ (randomized) [1]
Undir. BFS	$\mathcal{O}(\sqrt{V \cdot E/B} + \text{Sort}(E) + \text{ST}(E))$ [16]
Dir. BFS & DFS	$\mathcal{O}((V + E/B) \cdot \log V + \text{Sort}(E))$ [9]
Undir. SSSP	$\mathcal{O}(V + (E/B) \cdot \log(E/B))$ [15]

Table 1: I/O-bounds for some fundamental graph problems.

running Dijkstra’s internal-memory algorithm [12, 13] in external memory, which would take  $\mathcal{O}(V \log V + E)$  I/Os. On dense graphs, however, the algorithm is efficient. The algorithm of [15] is not cache-oblivious, because the applied external-memory priority queue based on the tournament tree is not cache-oblivious. Cache-oblivious priority queues do exist [4, 7]; but none of them efficiently supports a DECREASEKEY operation. Indeed, the tournament tree is also the only cache-aware priority queue that supports at least a weak form of this operation.

For bounded edge weights, an improved external-memory SSSP-algorithm has been developed recently [18]. This algorithm is an extension of the currently best external-memory BFS-algorithm [16], *Fast-BFS*, which performs  $\mathcal{O}(\sqrt{V \cdot E/B} + \text{Sort}(E) + \text{ST}(E))$  I/Os, where  $\text{ST}(E)$  is the number of I/Os required to compute a spanning tree (see Table 1). Again, the key data structure used in *Fast-BFS* is not cache-oblivious, which is why the currently best cache-oblivious BFS-algorithm is that of [19].

## Our Results

In this paper, we obtain the following results:

- In Section 2, we develop the first cache-oblivious priority queue, called *bucket heap*, that supports an UPDATE operation, which is a combination of the INSERT and DECREASEKEY operations supported by standard internal-memory priority queues. The amortized cost of operations UPDATE, DELETE, and DELETEMIN is  $\mathcal{O}((\log(N/B))/B)$  where  $N$  is the number of distinct elements in the priority-queue.
- In Section 3, we use the bucket heap to obtain a cache-oblivious shortest-path algorithm for undirected graphs with non-negative edge weights that incurs  $\mathcal{O}(V + (E/B) \cdot \log(E/B))$  I/Os, thereby matching the performance of the best cache-aware algorithm for this problem.
- In Section 4, we develop a new cache-oblivious algorithm for undirected BFS. The algorithm comes in two versions: The first one performs  $\mathcal{O}(\text{ST}(E) + \text{Sort}(E) + \frac{E}{B} \cdot \log V + \sqrt{V \cdot E/B})$  I/Os; the second one performs  $\mathcal{O}(\text{ST}(E) + \text{Sort}(E) + \frac{E}{B} \cdot \frac{1}{\varepsilon} \cdot \log \log V + \sqrt{V \cdot E/B} \cdot \sqrt{V \cdot B/E^\varepsilon})$  I/Os, where  $1 \geq \varepsilon > 0$ . Here,  $\text{ST}(E)$  denotes the cost of cache-obliviously finding a spanning tree.

## 2 The Bucket Heap

Our first contribution is a new cache-oblivious priority queue, called the *bucket heap*, which supports an UPDATE (a weak DECREASEKEY) operation and does so in the same I/O-bound as the tournament tree of [15]. In Section 3, we demonstrate that, using the bucket heap, the single-source shortest path algorithm of [15] becomes cache-oblivious and achieves the same performance as in the I/O-model:  $\mathcal{O}(V + (E/B) \cdot \log(E/B))$  I/Os. In analogy to the tournament tree, the bucket heap supports the following three operations:

UPDATE( $x, p$ ): If  $x$  is currently in the priority queue and has priority  $p'$ , its new priority becomes  $\min(p, p')$ . Otherwise, element  $x$  is inserted with priority  $p$ . (Operation UPDATE( $x, p$ ) is a combination of the INSERT( $x, p$ ) and DECREASEKEY( $x, p$ ) operations supported by standard internal-memory priority queues.)

DELETE( $x$ ): Element  $x$  is removed from the priority queue (provided it is currently in the priority queue).

DELETEMIN: The element with minimal priority is removed from the priority queue and reported.

Since every element  $x$  in the priority queue has an associated priority  $p$ , we refer to this element as element  $(x, p)$  unless the priority is irrelevant, in which case we denote the element simply as element  $x$ .

The bucket heap consists of  $2q + 1$  arrays  $\mathcal{B}_1, \dots, \mathcal{B}_q$  and  $\mathcal{S}_1, \dots, \mathcal{S}_{q+1}$ , where  $q$  varies over time, but is always at most  $\lceil \log_4 N \rceil$ . We call array  $\mathcal{B}_i$  the  $i$ -th *bucket* and array  $\mathcal{S}_i$  the  $i$ -th *signal buffer* (or *buffer* for short). The *capacity* of bucket  $\mathcal{B}_i$  is  $2^{2i}$ ; buffer  $\mathcal{S}_i$  has capacity  $2^{2i-1}$ . In order to allow for temporary overflow, we allocate twice as much space, that is,  $2^{2i+1}$  memory locations for bucket  $\mathcal{B}_i$  and  $2^{2i}$  memory locations for buffer  $\mathcal{S}_i$ . We store all buckets and buffers consecutively in memory, in the following order:  $\mathcal{S}_1, \mathcal{B}_1, \mathcal{S}_2, \mathcal{B}_2, \dots, \mathcal{S}_q, \mathcal{B}_q, \mathcal{S}_{q+1}$ .

Buffers  $\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_{q+1}$  store three types of signals, which we use to update the contents of the heap in a lazy manner: An  $\text{UPDATE}(x, p)$  signal is used to perform the  $\text{UPDATE}(x, p)$  operation above; a  $\text{DELETE}(x)$  signal performs the  $\text{DELETE}(x)$  operation above; a  $\text{PUSH}(x, p)$  signal is used to push elements from a buffer  $\mathcal{B}_i$  to a buffer  $\mathcal{B}_{i+1}$  when  $\mathcal{B}_i$  overflows. Every signal has a time stamp corresponding to the time when the operation posting this signal was performed. The time stamp of an element in a bucket is the time stamp of the  $\text{UPDATE}$  signal that led to its insertion.

The three priority queue operations are implemented as follows: A  $\text{DELETETEMIN}$  operation first uses the  $\text{FILL}$  operation described below to make sure that bucket  $\mathcal{B}_1$  is non-empty. Then Observation 1 below ensures that  $\mathcal{B}_1$  contains the element with lowest priority. This element is removed and returned. An  $\text{UPDATE}(x, p)$  operation inserts an  $\text{UPDATE}(x, p)$  signal into buffer  $\mathcal{S}_1$  and then applies the  $\text{EMPTY}$  operation below to  $\mathcal{S}_1$ . Similarly, a  $\text{DELETE}(x)$  operation inserts a  $\text{DELETE}(x)$  signal into  $\mathcal{S}_1$  and empties  $\mathcal{S}_1$ . Essentially, all the work to update the contents of the bucket heap is delegated to two auxiliary procedures: Procedure  $\text{EMPTY}(\mathcal{S}_i)$  empties the signal buffer  $\mathcal{S}_i$ , applies these signals to bucket  $\mathcal{B}_i$ , and inserts appropriate signals into buffer  $\mathcal{S}_{i+1}$ . If this leads to an overflow of buffer  $\mathcal{S}_{i+1}$ , the procedure is applied recursively to  $\mathcal{S}_{i+1}$ . Procedure  $\text{FILL}(\mathcal{B}_i)$  fills an underfull bucket  $\mathcal{B}_i$  with the smallest  $2^{2i}$  elements in buckets  $\mathcal{B}_i, \dots, \mathcal{B}_q$ .

## 2.1 Correctness

We say that a priority queue is correct if, given a sequence  $o_1, o_2, \dots, o_t$  of priority queue operations, every  $\text{DELETETEMIN}$  operation  $o_i$  returns the smallest element in the set  $\mathcal{O}_{i-1}$  constructed by operations  $o_1, \dots, o_{i-1}$  according to their definitions at the beginning of this section. We make the following two observations about the behaviour of the heap update operations:

**Observation 1** *For every  $1 \leq i < q$  and any two elements  $(x, p) \in \mathcal{B}_i$  and  $(y, p') \in \mathcal{B}_{i+1}$ ,  $p \leq p'$ .*



**EMPTY( $\mathcal{S}_i$ )**

- 1 If  $i = q + 1$ , increase  $q$  by one and create two new arrays  $\mathcal{B}_q$  and  $\mathcal{S}_{q+1}$ .
- 2 Scan bucket  $\mathcal{B}_i$  to determine the maximal priority  $p'$  of the elements in  $\mathcal{B}_i$ . (If  $i = q$  and  $\mathcal{S}_{q+1}$  is empty, then  $p' = \infty$ . Otherwise, if  $\mathcal{B}_i$  is empty, then  $p' = -\infty$ .)
- 3 Scan buckets  $\mathcal{S}_i$  and  $\mathcal{B}_i$  simultaneously and perform the following operations for each signal in  $\mathcal{S}_i$ :
  - UPDATE( $x, p$ ): If there is an element  $(x, p'')$  in  $\mathcal{B}_i$ , replace it with  $(x, \min(p, p''))$  and mark the UPDATE( $x, p$ ) signal in  $\mathcal{S}_i$  as handled. If  $x$  is not in  $\mathcal{B}_i$ , but  $p \leq p'$ , insert  $(x, p)$  into  $\mathcal{B}_i$  and replace the UPDATE( $x, p$ ) signal with a DELETE( $x$ ) signal. If  $x$  is not in  $\mathcal{B}_i$  and  $p > p'$ , do nothing.
  - PUSH( $x, p$ ): If there is an element  $(x, p'')$  in  $\mathcal{B}_i$ , replace it with  $(x, p)$ . Otherwise, insert  $(x, p)$  into  $\mathcal{B}_i$ . Mark the PUSH( $x, p$ ) signal as handled.
  - DELETE( $x$ ): If element  $x$  is in  $\mathcal{B}_i$ , delete it.<sup>3</sup>
- 4 If  $i < q$  or  $\mathcal{S}_{i+1}$  is non-empty, scan buffers  $\mathcal{S}_i$  and  $\mathcal{S}_{i+1}$  and insert all unhandled signals in  $\mathcal{S}_i$  into  $\mathcal{S}_{i+1}$ .  $\mathcal{S}_i \leftarrow \emptyset$
- 5 **if**  $|\mathcal{B}_i| > 2^{2i}$
- 6     **then** Find the  $2^{2i}$ -th smallest priority  $p$  in  $\mathcal{B}_i$ .
- 7         Scan bucket  $\mathcal{B}_i$  and buffer  $\mathcal{S}_{i+1}$  simultaneously and remove all elements with priority greater than  $p$  from  $\mathcal{B}_i$ ; for each removed element  $(x, p'')$ , insert a PUSH( $x, p''$ ) signal into  $\mathcal{S}_{i+1}$ .
- 8     **if**  $|\mathcal{B}_i| > 2^{2i}$
- 9         **then** Scan bucket  $\mathcal{B}_i$  and buffer  $\mathcal{S}_{i+1}$  simultaneously and remove  $|\mathcal{B}_i| - 2^{2i}$  elements with priority  $p$  from  $\mathcal{B}_i$ ; for each removed element  $(x, p)$ , insert a PUSH( $x, p$ ) signal into  $\mathcal{S}_{i+1}$ .
- 10 **if**  $|\mathcal{S}_{i+1}| > 2^{2i+1}$
- 11     **then** EMPTY( $\mathcal{S}_{i+1}$ )

**FILL( $\mathcal{B}_i$ )**

- 1 EMPTY( $\mathcal{S}_i$ )
- 2 **if**  $|\mathcal{B}_{i+1}| < 2^{2i}$  and  $i < q$
- 3     **then** FILL( $\mathcal{B}_{i+1}$ )
- 4 Find the  $(2^{2i} - |\mathcal{B}_i|)$ -th smallest priority  $p$  in  $\mathcal{B}_{i+1}$ .
- 5 Scan  $\mathcal{B}_i$  and  $\mathcal{B}_{i+1}$  and move all elements with priority less than  $p$  from  $\mathcal{B}_{i+1}$  to  $\mathcal{B}_i$ .
- 6 Scan  $\mathcal{B}_i$  and  $\mathcal{B}_{i+1}$  again and move the correct number of elements with priority  $p$  from  $\mathcal{B}_{i+1}$  to  $\mathcal{B}_i$  so that  $\mathcal{B}_i$  contains  $2^{2i}$  elements or  $\mathcal{B}_{i+1}$  is empty at the end.
- 7  $q \leftarrow \max\{j : \mathcal{B}_j \text{ or } \mathcal{S}_{j+1} \text{ is non-empty}\}$

**Observation 2** *When an element  $x \in \mathcal{B}_{i+1}$  moves to bucket  $\mathcal{B}_i$ , buffers  $\mathcal{S}_1, \dots, \mathcal{S}_{i+1}$  are empty.*

Using these two observations, we can prove the following two lemmas, which together establish the correctness of the bucket heap.

**Lemma 1** *Given a sequence  $o_1, o_2, \dots, o_t$  of priority queue operations and an index  $i$  such that  $o_i$  is a DELETETEMIN operation, the element returned by operation  $o_i$  is in  $\mathcal{O}_{i-1}$ .*

*Proof sketch.* Assume the contrary; that is, there exists a sequence  $o_1, \dots, o_t$  of priority queue operations and an index  $i$  such that operation  $o_i$  is a DELETETEMIN operation that returns an element  $(x, p)$  not in  $\mathcal{O}_{i-1}$ . In order to be returned by operation  $o_i$ , element  $(x, p)$  must have been inserted into a bucket by an UPDATE( $x, p$ ) operation. Let  $o_h$  be the latest such operation preceding  $o_i$ . Since we assume that  $(x, p) \notin \mathcal{O}_{i-1}$ , there has to exist an index  $j$ ,  $h < j < i$ , such that  $o_j$  is a DELETE( $x$ ) operation or a DELETETEMIN operation that returns an element  $(x, p')$ .

Assume that operation  $o_h$ , after percolating through buffers  $\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_k$ , inserts element  $(x, p)$  into bucket  $\mathcal{B}_k$ . Then, in order to reach bucket  $\mathcal{B}_1$ , element  $(x, p)$  has to be moved from bucket  $\mathcal{B}_k$  through buckets  $\mathcal{B}_{k-1}, \dots, \mathcal{B}_2$  to  $\mathcal{B}_1$ , using FILL operations. If there exists an index  $j$ ,  $h < j < i$ , such that operation  $o_j$  is a DELETE( $x$ ) operation, this operation inserts a DELETE( $x$ ) signal into buffer  $\mathcal{S}_1$ . Assume that, at this time,  $(x, p)$  is in a bucket  $\mathcal{B}_{k'}$ . Then, before  $(x, p)$  is moved to bucket  $\mathcal{B}_{k'-1}$  or returned by a DELETETEMIN operation, the DELETE( $x$ ) signal is applied to  $\mathcal{B}_{k'}$ , which leads to the deletion of element  $(x, p)$ , a contradiction.

So assume that none of  $o_{h+1}, o_{h+2}, \dots, o_{i-1}$  is a DELETE( $x$ ) operation. Then there has to exist an index  $j$ ,  $h < j < i$ , such that operation  $o_j$  is a DELETETEMIN operation that returns an element  $(x, p')$ . If  $p' > p$ , element  $(x, p)$  would have to reach bucket  $\mathcal{B}_1$  before  $(x, p')$  because there is no DELETE( $x$ ) signal pending for  $(x, p)$  and by Observation 1. Hence,  $p' \leq p$ . Now consider the order of the two operations  $o_h = \text{UPDATE}(x, p)$  and  $o_{h'} = \text{UPDATE}(x, p')$  that inserted elements  $(x, p)$  and  $(x, p')$  into the priority queue. If  $h' < h$ , then  $(x, p')$  is inserted into a bucket  $\mathcal{B}_{k'}$  before  $o_h$  is applied to bucket  $\mathcal{B}_k$ . Hence,  $o_h$  would find a copy of  $x$  with lower priority, and element  $(x, p)$  would never be inserted into any bucket, a contradiction. If  $h' > h$ , then either  $(x, p')$  is inserted into the bucket  $\mathcal{B}_k$  currently containing  $(x, p)$ , thereby replacing  $(x, p)$ ; or  $(x, p')$  is inserted into a bucket  $\mathcal{B}_{k'}$ ,  $k' < k$ , which leads to the insertion of a DELETE( $x$ ) signal into  $\mathcal{S}_{k'+1}$ . This prevents  $(x, p)$  from moving from  $\mathcal{B}_k$  to  $\mathcal{B}_{k-1}$  before being deleted. Hence, operation  $o_i$  cannot return element  $(x, p)$ .  $\square$

**Lemma 2** *Given a sequence  $o_1, o_2, \dots, o_t$  of priority queue operations and an index  $i$  such that  $o_i$  is a DELETEMIN operation, the element returned by operation  $o_i$  is the one with smallest priority in  $\mathcal{O}_{i-1}$ .*

*Proof sketch.* By Observation 1, the element returned by operation  $o_i$  is the one with the smallest priority in all buckets at the time of its deletion. By Lemma 1, the returned element  $(x, p)$  is in  $\mathcal{O}_{i-1}$ . If there is an element  $(y, p')$  with lower priority in  $\mathcal{O}_{i-1}$ , this element must have been inserted by an UPDATE operation  $o_j$ ,  $j < i$ . The only way this element is not returned by operation  $o_i$  is if there is an operation  $o_k$ ,  $j < k < i$ , that is a DELETE( $y$ ) operation or a DELETEMIN operation that returns  $y$ . But if this is the case,  $(y, p')$  is not in  $\mathcal{O}_{i-1}$ , a contradiction.  $\square$

## 2.2 Analysis

The efficiency of our data structure is based on the following *order invariant*:

The elements in every bucket or buffer are sorted primarily by their IDs and secondarily by their time stamps.<sup>4</sup>

This invariant allows us to perform updates by scanning buckets and buffers as in the description of procedures EMPTY and FILL. The amortized cost per scanned element is hence  $\mathcal{O}(1/B)$ . In our analysis of the amortized complexity of the priority queue operations, we assume that  $M = \Omega(B)$ , large enough to hold the first  $\log_4 B$  buckets and buffers plus one cache block per stream that we scan. Under this assumption, operations UPDATE, DELETE, and DELETEMIN, excluding the calls to EMPTY and FILL, do not cause any I/Os because bucket  $\mathcal{B}_1$  and buffer  $\mathcal{S}_1$  are always in cache. We have to charge the I/Os incurred by EMPTY and FILL operations to UPDATE and DELETE operations in such a manner that no operation is charged for more than  $\mathcal{O}((\log(N/B))/B)$  I/Os. To achieve this, we assign credits to the signals in buffers  $\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_{q+1}$  and define the following potential function:

$$\Phi = \sum_{i=\log_4 B}^q (|\mathcal{B}_i| \cdot (i - \log_4 B) + 2^{2i})$$

Every DELETE( $x$ ) signal inserted into buffer  $\mathcal{S}_1$  receives a credit of  $2c = 2\lceil \log_4(N/B) \rceil$ ; every UPDATE( $x, p$ ) signal receives a credit of  $5c$ . We prove that these credits can pay for the I/Os incurred by all updates of the priority queue, where  $B$  credits are required to pay for  $\mathcal{O}(1)$  I/Os. Hence, the amortized cost of every DELETEMIN operation is 0; an UPDATE or DELETE incurs  $\mathcal{O}((\log(N/B))/B)$  I/Os amortized.

---

<sup>4</sup>Since every piece of data must be encoded as a binary string in practice, it is reasonable to assume that there is a total order defined over the set of identifiers of data elements.

We divide EMPTY operations into two groups: A *regular* EMPTY operation is triggered by an EMPTY( $\mathcal{S}_1$ ) operation. An EMPTY operation triggered by a FILL( $\mathcal{B}_1$ ) operation is *early*. We maintain the invariant that every UPDATE signal in a buffer  $\mathcal{S}_i$  has credit at least  $5c - \max(0, i - \log_4 B)$ , every PUSH signal has credit at least  $3c - \max(0, i - \log_4 B)$ , and every DELETE signal has credit at least  $2c - \max(0, i - \log_4 B)$ .

Now consider a regular EMPTY( $\mathcal{S}_i$ ) operation, where  $i \leq q$ . If  $i < \log_4 B$ , the operation causes no I/Os because  $\mathcal{S}_i$ ,  $\mathcal{S}_{i+1}$ , and  $\mathcal{B}_i$  are in cache; and the potential does not change. If  $i \geq \log_4 B$ , the cost is bounded by  $\mathcal{O}(2^{2i-1}/B)$  because only buffers  $\mathcal{S}_i$  and  $\mathcal{S}_{i+1}$  and bucket  $\mathcal{B}_i$  are scanned. Let  $k$  be the increase of the size of bucket  $\mathcal{B}_i$ , let  $u$ ,  $p$ , and  $d$  be the number of UPDATE, PUSH, and DELETE operations in  $\mathcal{S}_i$ , and let  $u'$ ,  $p'$ , and  $d'$  be the number of such operations inserted into  $\mathcal{S}_{i+1}$ . Then we have  $u - u' + d = d'$  and  $k + p' \leq (u - u') + p$ . The total number of credits we need to pay for the cost of the EMPTY operation, for the potential increase, and for the credit of elements inserted into  $\mathcal{S}_{i+1}$  is  $2^{2i-1} + k(i - \log_4 B) + u'(5c - \max(0, i + 1 - \log_4 B)) + p'(3c - \max(0, i + 1 - \log_4 B)) + d'(2c - \max(0, i + 1 - \log_4 B))$ . Since  $u + p + d > 2^{2i-1}$  in the case of a regular EMPTY operation, this is easily bounded by  $u(5c - \max(0, i - \log_4 B)) + p(3c - \max(0, i - \log_4 B)) + d(2c - \max(0, i - \log_4 B))$ , which is the amount of credit carried by the signals in  $\mathcal{S}_i$ .

If  $i = q + 1$ , then the cost of a regular EMPTY( $\mathcal{S}_i$ ) operation remains the same, but there is an additional  $2^{2i}$  increase in potential. However, in this case, there are no signals that are inserted into  $\mathcal{S}_{i+1}$ , so that the amount of credit we would have assigned to these signals can be used to pay for this increase in potential. (A rigorous analysis appears in the full paper.)

The cost of FILL and early EMPTY operations will be paid by the resulting decrease of the potential  $\Phi$ . Consider a FILL( $\mathcal{B}_1$ ) operation, and let  $j$  be the highest index such that a FILL( $\mathcal{B}_j$ ) operation is triggered by this FILL( $\mathcal{B}_1$ ) operation. Then the cost of all EMPTY and FILL operations triggered by this FILL( $\mathcal{B}_1$ ) operation is  $\mathcal{O}(2^{2j}/B)$ . If there are new elements inserted into the buckets during the EMPTY operations, the resulting increase of potential can be paid from the credit of the corresponding UPDATE or PUSH signals. Hence, it suffices to show that, excluding this increase of potential due to element insertions, the potential decreases by  $\Omega(2^{2j})$ . We distinguish two cases. If  $q$  does not change, then the FILL( $\mathcal{B}_j$ ) operation moves at least  $3 \cdot 2^{2j-2}$  elements from  $\mathcal{B}_{j+1}$  to  $\mathcal{B}_j$ , which results in the desired potential change. If  $q$  decreases, then  $q > j$  before the FILL( $\mathcal{B}_1$ ) operation and  $q$  decreases by at least one. This results in a potential decrease of at least  $2^{2q} \geq 2^{2j}$ . Hence, we obtain the following result.

**Theorem 1** *The bucket heap supports operations UPDATE, DELETE, and DELETEMIN at an amortized cost of  $\mathcal{O}((\log(N/B))/B)$  I/Os.*

### 3 Cache-Oblivious Undirected Shortest Paths

The shortest-path algorithm of [15] is an adaptation of Dijkstra’s algorithm [12]. It uses a priority queue  $Q$  that stores the vertices of  $G$  whose distance from the source  $s$  is not known yet; the priority of a vertex  $v$  is equal to the length of the currently known shortest path from  $s$  to  $v$ ; it is  $\infty$  if no path from  $s$  to  $v$  is known. We call the vertices in  $G$  whose distances from  $s$  are known *settled*. As long as there is still a vertex in  $Q$ , the vertex  $v$  with lowest priority is removed, its current priority is fixed as the final distance from  $s$  to  $v$ , and all edges incident to  $v$  are relaxed by performing UPDATE operations on  $Q$ , for all of  $v$ ’s neighbours. A problem with this approach is that a settled neighbour  $w$  is re-inserted into  $Q$  by such an UPDATE operation. Kumar and Schwabe use a second priority queue to ensure that this re-inserted vertex  $w$  is removed from  $Q$  using a DELETE( $w$ ) operation before it can be settled a second time.

Kumar and Schwabe argue that their algorithm performs  $\mathcal{O}(E)$  priority queue operations. Apart from these, the algorithm accesses the adjacency list of every vertex  $v$  once, namely when  $v$  is settled, to retrieve the set of edges incident to  $v$  and relax them. Accessing this adjacency list incurs  $\mathcal{O}(1 + \deg(v)/B)$  I/Os. Summing over all vertices, the cost of accessing all adjacency lists is  $\mathcal{O}(V + E/B)$ . The  $\mathcal{O}(E)$  priority queue operations cause  $\mathcal{O}((E/B) \cdot \log(E/B))$  I/Os if we use the bucket heap as the priority queue. This proves the following result.

**Theorem 2** *There exists a cache-oblivious algorithm that solves the single-source shortest path problem on an undirected graph  $G = (V, E)$  with non-negative edge weights and incurs at most  $\mathcal{O}(V + (E/B) \cdot \log(E/B))$  I/Os.*

### 4 Cache-Oblivious Breadth-First Search

In this section, we develop a cache-oblivious version of the undirected BFS-algorithm from [16]. As in [16], the actual BFS-algorithm is the algorithm from [19], which generates the BFS levels  $L_i$  one by one, exploiting that, in an undirected graph,  $L_{i+1} = \mathcal{N}(L_i) \setminus (L_i \cup L_{i-1})$ , where  $\mathcal{N}(S)$  denotes the set of nodes<sup>5</sup> that are neighbours of nodes in  $S$ , and  $L_i$  denotes the nodes of the  $i$ ’th level of the BFS traversal, that is, the nodes at distance  $i$  from the source vertex  $s$ . The algorithm from [19] relies only on sorting and scanning and, hence, is straightforward to implement cache-obliviously; this gives the cache-oblivious  $\mathcal{O}(V + \text{Sort}(E))$  result mentioned in Table 1.

The speed-up in [16] compared to [19] is achieved using a data structure which, for a query set  $S$ , returns  $\mathcal{N}(S)$  in an I/O-efficient manner. To

---

<sup>5</sup>As shorthand for  $\mathcal{N}(\{v\})$  we will use  $\mathcal{N}(v)$ .

do so, the data structure exploits the fact that the query sets are the levels  $L_0, L_1, L_2, \dots$  of a BFS traversal. Our contribution is to develop a cache-obliviously version of this data structure.

#### 4.1 The Data Structure

To construct the data structure, we first build a spanning tree for  $G$  and then construct an Euler tour for the tree (which is a list containing the edges of the Euler tour in the order they appear in the tour) using the algorithm described in [4]. Next, we assign to each node  $v$  the rank in the Euler tour of the first occurrence of the node, which is done by a traversal of the Euler tour and a sorting step. We denote this value as  $r(v)$ . The Euler tour has length  $2V - 1$ ; so  $r(v) \in [0; 2V - 2]$ . A central observation used in [16] as well as in this paper is the following:

**Observation 3** *If for two nodes  $u$  and  $v$ , the values  $r(v)$  and  $r(u)$  differ by  $d$ , then a section of the Euler tour is a path in  $G$  of length  $d$  that connects  $u$  and  $v$ ; hence,  $d$  is an upper bound on the distance between their BFS levels.*

Let  $g_0 < g_1 < g_2 < \dots < g_h$  be an increasing sequence of  $h + 1$  integers where  $g_0 = 1$ ,  $g_{h-1} < 2V - 2 \leq g_h$ , and  $g_i$  divides  $g_{i+1}$ . We will later consider two specific sequences, namely  $g_i = 2^i$  and one for which  $g_i = \Theta(2^{(1+\varepsilon)^i})$ , where  $1 \geq \varepsilon > 0$ . For each integer  $g_i$ , we can partition the nodes into groups of at most  $g_i$  nodes each by letting the  $k$ 'th group  $V_{ki}$  be all nodes  $v$  for which  $kg_i \leq r(v) < (k + 1)g_i$ . We call a group  $V_{ki}$  of nodes a  $g_i$ -node-group and call its set of adjacency lists  $\mathcal{N}(V_{ki})$  a  $g_i$ -edge-group. Since  $g_i$  divides  $g_{i+1}$ , the groups form a hierarchy of  $h + 1$  levels, with level  $h$  containing one group with all nodes and level 0 containing  $2E - 1$  groups of at most one node. Note that there is no strong relation between  $g_i$  and the number of nodes in a  $g_i$ -node-group or the number of edges in a  $g_i$ -edge-group. In particular,  $g_i$ -edge-groups may be much larger than  $g_i$ .

The data structure consists of  $h$  levels  $G_1, \dots, G_h$ , where each level stores a subset of the adjacency lists of the graph  $G$ . Each adjacency list  $\mathcal{N}(v)$  will appear in exactly one of the levels, unless it has been removed from the structure. Recall that the query sets of the BFS-algorithm are the BFS-levels  $L_0, L_1, L_2, \dots$ ; so each node  $v$  is part of a query set  $S$  from the BFS-algorithm exactly once. Its adjacency list  $\mathcal{N}(v)$  will leave the data structure when this happens. Initially, all adjacency lists are in level  $h$ . Over time, the query algorithm for the data structure moves the adjacency list of each node from higher numbered to lower numbered levels, until the adjacency list eventually reaches level  $G_1$  and is removed.

A high-level description of our query algorithm is shown in Figure 1. It is a recursive procedure that takes as input a query set  $S$  of nodes and a

```

GETEDGELISTS( $S, i$ )
1   $S' = S \setminus \{v \in V \mid \mathcal{N}(v) \text{ is stored in level } G_i\}$ 
2  if  $S' \neq \emptyset$ :
3       $X = \text{GETEDGELISTS}(S', i + 1)$ 
4      for each  $g_i$ -edge-group  $g$  in  $X$ 
5          insert  $g$  in  $G_i$ 
6  for each  $g_{i-1}$ -edge-group  $\gamma$  in  $G_i$  containing  $\mathcal{N}(v)$  for some  $v \in S$ 
7      remove  $\gamma$  from  $G_i$ 
8      include  $\gamma$  in the output set

```

Figure 1: The query algorithm.

level number  $i$  to query. The output consists of the  $g_{i-1}$ -edge-groups stored at level  $G_i$  for which the corresponding  $g_{i-1}$ -node-group contains one or more nodes in  $S$ . The BFS-algorithm will query the data structure by calling `GETEDGELISTS( $S, 1$ )`, which will return  $\mathcal{N}(v)$ , for all  $v$  in  $S$ . (Recall that  $g_0 = 1$ ; so a non-empty  $g_0$ -edge-group is the adjacency list of a single node).

Next we describe how we represent a level  $G_i$  so that `GETEDGELISTS` can be performed efficiently. First of all, it follows from Figure 1 that at any time, the edges stored in level  $G_i$  constitute a set of  $g_i$ -edge-groups from each of which zero or more  $g_{i-1}$ -edge-groups have been removed. Since  $g_{i-1}$  divides  $g_i$ , the part of a  $g_i$ -edge-group residing at level  $G_i$  is a collection of  $g_{i-1}$ -edge-groups.

We store the adjacency lists of level  $G_i$  in an array  $B_i$ . Each  $g_{i-1}$ -edge-group of the level is stored in consecutive locations of  $B_i$ , and the adjacency lists  $\mathcal{N}(v)$  of the nodes  $v$  in a group occupy these locations in order of increasing ranks  $r(v)$ . The  $g_{i-1}$ -edge-groups of each  $g_i$ -edge-group are also stored in order of increasing ranks of the nodes involved, but empty location may exist between the  $g_{i-1}$ -edge-groups. However, the entire array  $B_i$  has a number of locations which is at most a constant times the number of edges it contains. This will require  $B_i$  to shrink and grow appropriately. The arrays  $B_1, B_2, \dots, B_h$  will be laid out in  $\mathcal{O}(E)$  consecutive memory locations. For purposes of exposition, we defer the discussion of how this shrinking and growing can be achieved at sufficiently low cost while maintaining this layout to Section 4.3.

Note that there are no restrictions on the order in which the  $g_i$ -edge-groups appear in among each other in  $B_i$ . To locate these groups, we keep an index  $A_i$  at every level, which is an array of entries  $(k, p)$ , one for every  $g_i$ -edge-group present in  $G_i$ . The value  $k$  is the number of the corresponding  $g_i$ -node-group  $V_{ki}$ , and  $p$  is a pointer to the start of the  $g_i$ -edge-group in  $B_i$ . The entries of  $A_i$  are sorted by their first components. The indexes  $A_1, A_2, \dots, A_h$  occupy consecutive locations of one of two arrays  $A'$  and  $A''$  of size  $\mathcal{O}(V)$ .



Finally, every  $g_i$ -edge-group  $g$  of a level will contain an index of the  $g_{i-1}$ -edge-groups it presently contains. This index consists of the first and last edge of each  $g_{i-1}$ -edge-group  $\gamma$  together with pointers to the first and last locations of the rest of  $\gamma$ . These edges are kept at the front of  $g$ , in the same order as the  $g_{i-1}$ -edge-groups to which they belong. For simplicity, we assume that the data type for an edge has room for such a pointer.

We now describe how each step of `GETEDGELISTS` is performed. We assume that every query set  $S$  of nodes is sorted by assigned rank  $r(v)$ , which can be ensured by sorting the initial query set before the first call. In Line 1 of the algorithm, we find  $S'$  by simultaneously scanning  $S$  and  $A_i$ , using that, if  $(k, p)$  is an entry of  $A_i$ , all nodes  $v \in S$  for which  $kg_i \leq r(v) < (k+1)g_i$  will have  $\mathcal{N}(v)$  residing in the  $g_i$ -edge-group pointed to by  $p$  (otherwise,  $\mathcal{N}(v)$  would have been found at an earlier stage of the recursion). In other words,  $S'$  is the subset of  $S$  not covered by the entries in  $A_i$ .

In line 5, when a  $g_i$ -edge-group  $g$  is to be inserted into level  $G_i$ , the index of the  $g_{i-1}$ -edge-groups of  $g$  is generated by scanning  $g$ , and  $g$  (now with the index) is appended to  $B_i$ . An entry for  $A_i$  is generated. When the for-loop in Line 4 is finished, the set of new  $A_i$  entries are sorted by their first components and merged with the current  $A_i$ . Specifically, the merging writes  $A_i$  to  $A''$  if  $A_i$  currently occupies  $A'$ , and vice versa, implying that the location of the entire set of  $A_i$ 's alternates between  $A'$  and  $A''$  for each call `GETEDGELISTS( $S$ , 1)`.

In Line 6, we scan  $S$  and the updated  $A_i$  to find the entries  $(k, p)$  pointing to the relevant  $g_i$ -edge-groups of the updated  $B_i$ . During the scan, we generate for each of these groups  $g$  each of the  $g_{i-1}$ -edge-groups  $\gamma$  inside  $g$  that contains one or more nodes from  $S$ , a pair  $(v, p)$ , where  $v$  is the first node in the group. These pairs are now sorted in reverse lexicographic order (the second component is most significant), so that the  $g_i$ -edge-groups can be accessed in the same order as they are located in  $B_i$ . For each such group  $g$ , we scan its index to find the relevant  $g_{i-1}$ -edge-groups and access these in the order of the index. Each  $g_{i-1}$ -edge-group is removed from its location in  $B_i$  (leaving empty positions) and placed in the output set. We also remove its entry in the index of  $g$ . The I/O-bound of this process is the minimum I/O-bound of a scan of  $B_i$  and a scan of each of the moved  $g_{i-1}$ -edge-groups.

## 4.2 Analysis

In the following we analyze the number of I/Os performed by our cache-oblivious BFS-algorithm (assuming that the management of the layout of the  $B_i$ 's can be done efficiently, which will be easier to discuss in Section 4.3 once we have the analysis of the main structure available).

The basic BFS-algorithm from [19] scans each BFS-level  $L_i$  twice: once while constructing  $L_{i+1}$  and once while constructing  $L_{i+2}$ , causing a total of



$\mathcal{O}(V/B)$  I/Os for all lists  $L_i$ . Edges extracted from the data structure storing the adjacency lists are sorted once and scanned once for filtering out duplicates and already discovered nodes, causing a total of  $\mathcal{O}(\text{Sort}(E) + E/B)$  I/Os. It follows that the basic algorithm requires a total of  $\mathcal{O}(\text{Sort}(E))$  I/Os.

We now turn to the I/Os performed during queries of the data structure storing the adjacency lists. The cost for constructing the initial spanning tree is  $\mathcal{O}(\text{ST}(E))$ ; the Euler tour can be constructed in  $\mathcal{O}(\text{Sort}(V))$  I/Os [4]. Assigning ranks to nodes and labelling the edges with the assigned ranks requires further  $\mathcal{O}(\text{Sort}(E))$  I/Os. The total preprocessing of the data structure hence costs  $\mathcal{O}(\text{ST}(E) + \text{Sort}(E))$  I/Os.

For each query from the basic BFS-algorithm, the query algorithm for the data structure accesses the  $A_i$  and  $B_i$  lists. We first consider the number of I/Os for handling the  $A_i$  lists. During a query, the algorithm scans each  $A_i$  list at most a constant number of times: to identify which  $g_i$ -edge-groups to extract recursively from  $B_{i+1}$ ; to merge  $A_i$  with new entries extracted recursively; and to identify the  $g_{i-1}$ -edge-groups to extract from  $B_i$ . The number of distinct  $g_i$ -edge-groups is  $2V/g_i$ . Each group is inserted into level  $G_i$  at most once. By Observation 3, when a  $g_i$ -edge-group is inserted into level  $G_i$ , it will become part of an initial query set  $S$  within  $g_i$  queries from the basic BFS-algorithm, that is, within the next  $g_i$  BFS-levels; at this time, it will be removed from the structure. In particular, it will reside in level  $G_i$  for at most  $g_i$  queries. We conclude that the total cost of scanning  $A_i$  during the run of the algorithm is  $\mathcal{O}((2V/g_i) \cdot g_i/B)$ , implying a total number of  $\mathcal{O}(h \cdot V/B)$  I/Os for scanning all  $A_i$  lists. This bound holds because the  $A_i$ 's are stored in consecutive memory locations, which can be considered to be scanned in a single scan during a query from the basic BFS-algorithm. Since each node is part of exactly one query set of the basic BFS-algorithm, the total I/O cost for scanning the  $S$  sets during all recursive calls is also  $\mathcal{O}(h \cdot V/B)$ .

We now bound the sorting cost caused during the recursive extraction of groups. The pointer to each  $g_i$ -edge-group participates in two sorting steps: When the group is moved from level  $i + 1$  to level  $i$ , the pointer to the group participates in the sorting of  $A_{i+1}$  before accessing  $B_{i+1}$ ; when the  $g_i$ -edge-group has been extracted from  $B_{i+1}$ , the pointer is involved in the sorting of the pointers to all extracted groups before they are merged into  $A_i$ . We conclude that the total sorting cost is bounded by  $\mathcal{O}(\sum_{i=1}^h \text{Sort}(2V/g_i))$  which is  $\mathcal{O}(\text{Sort}(V))$ , since  $g_i$  is at least exponentially increasing for both of the two sequences considered.

Finally, we need to argue about the I/O cost of accessing the  $B_i$  lists. For each query of the basic BFS-algorithm, these will be accessed in the order  $B_h, B_{h-1}, \dots, B_1$ , since the  $B_i$  lists are accessed after the recursive call to GETEDGELists. Let  $t$  be an integer, where  $1 \leq t \leq h$ . The cost of accessing  $B_t, \dots, B_1$  during a query is bounded by the cost of scanning

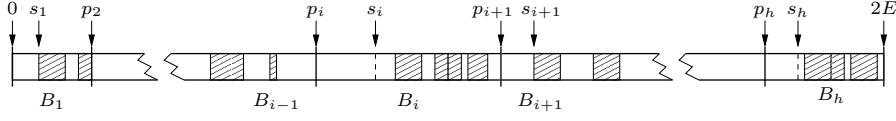


Figure 2: The memory layout of  $B_1, \dots, B_h$

$B_t, \dots, B_1$ . Since an edge in  $B_i$  can only remain in  $B_i$  for  $g_i$  queries from the basic BFS-algorithm, we get a bound on the total I/O cost for  $B_1, \dots, B_t$  of  $\mathcal{O}(\sum_{i=1}^t g_i \cdot E/B)$ , which is  $\mathcal{O}(g_t \cdot E/B)$  since  $g_i$  is at least exponentially increasing.

To bound the cost of accessing  $B_h, \dots, B_{t+1}$ , we note that the number of I/Os for moving a  $g_i$ -edge-group list containing  $k$  edges from  $B_{i+1}$  to  $B_i$  is bounded by  $\mathcal{O}(1 + k/B + g_{i+1}/(g_i \cdot B))$ , where  $g_{i+1}/g_i$  is the bound of the size of the index of a  $g_{i+1}$ -edge-group. Since the number of  $g_i$ -edge-groups is bounded by  $2V/g_i$ , it follows that the I/O cost for accessing  $B_h, \dots, B_{t+1}$  is bounded by  $\mathcal{O}(\sum_{i=t}^{h-1} (2V/g_i + E/B + (2V/g_i) \cdot g_{i+1}/(g_i \cdot B))) = \mathcal{O}(V/g_t + h \cdot E/B)$ , since  $g_{i+1} \leq g_i^2$  for both of the two sequences considered. The total cost of accessing all  $B_i$  is, hence,  $\mathcal{O}(g_t \cdot E/B + V/g_t + h \cdot E/B)$ , for all  $1 \leq t \leq h$ .

Adding all the bounds above gives a bound of  $\mathcal{O}(\text{ST}(E) + \text{Sort}(E) + g_t \cdot E/B + V/g_t + h \cdot E/B)$  on the total number of I/Os incurred by the query algorithm, for all  $1 \leq t \leq h$ .

For  $g_i = 2^i$ , we select  $g_t = \Theta(\sqrt{V \cdot B/E})$  and have  $h = \Theta(\log V)$ , so the I/O-bound becomes  $\mathcal{O}(\text{ST}(E) + \text{Sort}(E) + \frac{E}{B} \cdot \log V + \sqrt{V \cdot E/B})$ . For  $g_i = \Theta(2^{(1+\epsilon)^i})$ , we select the smallest  $g_t \geq \sqrt{V \cdot B/E}$ , i.e.  $g_t \leq \sqrt{V \cdot B/E}^{1+\epsilon}$ , and have  $h = \Theta(\frac{1}{\epsilon} \cdot \log \log V)$ , so the I/O-bound becomes  $\mathcal{O}(\text{ST}(E) + \text{Sort}(E) + \frac{E}{B} \cdot \frac{1}{\epsilon} \cdot \log \log V + \sqrt{V \cdot E/B} \cdot \sqrt{V \cdot B/E}^\epsilon)$ .

### 4.3 Layout Management

In this section, we describe how to efficiently maintain a layout of  $B_1, \dots, B_h$  in a single array of size  $2E$ . We will maintain a sequence of pointers  $0 = p_1 \leq s_1 \leq p_2 \leq s_2 \leq p_3 \leq \dots \leq p_h \leq s_h \leq p_{h+1} = 2E$  into the array, where  $p_2 \geq 2|B_1|$  and  $p_{i+1} - p_i = 2|B_i|$  for  $1 \leq i \leq h$ . Here,  $|B_i|$  denotes the total number of edges stored in  $B_i$ .

We maintain the invariant that all adjacency lists within  $B_i$  are stored within positions  $[s_i; p_{i+1})$ . The  $g_{i-1}$ -edge-groups will always be stored as a consecutive list of positions within  $B_i$ ; but between the individual  $g_{i-1}$ -edge-groups, there can be unused array positions. In Figure 2, we show the memory layout of the sets  $B_1, \dots, B_h$ . Each of the shaded areas within  $B_i$  represents a  $g_{i-1}$ -edge-group list.

In the initial state  $0 = p_1 = s_1 = p_2 = s_2 = \dots = p_h, s_h = E$  and  $p_{h+1} =$

$2E$ , and all adjacency lists are contained within  $B_h$  are stored consecutively within array positions  $[E; 2E)$ .

During a query,  $g_{i-1}$ -edge-groups will be moved from  $B_i$  to  $B_{i-1}$ , leaving unused array positions in  $B_i$ . When moving a total of  $x$  edges from  $B_i$  to  $B_{i-1}$ , we place these consecutively at array positions  $[p_i; p_i + x)$ . We then increase  $p_i$  by  $2x$  to maintain that  $p_{j+1} - p_j = 2|B_j|$  for all  $j$ . If this makes  $p_i$  larger than  $s_i$ , we compress  $B_i$  by scanning  $[s_i; p_{i+1})$ , moving  $B_i$  into the consecutive positions  $[p_{i+1} - |B_i|; p_{i+1})$ , and setting  $s_i = p_{i+1} - |B_i|$ . Since for the previous value of  $s_i$ , we had  $s_i < p_i = p_{i+1} - 2|B_i|$ , the value of  $s_i$  has increased by at least half of the scanned area. Charging the scanning cost to the advance of  $s_i$ , and using that each of the  $s_i$ , for  $i = 1, 2, \dots, h$ , can increase by  $2E$  in total, this cost is bounded by  $\mathcal{O}(h \cdot E/B)$  I/Os in total. Since each compression is triggered by the movement of  $x$  edges into  $[p_i; p_i + x)$ , and the area to be scanned starts below  $p_i + 2x$ , the single I/O initiating the scan can be included in the cost of moving the  $x$  edges, a cost which was bounded in the previous section.

After the compression of  $B_i$ , we will need to update  $A_i$ . This can be done efficiently if we add another index  $C_i$  to the structure. For level  $G_i$ , index  $C_i$  stores the pairs  $(k, p)$  previously discussed for the  $A_i$  array, but stores them in the same order as the  $g_i$ -edge-groups appear in  $B_i$ . The pointers in the pairs in  $A_i$  now point to the entries with the same  $k$  value in  $C_i$ . In short, information about the permutation of the  $g_i$ -edge-groups is stored in  $A_i$ , while information about their locations is stored in  $C_i$ . The  $C_i$  index can be updated during compression of  $B_i$  using a scan (whereas updating  $A_i$  in the absence of  $C_i$  would require sorting). Since  $C_i$  and  $A_i$  have the same size and the  $C_i$ 's are accessed in the order  $h, h-1, \dots, 2, 1$  during a query, the analysis of the cost of accessing the  $A_i$ 's from the previous section applies also to the  $C_i$ 's.

In summary, the desired layout of the  $B_i$ 's can be achieved without changing the I/O-bounds from the previous section.

In the described layout, we have chosen to store  $B_0, \dots, B_h$  in an array of size  $2E$ . We could also have chosen a smaller array of size  $(1 + \varepsilon) \cdot E$ , for some constant  $\varepsilon > 0$ , such that  $p_{i+1} - p_i = (1 + \varepsilon)|B_i|$ . The result would be an increased number of compression steps, implying a factor  $1/\varepsilon$  more I/Os for the repeated compression of the  $B_i$ 's.

## 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] A. Aggarwal and J. S. Vitter. The input/output complexity of sorting and related problems. *Communications of the ACM*, 31(9):1116–1127,

Sept. 1988.

- [3] L. Arge. External memory data structures. In *Proc. 9th Annual European Symposium on Algorithms*, volume 2161 of *LNCS*, pages 1–29. Springer Verlag, 2001.
- [4] L. Arge, M. A. Bender, E. D. Demaine, B. Holland-Minkley, and J. I. Munro. Cache-oblivious priority queue and graph algorithm applications. In *Proc. 34th Annual ACM Symposium on Theory of Computing*, pages 268–276. ACM Press, 2002.
- [5] L. Arge, G. Brodal, and L. Toma. On external-memory MST, SSSP and multi-way planar graph separation. In *Proc. 8th Scandinavian Workshop on Algorithmic Theory*, volume 1851 of *LNCS*, pages 433–447. Springer Verlag, 2000.
- [6] G. S. Brodal and R. Fagerberg. Cache oblivious distribution sweeping. In *Proc. 29th International Colloquium on Automata, Languages, and Programming*, volume 2380 of *LNCS*, pages 426–438. Springer Verlag, 2002.
- [7] G. S. Brodal and R. Fagerberg. Funnel heap - a cache oblivious priority queue. In *Proc. 13th Annual International Symposium on Algorithms and Computation*, volume 2518 of *LNCS*, pages 219–228. Springer Verlag, 2002.
- [8] G. S. Brodal and R. Fagerberg. On the limits of cache-obliviousness. In *Proc. 35th Annual ACM Symposium on Theory of Computing*, pages 307–315. ACM Press, 2003.
- [9] A. Buchsbaum, M. Goldwasser, S. Venkatasubramanian, and J. Westbrook. On external memory graph traversal. In *Proc. 11th ACM-SIAM Symposium on Discrete Algorithms*, pages 859–860. ACM-SIAM, 2000.
- [10] Y. Chiang, M. T. Goodrich, E. F. Grove, R. Tamassia, D. E. Vengroff, and J. S. Vitter. External-memory graph algorithms. In *Proc. 6th Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 139–149. ACM-SIAM, 1995.
- [11] E. D. Demaine. Cache-oblivious data structures and algorithms. In *Proc. EFF summer school on massive data sets*, LNCS. Springer Verlag, 2004, to appear.
- [12] E. W. Dijkstra. A note on two problems in connexion with graphs. *Num. Math.*, 1:269–271, 1959.

- [13] M. L. Fredman and R. E. Tarjan. Fibonacci heaps and their uses in improved network optimization algorithms. *Journal of the ACM*, 34:596–615, 1987.
- [14] M. Frigo, C. E. Leiserson, H. Prokop, and S. Ramachandran. Cache-oblivious algorithms. In *40th Annual Symposium on Foundations of Computer Science*, pages 285–297. IEEE Computer Society Press, 1999.
- [15] V. Kumar and E. J. Schwabe. Improved algorithms and data structures for solving graph problems in external memory. In *Proc. 8th Symposium on Parallel and Distrib. Processing*, pages 169–177. IEEE Computer Society Press, 1996.
- [16] K. Mehlhorn and U. Meyer. External-memory breadth-first search with sublinear I/O. In *Proc. 10th Annual European Symposium on Algorithms*, volume 2461 of *LNCS*, pages 723–735. Springer Verlag, 2002.
- [17] U. Meyer, P. Sanders, and J. F. Sibeyn, editors. *Algorithms for Memory Hierarchies*, volume 2625 of *LNCS*. Springer Verlag, 2003.
- [18] U. Meyer and N. Zeh. I/O-efficient undirected shortest paths. In *Proc. 11th Annual European Symposium on Algorithms*, volume 2832 of *LNCS*, pages 434–445. Springer Verlag, 2003.
- [19] K. Munagala and A. Ranade. I/O-complexity of graph algorithms. In *Proc. 10th Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 687–694. ACM-SIAM, 1999.
- [20] J. S. Vitter. External memory algorithms and data structures: Dealing with massive data. *ACM Computing Surveys*, 33(2):209–271, June 2001.

## Recent BRICS Report Series Publications

- RS-04-2 Gerth Stølting Brodal, Rolf Fagerberg, Ulrich Meyer, and Norbert Zeh. *Cache-Oblivious Data Structures and Algorithms for Undirected Breadth-First Search and Shortest Paths*. February 2004. 19 pp.
- RS-04-1 Luca Aceto, Willem Jan Fokkink, Anna Ingólfssdóttir, and Bas Luttik. *Split-2 Bisimilarity has a Finite Axiomatization over CCS with Hennessy's Merge*. January 2004. 16 pp.
- RS-03-53 Kyung-Goo Doh and Peter D. Mosses. *Composing Programming Languages by Combining Action-Semantics Modules*. December 2003. 39 pp. Appears in *Science of Computer Programming*, 47(1):2–36, 2003.
- RS-03-52 Peter D. Mosses. *Pragmatics of Modular SOS*. December 2003. 22 pp. Invited paper, published in Kirchner and Ringeissen, editors, *Algebraic Methodology and Software Technology: 9th International Conference, AMAST '02 Proceedings*, LNCS 2422, 2002, pages 21–40.
- RS-03-51 Ulrich Kohlenbach and Branimir Lambov. *Bounds on Iterations of Asymptotically Quasi-Nonexpansive Mappings*. December 2003. 24 pp.
- RS-03-50 Branimir Lambov. *A Two-Layer Approach to the Computability and Complexity of Real Numbers*. December 2003. 16 pp.
- RS-03-49 Marius Mikucionis, Kim G. Larsen, and Brian Nielsen. *Online On-the-Fly Testing of Real-time Systems*. December 2003. 14 pp.
- RS-03-48 Kim G. Larsen, Ulrik Larsen, Brian Nielsen, Arne Skou, and Andrzej Wasowski. *Danfoss EKC Trial Project Deliverables*. December 2003. 53 pp.
- RS-03-47 Hans Hüttel and Jiří Srba. *Recursive Ping-Pong Protocols*. December 2003. 21 pp. To appear in the proceedings of 2004 IFIP WG 1.7, ACM SIGPLAN and GI FoMSESS Workshop on Issues in the Theory of Security (WITS'04).