

Algorithmic Building Blocks for Asymmetric Memories

Yan Gu
Carnegie Mellon University
yan.gu@cs.cmu.edu

Yihan Sun
Carnegie Mellon University
yihans@cs.cmu.edu

Guy E. Blelloch
Carnegie Mellon University
guyb@cs.cmu.edu

Abstract

The future of main memory appears to lie in the direction of new non-volatile memory technologies that provide strong capacity-to-performance ratios, but have write operations that are much more expensive than reads in terms of energy, bandwidth, and latency. This asymmetry can have a significant effect on algorithm design, and in many cases it is possible to reduce writes at the cost of reads. In this paper we study which algorithmic techniques are useful in designing practical write-efficient algorithms. We focus on several fundamental algorithmic building blocks including unordered set/map implemented using hash tables, ordered set/map implemented using various binary search trees, comparison sort, and graph traversal algorithms including breadth-first search and Dijkstra’s algorithm. We introduce new algorithms and implementations that can reduce writes, and analyze the performance experimentally using a software simulator. Finally we summarize interesting lessons and directions in designing write-efficient algorithms.

1 Introduction

The future of main memory appears to lie in the non-volatile memory technologies that promise persistence, significantly lower energy costs, and higher density than the DRAM technology used in today’s main memories [HP15, Int15, MSCT14, Yol13]. Despite the advantages, a key property of such memory technologies, however, is their asymmetric read-write costs: compared to reads, writes can be much more expensive in terms of latency, bandwidth, and energy. Because bits are stored in these technologies as at rest “states” of the given material that can be quickly read but require physical change to update, this asymmetry appears fundamental. This motivates the need for *write-efficient* algorithms that largely reduce the number of writes compared to existing algorithms.

In the related work section, we review the literature on studying this read-write asymmetry on NAND Flash chips [BAT06, EGMP14, GT05, PS09] and algorithms targeting database operators [CGN11, Vig12, Vig14]. These works provide novel aspects on rethinking algorithm design. However, most of the papers either treat NVMs as external memories, or are based on hardware simulators for existing architecture, which may have many concerns that we will further discuss in the related work section.

Blelloch et al. [BDBF⁺16, BFG⁺15, BGST16] formally defined and analyzed several sequential and parallel computation models with good caching and scheduling guarantees. The models abstract such asymmetry between reads and writes, and can be used to analyze algorithms on future memory. The basic model, which is the Asymmetric RAM (ARAM), extends the well-known external-memory model [AV88] and parameterizes the asymmetry using ω , which corresponds to the cost of a write relative to a read to the non-volatile main memory. The cost of an algorithm on the ARAM, the **asymmetric I/O cost**, is the number of write transfers to the main memory multiplied by ω , plus the number of read transfers. This model captures different system consideration (latency, bandwidth, or energy) by simply plugging in different values of ω ,

and also allows algorithms to be analyzed theoretically. Based on this idea, many interesting algorithms (and lower bounds) are designed and analyzed by various recent papers [BDBF⁺16, BDBF⁺18, BFG⁺15, BFG⁺16, BGSS18, JS17].

Unfortunately, all of the analyses of such write-efficient algorithms are asymptotic, showing the upper and lower bounds on the complexity of these problems. Also, to prove the bounds, the theoretical models simplify the real architecture (e.g., without considering blocking of cache-lines or cache policies). It still remains unknown what the performance of these algorithms are in practice. In this paper, our goal is to show such performance on a number of fundamental algorithmic building blocks. We believe the lessons in designing and implementing them are useful for our community to use new memory in the future.

Contribution of this paper.

In this work, our goal is to bridge the gap between theory and practice. We try to study and understand which algorithmic techniques are useful in designing practical write-efficient algorithms. As the first paper of this kind, we focus on several of the most commonly-seen algorithmic building blocks in modern programming: unordered set/map implemented using **hash tables**, set/map implemented using **balanced binary search trees**, **comparison sort**, and graph traversal algorithms: **breadth-first search** for unweighted graphs and **Dijkstra’s algorithm** for weighted graphs¹.

Unfortunately, no non-volatile main memory is currently available, making it impossible to get real timings. Furthermore, details about latency and other parameters of the memory and how they will be incorporated into the architecture are also not available. This makes detailed cycle-level simulation (e.g., PTLsim [PTL], MARSSx86 [Mic] or ZSim [SK13]) of questionable utility. However, it is quite feasible to count the number of reads and write to main memory while simulating a variety of cache configurations. For I/O-bounded algorithms, these numbers can be used as reasonable proxies for both running time (especially when implemented in parallel) and energy consumption.² Moreover, conclusions drawn from these numbers can likely give insights into tradeoffs between reads and writes among different algorithms.

For these reasons, we propose a framework based on a software simulator that can efficiently and precisely measure the number of read and write transfers of an algorithm using different caching policies. We also consider variants in caching policies that might lead to improvements when read and write are not the same.

We also note that designing write-efficient algorithms falls in a high dimensional parameter space since the asymmetries on latency, bandwidth, and energy consumption between reads and writes are different. Here we abstract this as a single value ω . This value together with the cache size M and cache-line size B (set to be 64 bytes in this paper) form the parameter space of an algorithm.

Our framework provides a simple, clean and hardware-independent method to analyze and experiment the performance on the asymmetric memory. We investigate the algorithmic techniques and learn lessons from the experiments that generally apply for a reasonably large parameter space of ω , M and B . This framework also allows monitoring, reasoning and debugging the code easily, so it can remain useful even after the new hardware is available.

With the framework, we design, implement and discuss many different algorithms and data structures and their write-efficient implementations. Although some of the implementations are standard, like quicksort and the classic hash tables, many others, including the k -level hash tables, sample sort and phased Dijkstra, require careful algorithmic design, analysis, and coding. Under our cost measure, which is the asymmetric

¹The algorithms that are not I/O-bounded or use much fewer writes than reads are not discussed in this paper (e.g., matrix multiplication and other algorithms with a similar computation pattern [BFG⁺15]).

²The energy consumption of main memory is a key concern since it costs 25-50% energy on data centers and servers [LRR⁺03, MSG⁺12, LK14b].

I/O cost, we show better approaches on all problems we study in this paper, compared to the most basic and commonly-used ones on symmetric memories. We understand that there are more advanced versions of the algorithms and data structures discussed in this paper on some specific applications, and how to implement them write-efficiently is an interesting topic for future work.

With the algorithms and their experimental results, we draw many interesting algorithmic strategies and guidance in designing write-efficient algorithms. A common theme is to trade (more) reads for (fewer) writes (apparently it is hard to directly decrease the writes since this can improve the performance on symmetric memory as well and should have been investigated already). Some interesting lessons we learned and can be valuable to share are listed as follows, which can suggest some potential directions to design and engineer write-efficient algorithms in the future.

1. Indirect addressing is less problematic. In the classic setting, indirect addressing should be avoided if possible, since each addressing can be a random access to the memory. However, when writes are expensive, moving the entire data is costly, while indirect addressing only modifies the pointers (at the cost of a possible random access per lookup).
2. Multiple candidate positions for a single entry in a data structure can help. It can be a good option to use more reads per lookup but apply less frequent data movements, when the size of a data structure changes significantly. This is a common strategy we have applied in this paper to provide an algorithmic tradeoff between reads and writes.
3. It is usually worth to investigate existing algorithms that move or modify the data less. These algorithms can be less efficient in the symmetric setting due to various reasons (e.g., more random accesses, less balanced), but the property that they use fewer writes can be useful in the asymmetric setting (like samplesort vs. quicksort, treap vs. AVL or red-black tree).
4. In-cache data structures should draw more attention³. Since the data structures are kept in the cache (or small symmetric memory), the algorithm requires significantly less writes to the large asymmetric memory, although may require extra reads to compensate for less information we can keep within the data structure. In this paper, we discuss Dijkstra’s algorithm on shortest-paths as an example, and such idea can also be applied to computing minimum spanning tree, sorting, and many other problems.

2 Related Work

There exist a rich literature to show the read-write asymmetry on the new memories [ACM⁺11, ABCR12, BFG⁺15, BFG⁺16, CDG⁺16, CL09, DJX09, DWS⁺08, HZX⁺14, IBM14, KSDC14, LIMB09a, MDS⁺15, QGR11, XDJX11, YLK⁺07, ZZYZ09, ZWT13]. Regarding adapting softwares for such read-write asymmetry, some work has studied the system aspect. For example, there exist many papers on how to balance the writes across the chip to avoid uneven wear-out of locations in the context of NAND Flash chips [BAT06, EGMP14, GT05, PS09].

More closely-related papers targeting database operators: Chen et al. [CGN11] and Viglas [Vig12, Vig14] presented several interesting and write-efficient sequential algorithms for searching, hash joins and sorting. These are the early and inspirational attempts to design algorithms with fewer writes. Instead of formally proposing new computational models and analyzing the asymptotic cost, they mainly show the performance

³Similar ideas appear in many existing models already, like the external-memory model or the streaming model. However, the motivations are different: these models restrict the amount of space that can be used in the computation, while in our case the data structures are used to reduce the writes to the asymmetric main memory without including too many extra reads.

by the experiment results assuming external memories rather than main memories, or on the cycle-based simulators for existing architecture. For the latter case however, the prototypes of the new memories are still under development, and yet nobody actually knows the exact parameters of the new memories, or how they are incorporated into the actual architecture which is required for the setup of the cycle-based simulator. As far as we know, there is no available cycle-based simulator at the present time for the new memories. In the meantime, the asymmetries on latency, bandwidth, and energy consumption between reads and writes are different, and any of these constraints can be the bottleneck of an algorithm. Hence, designing algorithms on asymmetric memory are in a multiple-dimension parameter space, rather than just recording the running time from a simulator. Therefore, it is essential to develop theoretical models and tools that accounts for, and abstract this asymmetry and use them to analyze algorithms on future memory.

Blelloch et al. [BDBF⁺16, BFG⁺15, BGST16] formally defined several sequential and parallel computation models that take asymmetric read-write costs into account. Based on the computational models, many interesting algorithms (and lower bounds) are designed and analyzed in both sequential and parallel settings, which includes sorting, permuting, matrix multiplication, FFT, list/tree contraction, BFS/DFS and other graph algorithms, and many computational geometry and dynamic programming problems [BDBF⁺16, BDBF⁺18, BFG⁺15, BFG⁺16, BGSS18, JS17]. Carson et al. [CDG⁺16] also presented write-efficient sequential algorithms for a similar model, as well as write-efficient parallel algorithms (and lower bounds) on a distributed memory model with asymmetric read-write costs, focusing on linear algebra problems and direct N-body methods. Although many problems under the asymmetric setting have been studied, all the analyses are asymptotic and only show the upper and lower bounds on the complexity of these problems.

3 Our Model and Simulator

To start with, we discuss how to measure the performance of algorithms on asymmetric memories. We begin with the computational model that estimates the cost of an algorithm. This model requires the numbers of read and write transfers between the non-volatile memory and the cache, so later we introduce how the numbers of an algorithm can be simulated. Unlike the existing symmetric memories, a simple cache policy like LRU does not work on some asymmetric settings. Thus in Section 3.2 we briefly summarize the solutions to fix it, and then the cache simulator given in Section 3.3 captures this number with different cache policies.

3.1 The Cost Model for Asymmetric Memory

The most commonly-used cost measure of an algorithm is the time complexity based on the RAM model, which is the overall number of instructions and memory accesses executed in this algorithm. Nowadays, since the actual latency of an access to the main memory is at least two orders of magnitudes more expensive than a CPU instruction, the *I/O cost* based on the external-memory model [AV88] is widely used to analyze the cost of an *I/O*-bounded algorithm. This model assumes a *small-memory* (cache) of size $M \geq 1$, and a *large-memory* of unbounded size. Both memories are organized in blocks (cache-lines) of B words. The CPU can only access the small-memory (with no cost), and it takes unit cost to transfer a single block between the small-memory and the large-memory. This cost measure estimates the running time reasonably well for *I/O*-bounded algorithms, especially in multi-core parallelism. An efficient algorithm in practice should achieve optimality in both the time complexity and the *I/O* cost.

To account for more expensive writes on future memories, here we adopt the idea of an (M, ω) -Asymmetric RAM (ARAM) [BFG⁺16]: similar to the external-memory model, transferring a block from large-memory to small-memory takes unit cost; on the other direction, the cost is either 0 if this block

is clean and never modified, or $\omega \gg 1$ otherwise. The **asymmetric I/O cost** Q^4 of an algorithm is the overall costs for all memory transfers. Theoretical results on this new model have been studied in [BDBF⁺16, BDBF⁺18, BFG⁺15, BFG⁺16, BGSS18, JS17].

3.2 Cache Policies

Either the classic external-memory model or the new ARAM assumes that we can explicitly manipulate the cache in the algorithm. This largely simplifies the analysis, and in many cases is provably within a constant factor of a more realistic cache’s performance. For example, the standard least-recent used (LRU) policy is 2-competitive against the optimal offline cache-replacement sequence.

Interestingly, the competitive ratio does not hold in the asymmetric setting. Consider a cache with $k = M/B$ cache-lines and a memory access pattern that repeatedly writes to $k - 1$ cache-lines and read from other $k - 1$ cache-lines. An ideal cache policy will keep all $k - 1$ cache-lines associated to writes, so the I/O cost of each round is $k - 1$ for $k - 1$ read misses. An LRU policy however causes a cache miss for every single memory access, leading the I/O cost of each round to $\omega(k - 1) + k - 1$. This overhead is proportional to ω , which can be significant and problematic.

The solution is affected by the architecture, depending on whether software explicitly controls a DRAM buffer or not [CGN11, CNF⁺09, LIMB09b, QSR09]. If so, then the cost measures on these models are just the costs in practice, but programmers are responsible for managing what to put on the small-memory and guaranteeing correctness. The other option is to leave the hardware to control the small-memory. In this case, Blelloch et al. [BFG⁺15] show that if the small-memory is partitioned into two equal-size pools and each of them is maintained using LRU policy, the performance is 3-competitive against the optimal offline cache-replacement sequence (e.g. using $3\times$ space and incurring no more than $3\times$ cost).

We consider three different cache policies in this paper. The **Classic** policy maintains the small-memory as one memory pool and uses the LRU policy for replacement. The **SplitPool** policy keeps two separate memory pools and each runs the LRU policy. The **Static** policy allows static allocation in one memory pool, and the unallocated memory space is maintained using the LRU policy.

3.3 The Cache Simulator

The goal of this paper is to discuss new algorithmic approaches that minimize the asymmetric I/O cost to the main memory on a variety of fundamental data structures and algorithms. To capture the number of reads and writes to the main memory, we developed a software simulator that can adapt to different cache policies introduced in Section 3.2. The cache simulator is composed of an ordered map that keeps tracks of the time stamp of the last visit to each cache-line in the current cache, and an unordered map that stores the mapping from each cache-line to the corresponding location in the ordered map if this cache-line is currently in the cache. Interestingly, the implementation of this cache simulator is a natural application of the techniques discussed in this paper.

The cache simulator encapsulates a new structure `ARRAY` that is used in coding algorithms in this paper. It is like a regular array that can be dynamically allocated and freed, and supports two functions: `READ` and `WRITE` to a specific location in this array. The `ARRAYS` are responsible for reporting the memory accesses of the algorithm to the cache simulator, and the cache simulator will update the state of the cache accordingly. Therefore, coding using the `ARRAYS` is not different from regular programming much.

The memory accesses to loop variables and temporary variables are ignored, as well as the call stack. This is because the number of such variables is small in all of the algorithms in this paper (usually no more

⁴Throughout the paper, we abbreviate it as the *I/O cost*, unless stated otherwise explicitly.

than 10). Meanwhile, the call stack of all algorithms in this paper has size $O(\log n)$. The overall amount of uncaptured space is orders of magnitudes smaller than the amount of fast memory in our experiments.

The cache consists of one or two memory pools, depending on different cache policies discussed in Section 3.2. We will explicitly indicate the cache policy used in each of our experiments. The cache simulator maintains two counters in each memory pool: the number of **read transfers**, and the number of **write transfers**. When testing each algorithm on a specific input instance, the cache is emptied at the beginning and flushed at the end. A read or write is free if the location is already in the cache; otherwise the corresponding cache-line is loaded, the counter of read transfer increments by 1, and the least-recently-used cache-line in this pool is evicted. Also, a write will mark the dirty-bit of the cache-line to be `true`. When evicting a dirty cache-line, the counter of write transfer increments by 1. Notice that memory reads can cause write transfers, and memory writes can lead to read transfers.

When simulating the **Classic** policy (i.e., the standard one), we also verified our simulated results to ZSim (cycle-level simulator for current architecture), and the numbers always differ by no more than 10% when the parameters are set correctly.

4 Sets and Maps

Sets and maps are two of the most commonly-used data types in modern programming. Most programming languages either have them built in as basic types (e.g. python) or supply them as standard libraries (C++, C#, Java, Scala, Haskell, ML). In this section we discuss efficient implementations of unordered sets and maps implemented using hash tables, and due to the page limit, in ordered sets and maps implemented using balanced binary search trees are introduced in the full version of this paper.

4.1 Unordered Sets and Maps

Our implementation of unordered sets and maps is based on hash tables that support **lookup**, **insertion**, and **deletion**. The hash tables discussed in this section use open addressing and linear probing, since the goal of the data structure is to try to minimize the I/O cost focusing on smaller entries (accessing and reading larger entries are costly anyway so different hash-table implementations make minor differences). For simplicity, we assume no duplicate keys, and it is straightforward to handle the duplicates with minor modifications. In this setting, each operation of the hash table reads a small number of cache-lines, and an insertion or deletion will modify exactly one cache-line that contains the location of the key and will be eventually written back to the large-memory.

The challenge emerges when the set size changes dynamically. For an efficient implementation, we hope the overall size of the hash table to be neither too large nor too small. If the load factor passes 80%, linear probing's performance drastically degrades. On the other hand, we want the hash table size to be reasonably small to better utilize the small-memory (cache), since each cache-line holds more entries in this case. In practice, some implementations keep the load factor up- and lower-bounded by some constant. For example, a typical implementation keeps the occupancy of the hash table between 1/8 and 1/2, and the size doubles or shrinks by half if the number of entries exceeds this range. Such resizing reinserts p entries before at least $p/2$ insertions and deletions (where p is the set/map size). When reads and writes have approximately the same cost, the extra cost for such resizing is small compared to the query and update costs (e.g., the queries read from lots of memory locations). In the asymmetric setting however, the reads cost much less, but the extra writes in resizing can be significant: the resizing can incur at most twice ($p/(p/2) = 2$) the writes

Algorithm 1: The k -level hash table

Input: Parameter k , occupancy range l and r

```
1 function LOOKUP( $x$ )
2   for  $i \leftarrow 1$  to  $k$  do
3      $p \leftarrow HashTable_i.LOOKUP(x)$ 
4     if  $p \neq \text{null}$  then return ( $i, p$ )
5   return null

6 function INSERT( $x$ ) //  $x$  is not in  $HashTable$ 
7   for  $i \leftarrow 1$  to  $k$  do
8     if  $HashTable_i.occupancy < r$  then
9        $HashTable_i.INSERT(x)$ 
10    return
11  Allocate  $HashTable_{k+1}$  of size  $2 \cdot HashTable_k.size$ 
12  Relabel the hash tables with indices from 0 to  $k$ 
13  foreach  $y \in HashTable_0$  do
14     $INSERT(y)$ 
15  Free  $HashTable_0$ 

16 function DELETE( $x; i, p$ ) //  $x$  is located  $p$ -th in  $HashTable_i$ 
17    $HashTable_i.DELETE(x, p)$ 
18   if Overall occupancy is less than  $l$  (and  $HashTable_1.size > 1$ ) then
19     Allocate  $HashTable_0$  of size  $HashTable_1.size/2$ 
20     Relabel the hash tables with indices between 1 to  $k + 1$ 
21     foreach  $y \in HashTable_{k+1}$  do
22        $INSERT(y)$ 
23     Free  $HashTable_{k+1}$ 
```

compared to the initial insertions ($3\times$ writes in total). Hence, our goal is to discuss an alternative approach that optimizes such extra writes.

4.1.1 The k -level Hash Table

Instead of keeping one hash table, our main idea is to maintain a small number k of hash tables simultaneously, where k is a pre-determined parameter. In particular, the k -level hash table $HashTable$ is initialized with k arrays $HashTable_{1,\dots,k}$ with size $2^{c'+i}$ for $1 \leq i \leq k$ (or smaller in specific applications) and a constant c' . In practice we set c' to be 5.

For insertions, when the overall load factor exceeds some threshold r , we allocate a new chunk of memory with the double size of the largest current array, and the smallest hash table is discarded after all elements in it have been reinserted back. Similarly for deletions, if the occupancy of the hash tables drops below a threshold l , a small array with half of the size of the current smallest hash table is allocated, and the largest table is freed after the entries in it being reinserted. For instance, a valid k -level hash table may contain two arrays of size $2^{15} = 32768$ and $2^{16} = 65536$, when $k = 2$ and 30000 entries in the current configuration. The pseudocode of the k -level hash table is given in Algorithm 1. The occupancy range $0 < l < r < 1$ indicates when the resizing happens (a valid set of parameters can be $1/8$ and $1/2$). A classic implementation can be viewed as the special case of the k -level hash table when $k = 1$.

We now analyze the I/O cost Q of the k -level hash table. Here we assume that the size of the k -level hash table is larger than the small-memory and $1 - r < 1/B$, so that one single lookup, insertion or deletion in a single level in the hash table on average requires no more than $c < 2$ cache-line loads to find the location.

Lookup. In a k -level hash table, a lookup requires ck instead of c read transfers (c is the constant just defined) in the worst case (can quit earlier once the entry is found). The cost increases by a factor of k at most.

Insert. There are two definitions of insertions: an insertion that the key is known to be not in the set/map, or an insertion that it is unknown whether the key is in this set/map. Both cases are commonly-used. In this paper, we take the first definition and analyze the cost of this type of insertions. The second type of insertion can be viewed as a lookup first, then an insert if the lookup fails.

When inserting an element in a k -level hash table, we always try the smaller tables first. Once all tables are full, we resize it. More details can be found in Algorithm 1.

The I/O cost Q of an insertion comes in two parts: the cost of the initial insertion to the hash table, and the cost of this entry in future hash-table resizings. The cost of the initial insertion is no more than $c + \omega$, where c is the number of cache-line reads to find the position to insert, plus ω , one cache-line write for the actual insertion. The cost of resizing is more complicated to analyze.

We note that although a specific entry can be reinserted multiple times during different resizing processes, the overall number of element reinsertion is bounded, and thus we can amortize the work. A resizing occurs when an insertion comes in and the hash table contains exactly $r \cdot 2^p (2^k - 1)$ elements for some positive integer p . In this case, at most $r \cdot 2^p$ entries (the size of the smallest hash table), are reinserted during the resizing. The total number of insertions from the last resizing is at least $r \cdot 2^{p-1} (2^k - 1)$ (assuming $4l \leq r$), so the amortized I/O cost Q of reinsertion for each insertion is upper bounded by $\frac{(c + \omega)r \cdot 2^p}{r \cdot 2^{p-1} (2^k - 1)} = (c + \omega) \cdot 2 / (2^k - 1)$.

In the asymmetric setting when $\omega \gg 1$, the I/O cost of each insertion is approximately $\omega \cdot (1 + 2 / (2^k - 1))$, indicating that compared to the classic implementation where $k = 1$, in the worst-case the improvement when $k = 2, 3, 4$ is about 44%, 57% and 62% respectively. The asymptotic improvement when $k \rightarrow +\infty$ is 67%.

Delete. A deletion in the k -level hash table is similar to an insertion except that a lookup for the location is required (details in Algorithm 1). The cost of the initial deletion is $ck + \omega$. A resizing of the hash table can occur after at least $l \cdot 2^p (2^k - 1)$ deletions for some positive integer p , and the current hash table keeps $l \cdot 2^p (2^k - 1)$ entries. However, it is possible that all of these entries are in the last hash table so they are all reinserted. We note that when reinserting the elements from the discarded array, we always try smaller arrays first. This means that a reinserted entry, if not being deleted in the future, will not be reinserted again in the next $\min(k - 1, \log_2 r / 2l)$ shrinking resizings. Namely, the amortized extra cost of a deletion in future resizings is about ω / k if l is set to be about $2^{-k}r$. The overall I/O cost for a deletion is $Q = ck + \omega(1 + 1/k)$.

We have bounded of the I/O cost of each lookup, insertion or deletion, and the overall cost Q can be estimated by summing the amount of each operation multiplied by the cost of this operation. In practice, insertions and deletions can interleave. For example, when a deletion comes after an insertion, the number of entries remains the same, which leads to no further cost for these two updates afterward. The exact cost is also affected by the pattern of the sequence of the operations, and we will show by experiments.

4.1.2 Experiments

In the experiment we test the performance of our k -level hash table, including the numbers of read transfers and write transfers, I/O costs, and wall-clock running time, on different patterns of queries. In all experiments, we insert 1 million elements to an empty hash table, each of which is a 4-byte integer, into the hash table, and we vary the number of queries. The simulated cache contains 10,000 cache-lines and uses the *classic* policy, and for wall-clock running time we run the code on a PC with Intel i7-2600 CPU and 8GB RAM. The occupancy rate is set to be $l = 0.2$ and $r = 0.8$. We have tried other parameters (r between 0.6 and 0.8 and $l = r/4$). The results slightly vary, but all general conclusions in this section still hold.

Non-deletion cases.

Many applications, like webpage caching or the breadth-first searches, only insert but never delete elements in a hash table. Our experiment starts with this simpler case. We first show the relationship between k (the number of hash tables) and the numbers of read transfers and write transfers for a variety of insertion/query ratios, and the results are shown in Table 1. We fix the number of insertions to be one million, and query α times after each insertion. We vary α from 0, 1/8, to 8 ($\alpha < 1$ indicates one query per $1/\alpha$ insertions). About 50% query keys are in the hash table (this ratio affects the I/O cost since a successful query can terminate earlier). The number of levels k varies from 1 to 4. In Table 2, we show the overall I/O costs, which are the weighted sums assuming two typical values of the write-read ratio ω , 10 and 100.

We first look at the number of write transfers. When there is no query (i.e., the first column, just inserting 1 million entries), the numbers of writes are consistent with our analysis for insertions in Section 4.1.1. The only exception here is that cache can hold a constant fraction of the elements, which batches the writes and reduces the number of memory transfers. However, the relative trend in each column remains unchanged. Namely, the number of writes always decreases as the increase of k regardless of the ratio between queries and updates. The number of writes is reduced by 33%, 40% and 43% when $k = 2, 3, 4$ respectively. Such improvement also shows up in the overall I/O cost in Table 2.

We note that more queries cause more reads, and larger k also leads to more reads. Since these reads flush the cache-lines, the numbers of writes in these cases also marginally increase. The optimal choice of k is decided by the update/query distribution as well as the write-read ratio ω . In general, more queries lead to worse performance with larger k , and larger ω prefers larger k . In Table 2, we underline the numbers indicating the best choice of k in that specific setting. The experiment results indicate that picking k to be 2 or 3 is always a good choice when $\omega = 10$, and 3 or 4 when $\omega = 100$.

Wall-clock running time.

We also measure the actual running time of the previously stated operations on a real machine. In the current platform with symmetric memory-access costs, the write-read ratio is close to 1. Here we show that, the weighted sums of the I/O measures in Table 1 almost match the actual running times shown in Table 3, when plugging $\omega = 1$. Therefore, it is reasonable to believe that, the I/O cost shown in Table 2 can reasonably well project the performance for the future memory, when the bandwidths for reads and writes become asymmetric. (The argument for energy consumptions holds independently with this running time and other architectural issues.) Meanwhile, it is interesting to point out that, when insertions are more than queries, using k -level hash table with $k = 2$ is actually faster even in the current platform.

Our implementation has a much lower I/O cost compared to separate chaining. We run the same experiment using the STL unordered set (with the same hash function and other setups), and our hash table is at least 3-4 times faster in all cases.

Insertion and deletion cases.

We also tested the performance of k -level hash table on deletions. We first insert 1 million elements and then remove them all. After each insertion or deletion, we query α times. The other settings are the same as the non-deletion case. The results on the numbers of read transfer and write transfer are shown in Table 4, and the overall I/O cost with write-read ratio ω to be 10 and 100 are shown in Table 5.

From the results, we get almost the same trend as the non-deletion cases. Compared to the classic implementation (i.e., $k = 1$), the overall number of write transfer is reduced by 25%, 32%, and 36% when no queries are involved, and the improvement is slightly decreased when more queries come in, since more read accesses flush out the cache-lines. We note that the number of read transfers required by a deletion is more than that for an insertion, since in each deletion we need to locate the element in the hash table, which requires to look up in most k hash table levels. Hence, compared to the non-deletion cases, slightly

10^6 insertions, $\alpha \times 10^6$ queries where α is from 0 to 8, cache size is 10,000 cache-lines.

α	0		1/8		1/4		1/2		1		2		4		8	
	RT	WT	RT	WT	RT	WT	RT	WT	RT	WT	RT	WT	RT	WT	RT	WT
k=1	1.35	1.17	1.44	1.18	1.52	1.19	1.69	1.21	2.02	1.24	2.68	1.27	4.00	1.31	6.64	1.34
k=2	0.85	0.79	1.06	0.84	1.23	0.87	1.54	0.91	2.09	0.96	3.11	1.00	5.07	1.03	8.94	1.05
k=3	0.76	0.72	1.08	0.80	1.32	0.85	1.73	0.90	2.44	0.95	3.76	0.99	6.31	1.02	11.32	1.05
k=4	0.70	0.67	1.11	0.78	1.40	0.82	1.89	0.88	2.74	0.93	4.30	0.97	7.33	1.00	13.31	1.03

Table 1: Numbers of read and write transfers of k -level hash tables with different query/insert ratios. Numbers of read and write transfers are divided by 1M.

The I/O costs of k -level hash tables with the same settings in Table 1.

α	$\omega = 10$								$\omega = 100$							
	0	1/8	1/4	1/2	1	2	4	8	0	1/8	1/4	1/2	1	2	4	8
k=1	13.0	13.2	13.4	13.8	14.4	15.4	17.1	20.0	117.9	119.3	120.5	122.7	125.8	129.9	134.8	140.5
k=2	<u>8.8</u>	<u>9.5</u>	<u>10.0</u>	<u>10.7</u>	<u>11.7</u>	<u>13.1</u>	<u>15.4</u>	<u>19.5</u>	79.9	85.1	88.4	92.8	97.7	102.9	108.2	<u>114.4</u>
k=3	<u>8.0</u>	<u>9.1</u>	<u>9.8</u>	<u>10.7</u>	11.9	13.7	16.5	21.8	73.1	81.4	85.8	91.3	97.0	102.7	108.7	116.3
k=4	<u>7.4</u>	<u>8.9</u>	<u>9.6</u>	<u>10.7</u>	12.0	14.0	17.4	23.6	<u>67.9</u>	<u>78.6</u>	<u>83.8</u>	<u>89.6</u>	<u>95.5</u>	<u>101.3</u>	<u>107.7</u>	116.1

Table 2: The I/O costs of k -level hash tables with different query/insert ratios. The write-read ratio ω are selected to be typical projected values 10 (latency, bandwidth) and 100 (energy). Results are based on the numbers in Table 1. Numbers in red with underlines indicate the best choice of k that minimizes the I/O cost in this setting, and numbers in blue indicate better I/O costs comparing to the classic hash table implementation (i.e. $k = 1$).

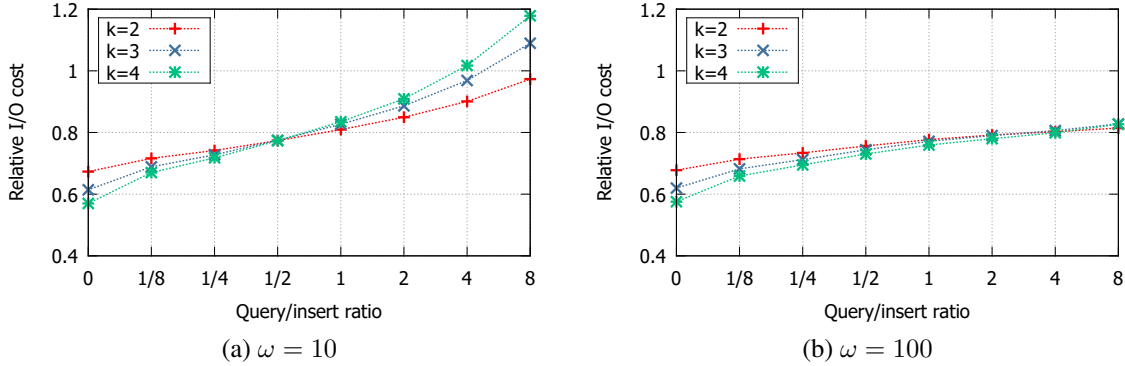


Figure 1: Relative I/O cost of k -level hash table with different k . The I/O cost is divided by the $k = 1$ case, so every data point below 1 indicates an improvement in such case. Numbers are from Table 2.

smaller values for k are more preferable. The best choice of k in each case is underlined and shown in Table 5. For smaller $\omega = 10$ (indicating bandwidth and latency), the best choice of k varies based on different query/update ratios, but $k = 2$ is always an acceptable choice. When considering the energy consumption ($\omega \approx 100$), a larger k , like 3 or 4, is more desirable in all cases.

Wall-clock running time (milliseconds), 10^6 insertions, $\alpha \times 10^6$ queries where α is from 0 to 8.

k	0	1/8	1/4	1/2	1	2	4	8
1	61	65	69	77	95	129	191	321
2	36	45	53	69	96	163	277	512
3	31	45	57	80	124	213	392	739
4	30	51	66	95	156	278	519	999

Table 3: Wall-clock running time (milliseconds) of k -level hash table with different combinations of k and query/insert ratios. Numbers in blue indicate a better performance comparing to the classic implementation (i.e. $k = 1$).

1 million insertions then 1 million deletions, α times 2 million queries where α is from 0 to 4, 10,000 cache-lines.

k	0		1/4		1		4	
	RT	WT	RT	WT	RT	WT	RT	WT
1	2.55	2.32	2.95	2.36	4.14	2.44	8.91	2.52
2	2.28	1.75	3.15	1.88	5.45	2.04	14.14	2.18
3	2.47	1.59	3.79	1.80	7.11	1.98	19.61	2.12
4	2.72	1.50	4.44	1.75	8.68	1.95	24.67	2.08

Table 4: Numbers of read and write transfers of k -level hash tables with different query/(insert+delete) ratios. Numbers of read and write transfers are divided by 1M.

The I/O costs of k -level hash tables with the same settings in Table 4.

k	$\omega = 10$				$\omega = 100$			
	0	1/4	1	4	0	1/4	1	4
1	25.8	26.6	28.5	34.2	235.0	239.4	247.6	261.4
2	19.7	22.0	25.9	35.9	176.8	191.6	209.6	232.0
3	18.4	21.8	26.9	40.8	161.4	183.6	205.1	231.4
4	17.7	22.0	28.2	45.5	152.3	179.6	203.5	233.1

Table 5: The I/O costs of k -level hash tables with different query vs. insert/delete ratios. The write-read ratio ω are 10 and 100. Results are based on the numbers in Table 4. Numbers in red with underlines indicate the best choice of k that minimizes the I/O cost, and numbers in blue indicate better I/O costs comparing to the classic hash table implementation (i.e. $k = 1$).

4.1.3 Conclusions

We proposed a new data structure, the k -level hash table, to implement unordered set and map, that has the same space utilization compared to the classic open-addressing hash tables. The key idea in the k -level hash table is to keep multiple instead of one level of hash tables. As a result, the algorithm uses fewer writes during resizings, at the cost of more reads in other operations.

The best choice of k is decided by the ratio of updates and queries. Our experiment shows that $k = 2$ always leads to a lower or similar I/O cost when the query/insert ratio is no more than 8, compared to the classic $k = 1$ setting. For the ratio of write/read cost is larger (like 100), larger values of k , like 3 or 4, are even more preferable than the $k = 2$ case.

4.2 Ordered Sets and Maps

The implementations of ordered sets and maps are based on some form of balanced tree (or tree-like) data structure and, at minimum, support **lookup**, **insertion**, and **deletion** in logarithmic time. Some set-set functions such as **union**, **intersection**, and **difference** are also required in many scenarios.

One commonly-used data structure to maintain ordered sets and maps is the self-balancing binary search tree (BST). Most languages and libraries use either AVL tree [AL63] or red-black tree [Bay72, GS78] to implement them, while some other implementations of weight-balanced trees [NR73] and treaps [AS89, SA96] also exist. Here we first analyze the I/O cost on some existing solutions.

4.2.1 I/O cost on BSTs

For simplicity, here we assume that the small-memory size is $M = O(1)$. Locating the key for a lookup, insertion or deletion requires to load and compare to $\Theta(\log n)$ tree nodes on the balanced binary search trees (BSTs), The I/O cost is $\Theta(\log n)$, which is also the lower bound of such operations.

For an insertion or deletion, we also need to modify the tree accordingly. In the asymmetric setting, weight-balanced trees [NR73] are not a good option since we have to update the subtree sizes all the way to the root. This update leads to $\Theta(\log n)$ writes to the large-memory per update. For the other types of BSTs, we show their I/O costs on insertions and deletions individually.

Red-black trees. Red-black trees [Bay72, GS78] have the simplest update rules among these balanced BSTs. With the classic rebalancing rules and careful implementation, it requires only $O(1)$ amortized time per update (insertion or deletion) after locating the key [Tar83]. Also, red-black trees require no extra cost to update balancing information except for the tree nodes involved in rotations (unlike the case in AVL trees). As a result, red-black trees have an optimal amortized I/O cost $Q = \Theta(\log n)$ per a lookup and $Q = \Theta(\omega + \log n)$ per insertion or deletion on the (M, ω) -ARAM.

AVL trees. An insertion in AVL trees requires at most two rotations (a double rotation) [AL63]. Unlike the red-black trees however, we need to track and update the balance factors along the path from the root to the modified tree node. We now bound the number of updated balance factors to be a constant. If we store the height of the subtree in each tree node, the difference of the two subtree heights can be checked in constant time. Since the number of subtrees of height d in an AVL tree is no more than $n/c^{\lfloor d/2 \rfloor}$ for a tree with n nodes and some constant $c > 1$ [BFS16], the number of increments of the counts for n nodes is $\sum_{d \geq 1} d \cdot (n/c^{\lfloor d/2 \rfloor}) = O(n)$. On average, an insertion needs $O(n)/n = O(1)$ writes.

The deletions of AVL trees is more complicated and $\Theta(\log n)$ rotations can be applied on every deletion. Amani et al. [ALT16] recently showed that there exists such a sequence of $3n$ intermixed insertions and deletions on an initially empty AVL tree that takes $\Theta(n \log n)$ rotations. This instance indicates that the classic implementation of AVL trees has an I/O cost $Q = \Theta(\omega \log n)$ per deletion in the worst case.

Treaps. A treap, also called as a randomized search tree [AS89, SA96], is a Cartesian tree in which each key is given a randomly chosen numeric priority, and the inorder traversal order of the nodes is the same as the sorted order of the keys. The priority for any non-leaf node must be greater than or equal to the priority of its children.

When inserting an element into a treap with $n - 1$ elements or removing an element from n elements, The updated element only compares to the elements that each has a higher priority than the other elements between

this element and the updated element. The number of such comparisons is $\sum_{j \in [n] \setminus \{i\}} 1/|i - j| = O(\log n)$ in expectation⁵, where i is the position of the updated element in the total order of n elements [SA96].

The number of rotations can be computed similarly. For an insertion, a rotation happens once the inserted element has a higher priority than all elements in the entire subtree. Again if we assume the inserted element ranked i -th in the total order, the probability that it rotates up for the j -th ranked tree node is $1/|i - j|^2$ (i.e., the j -th element has higher priority than all elements between, and lower than the priority of the inserted node). The overall expected number of rotations per insertion is $\sum_{j \in [n] \setminus \{i\}} 1/|i - j|^2 = O(1)$ in expectation. We can show the constant writes per deletion accordingly.

We note that unlike an AVL tree or a red-black tree, a treap does not require updates to the balancing criteria, which means that we never need to modify the information in each tree node after it is inserted. As a result, an insertion, deletion or query on a treap requires $O(\log n)$ reads *whp* and an insertion or deletion requires $O(1)$ writes in expectation.

4.2.2 Join-based implementation

Blelloch et al. [BFS16, SFB18] recently proposed a framework that supports fast bulk operations. This framework supports efficient single or multiple insertions or deletions on AVL trees, red-black trees, treaps and weight-balanced trees, as well as union, intersection and set difference on two BSTs (of the same kind). This framework is very concise: each function can be implemented within a dozen lines of code and are independent with the specific balancing criteria in different types of BSTs. The functions rely on only one helper operation JOIN [Ada92, Ada93, ST85, Tar83, BFS16], which deals with the tree balancing and can be treated as a black box. Under this framework, the implementation of these functions is highly efficient and parallelized.

JOIN(T_L, k, T_R): The JOIN function takes two balanced BSTs (of the same balancing schemes) T_L and T_R and a key k , and returns a new balanced BST for which the in-order values are a concatenation of the in-order values of T_L , then k , and then the in-order values of T_R . In order to keep the ordering in the tree nodes, the JOIN function assumes k to be greater than all keys in T_L and smaller than all keys in T_R . JOIN handles all issues related to rebalancing, and is the only function that knows about and maintains the balance invariants. The JOIN algorithms for each of the balancing schemes are given in [BFS16]. Meanwhile, JOIN is the only function that *writes* to the memory. More specifically, all operations that change the attributes of a tree node (e.g., linking to new children, height-maintaining for AVL or red-black trees, color-changing in red-black trees, etc.) are restricted in JOIN. This property also greatly simplifies the counting of writes in our simulator and makes the optimizations to reduce writes easier.

SPLIT(T, k) $\rightarrow (T_L, T_R)$: SPLIT can be viewed as the inverse of JOIN that takes a balanced BST and a key k , and splits the tree into two smaller trees T_L that contains keys smaller than key k , and T_R for keys larger than key k . SPLIT can be trivially implemented using JOIN by a recursive approach.

Bulk Operations. Using JOIN as a black box, either single or bulk updates can be implemented simply and generally across different balancing schemes. As an example, we give the algorithm of taking the union of two BSTs (or sets/maps) in Algorithm 2. The algorithm is based on divide-and-conquer: T_1 's root is used to partition all elements in two trees into two disjoint set, one with T_1 's left subtree and T_l , and the other with T_1 's right subtree and T_r . Then we union the two subproblems recursively and independently, and join the two output trees using T_1 's root. The correctness and running time are shown in [BFS16].

⁵Also with high probability.

Algorithm 2: The union function

```
1 function UNION( $T_1, T_2$ )
2   if  $T_1 = \text{NIL}$  then return  $T_2$ 
3   if  $T_2 = \text{NIL}$  then return  $T_1$ 
4    $\langle T_l, T_r \rangle \leftarrow \text{SPLIT}(T_2, T_1 \text{'s root})$ 
5   return JOIN( $T_1$ 's root, UNION( $T_1$ 's left subtree,  $T_l$ ), UNION( $T_1$ 's right subtree,  $T_r$ ))
```

There are several benefits of using this framework for our implementation and experiment. First, as we just explained, different updates have the uniform code on different types of BSTs (except for the JOIN), which justifies the performance by the different balancing criteria for the BSTs, instead of the different implementations for different trees. Second, although the joined-based implementation operates on two trees, one can check that when one tree contains only a singleton element, the algorithm runs the same as the algorithm of a single insertion on each type of the BSTs. The deletion can also be implemented by taking the difference by the original tree and a tree with a single element. As a result, the joined-based implementation is strictly more powerful. Lastly, we can also run interesting experiments on more operations like bulk updates, and compare the results on different BSTs.

4.2.3 Experiments

In this section, we show our experimental results on the counts of read/write transfers for different settings. Due to the page limit, our experiment mainly focuses on the performance of various binary search trees (AVL trees, red-black trees, and treaps) with different batch sizes.

In the experiment, we first insert $m = 10^6$ million integers as keys to a tree T (empty at the beginning), drawn from a uniform distribution from 32-bit unsigned integers, and then delete them in a uniformly random order. The insertion and deletion are grouped in batches of size s , indicating that the insertions are $\lceil m/s \rceil$ unions on the main tree T with trees of size s . The deletions are also batched in $\lceil m/s \rceil$ bulks of size s . In our experiment, we construct a smaller tree for each batch and then call the union or difference function we just mentioned. We note that if all update elements are given in advance, we can also sort them and put them in a list. Since this is a more specific case, our experiment is based on the tree-tree updates.

The node size in all different types of trees is 16 bytes. Each tree node stores four 4-byte data blocks to hold the key, the left and right pointers, and the balancing information. The cache contains 10,000 cache-lines, similar to the setting for unordered sets.

Table 6 shows the experimental results on BSTs with different balancing schemes with various batch sizes. The numbers are read and write transfers per update.

Batch size 1.

We first look at the case corresponding the single insertions and deletions, and the results are shown in Figure 2 and 3. Regarding the number of writes, treaps show the best performance. This is easy to understand since treaps do not modify any information for rebalancing during insertions/deletions. The structure of treaps is deterministic once the priorities are decided. The priorities are set before the merging (deleting), and never changed later. For such reasons, treaps require much fewer writes per update compared to AVL and red-black trees.

The AVL and red-black trees maintain the balancing information on each tree node that needs to be updated when the subtree height changes. Then more writes are used due to these updates. Between these two types of BSTs, red-black trees require more writes, since there are also some color flips on the siblings of

(a) Insertion / Union

Batch Size	AVL		RB-Tree		Treap		$\omega = 10$			$\omega = 100$		
	RT	WT	RT	WT	RT	WT	AVL	RB-T	Treap	AVL	RB-T	Treap
1	11.28	2.83	12.00	3.48	16.61	1.79	39.6	46.8	34.5	295	360	196
1k	12.67	2.89	13.29	3.66	17.18	1.89	41.6	49.9	36.1	302	379	206
10k	6.18	2.68	6.40	3.01	7.33	1.85	33.0	36.5	25.8	274	308	192
100k	2.54	1.84	2.54	1.86	2.56	1.66	20.9	21.2	19.2	186	189	169

(b) Deletion / Difference

Batch Size	AVL		RB-Tree		Treap		$\omega = 10$			$\omega = 100$		
	RT	WT	RT	WT	RT	WT	AVL	RB-T	Treap	AVL	RB-T	Treap
1	13.83	2.72	16.17	5.17	17.85	1.98	41.1	67.8	37.7	286	533	216
1k	15.11	2.79	17.65	5.21	18.52	2.08	43.0	69.8	39.3	294	539	227
10k	8.17	2.69	10.43	3.44	9.09	2.06	35.1	44.9	29.7	278	355	215
100k	3.22	1.99	3.54	2.43	3.23	1.78	23.2	27.8	21.1	203	246	182

(c) Average Tree Depth

AVL	RB-Tree	Treap
19.39	19.62	26.48

Table 6: Numbers of read and write transfers and asymmetric I/O costs of different BSTs with various batch sizes. The numbers are divided by 10^6 (i.e., per inserted/deleted elements). The write-read ratio ω are selected to be typical projected values 10 (latency, bandwidth) and 100 (energy).

the updated tree path. Such flips are extremely cheap in the classic symmetric setting but will cost much in the asymmetric setting when writes become expensive.

The fewer writes for treaps come together with the extra cost on more reads. Treaps are less strictly balanced compared to the other two trees, which saves the writes to maintaining such balancing, but leads to larger average tree depth (shown in Table 6(c), about 30% deeper than the other trees). The average depths for AVL and red-black trees are close to optimal, which is 18.96 for a perfectly balanced tree with 10^6 nodes. Because of the shallower average depth, the number of reads required in either updates or queries is much small on these two trees.

Larger batch sizes.

We now discuss how does the batch size affect the numbers of read and write transfers. The numbers are given in Table 6.

When the batch size increases but remains small, the reads and writes almost remain the same and slightly increase. When the batch size is smaller, the elements in a batch and the paths to visit them are always loaded into the cache once and always stay there. Therefore, the different batch sizes less than 100 do not affect the performance much.

The peak I/O cost per update is around batch size 1000, as we show here. In such case, the overall footprint of a batch update no longer fit into the cache, which may lead to two loads per tree node (once in the split phase and the other in the join phase).

However, the cost turns down when the batch size grows over 1000. The reason is that, the number of tree nodes visited during each bulk update is $\Theta(s \log(m/s))$ (recall that s is the bulk size), which is also the time complexity [BFS16] of this process. Namely, we need to visit $O(\log(m/s))$ tree nodes per inserted node, so

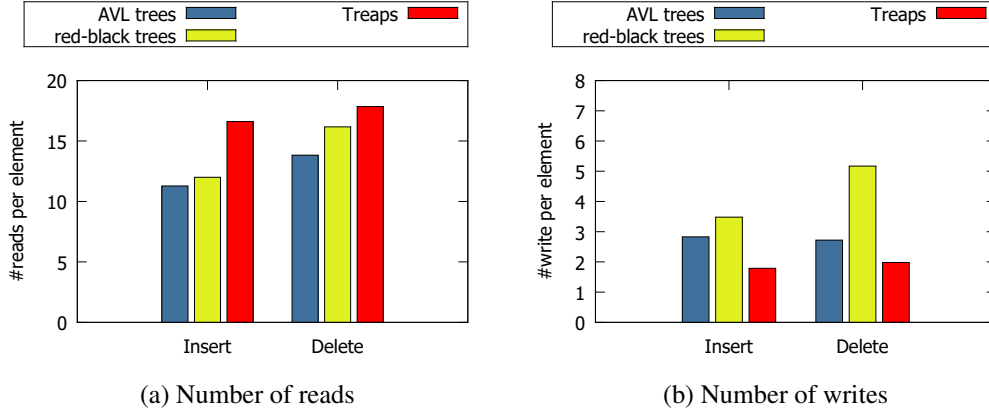


Figure 2: The number of read and write transfers of different BSTs on single insertion/deletion. Data are from the first rows in Table 6(a) and 6(b).

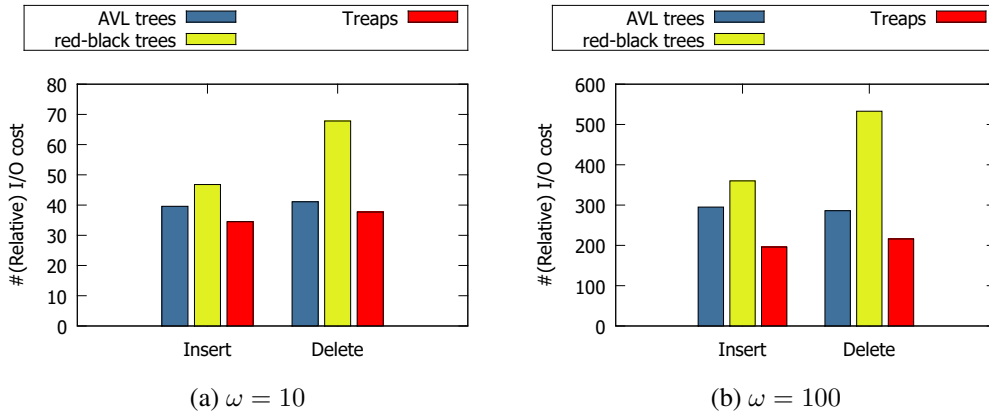


Figure 3: The I/O cost of different BSTs on single insertion/deletion. Data are from the first rows in Table 6(a) and 6(b).

we touch fewer nodes as s goes up. Compared to the single updates, each node on multiple tree paths⁶ is only looked at and compared with once in the joined-based bulk updates, and this is from which the improvement comes. Since the top levels in the trees are visited more frequently, they usually stay in the cache. As a result, the saved memory accesses for small batch size are hidden by the function of the cache. However, when the batch size keeps growing and exceed the cache size, then the saved memory accesses lead to lower reads, as shown in Table 6.

The number of writes also decreases for larger batch sizes in AVL and red-black trees, because of the previously stated reason. When the elements are inserted or deleted one by one, the balancing factor of a node on by multiple tree paths can be updated multiple times. On the other hand, when the tree is updated in a bulk, such information will be updated by at most once at the join point, which saves the number of writes to the asymmetric memory. However, since treaps do not need to update such information, we cannot observe a significant drop-off on writes for larger batch sizes.

⁶Usually on the top part of the tree. For example, every single insertion visits the root node.

Queries.

We have not explicitly tested the I/O costs for queries, since most queries (like finding or checking a key, locating the k -th element) have the same memory access pattern as the insertions. The only difference is that they do not modify the tree, so there will be no writes. These updates may flush the dirty cache lines, and thus slightly increase the writes. Our experiment shows that if we have the same number of queries and insertions, the number of writes increases by no more than 5%, which is insignificant. Therefore, we believe that we can ignore such changes in the most cases.

4.2.4 Conclusions

In this section, we theoretically analyze the asymmetric I/O costs of different types of binary search trees. We show that red-black trees, the insertions for AVL trees, and treaps on expectation have an optimal asymptotic cost ($\Theta(\omega + \log n)$ per update).

We then test the actual performance by conducting experiments based on the join-based implementation, and show that treaps have the best update cost in most cases. The advantage comes from a looser balancing constraint, which also leads to a larger tree depth and query costs. As a result, AVL tree will be a better option if the queries are much more than the updates.

5 Sorting

Sorting is one of the most fundamental algorithms and building blocks in algorithm design and programming. In this section, we analyze, performance engineer and experiment the performance of the existing sorting algorithms in the asymmetric setting, which include quicksort, mergesort, BST sort and samplesort.

5.1 Algorithms and Implementations

We discuss four sorting algorithms, which are either used as the state-of-the-art implementations or have efficient theoretical guarantees on asymmetric I/O cost. All four algorithms requires $O(n \log n)$ (optimal) comparisons.

Quicksort. Quicksort is one of the commonly-used sorting algorithms in both sequential and parallel setting. It picks a random pivot (or the median among several random samples) that partitions the array into three contiguous subsets containing the input values that are less than, equal to, and larger than the pivot respectively, based on a scan-based approach. Then two recursive calls are made to the less-than and larger-than subsets, until the base case (we set it to be 16 elements) is reach and the algorithm switch to an insertion sort. In our implementation, the pivot is selected to be the median of three random positions.

We now analyze the I/O cost of quicksort. Once the subproblem size is no more than M , the small-memory size, the subproblem will fit into the cache and no more read and write transfers are required for this subtask thereafter. The I/O cost Q is thus $O(\omega n/B \cdot \log(n/M))$, where the number of memory transfers to partition the array in the same recursive level add up to $O(n/B)$ and w.h.p. the algorithm requires $O(\log(n/M))$ levels to reach the subtask with size no more than M .

Mergesort. Mergesort is another textbook sorting algorithm that is stable and easy to parallel. The implement of mergesort has different versions, and here we discuss an I/O-efficient version. The algorithm partitions the array into two equal-length parts and recursively sort every subproblem individually. After the computation of both subtasks is finished, a scan-based process merges the two sorted subsequences into the final output.

To implement the algorithm efficiently, we use the rotating arrays so that every element is moved only once in one merge process. An extra round of data copy is applied if the latest merge leaves the sorted result in the temporal array.

Similarly to quicksort, no further memory transfers are used in subproblems with size no more than M . Therefore the I/O cost Q of mergesort is also $O(\omega n/B \cdot \log(n/M))$. More specifically, $Q = (1 + \omega)n/B \cdot \lceil \log_2(n/M) \rceil$ when $\lceil \log_2(n/M) \rceil$ is even, and otherwise there is an extra $(1 + \omega)n/B$.

BST sort. BST sort treat the input as an ordered set and all elements are maintained in a binary search tree. Then the in-order traversal of the tree is the sorted output. The BST can either be balanced, or does not apply any rotation but the elements are inserted in a random order. Both cases give I/O cost of $O(n \log n + n\omega)$ on n input elements. It is typically used when the input is adaptive: we can insert and delete elements dynamically, while the sorted result are maintained at any time. In this case, some balancing schemes are required if the insertions and deletions are biased.

In our implementation elements are inserted in a random order, which is done in two phases: the first phase generates n uniformly random integers a_1, \dots, a_n between 1 to n using some hash functions, and inserting the a_i -th elements into the tree with *no* rotations (check whether this element is already inserted before the actual insertion); then in the second phase we insert every uninserted element into the tree in the input order. This guarantees that w.h.p. the gap between the elements inserted in the first phase is $O(\log n)$ in the output order. Hence the tree depth is $O(\log n)$ disregarding any input order, and thus inserting one element requires $O(\log n)$ reads and two random writes (the boolean flag that this node is inserted, and to change the parent's pointer).

Samplesort. Samplesort is a class of sorting algorithms that are based on divide-and-conquer paradigm with multiple pivots and especially commonly-used in the multicore setting. To minimize the I/O cost and number of memory transfers, in this paper, we discuss a cache-aware implementation of samplesort.

Given the cache size M and block size B , if the input size is smaller than M then we simply call quicksort. Otherwise, we pick M random samples from the input, sort the samples, and pick the B -th, $2B$ -th, \dots , $(M - B)$ -th samples as the **pivots**⁷ and they form M/B buckets. Then for each input element, we binary search (an implicit search tree) its associated bucket. Finally, we again samplesort the elements in each bucket individually and recursively.

The algorithm partition the input into M/B almost equal-size subproblems with $O(n/B)$ read and write transfers. The I/O cost of this algorithm on average is $O\left(\frac{\omega n}{B} \log_{\frac{M}{B}} \frac{n}{B}\right)$.

The number of write transfer of samplesort in our implementation is highly optimized so that one round of samplesort only requires three sequential reads and one sequential write per input element. After picking the pivots, we first loop over the input to binary search the associate bucket of each element. However, we do not store this value; instead, we only modify the counters of each bucket. After that, we have known the number of elements in each bucket, and we apply a prefix sum compute the offset of each element. Finally, all elements are distributed to their associated buckets based on the offsets. The algorithm only requires three reads and one write for each element, and all other operations are all within the cache. Notice that after one round, the data are stored in another array, so the final results will be moved back to the original array if they happen to be in the other one.

⁷If $M < n/B$ then we pick n/B samples, and this will happen only once w.h.p. in the algorithm. In practice we replace M to be M/c' for some small constant $c' > 1$ to ensure everything that should be in the cache is actually in the cache.

Algorithm	I/O cost
Quicksort	$O(\omega n/B \cdot \log(n/M))^*$
Mergesort	$O(\omega n/B \cdot \log(n/M))$
BST sort	$O(n \log n + n\omega)$
Samplesort	$O(\omega n/B \cdot \log_{M/B}(n/B))^*$

Table 7: List of I/O costs on sorting algorithms. (*) indicates the bounds hold with high probability ($1 - n^{-c}$ for arbitrary $c > 0$).

There do exist other sorting algorithms like heapsort, shellsort, bitonic sort, etc. Their performances regarding I/O cost however, are less competitive against the previous sorting algorithms due to either more work or inefficient memory-access pattern. We did not compare with previous work in [Vig14, BFG⁺15] in this paper, since they are not optimal in terms of comparisons and have tunable parameters in the algorithm, that makes the comparison among them inconclusive. Instead as a very first paper on this topic, we focus on simple algorithms and draw interesting conclusions that can also be useful in implementing these more complicated approaches in the future.

5.2 Experiments

In the experiment we test our implementations of the four sorting algorithms on the numbers of read and write transfers with various cache sizes. We set the input into two categories: one that we indeed move the data, and the other that the algorithm sorts the indirect pointers pointing to the data fields. The first case is when we are sorting integers, real numbers or any structures with small and fixed size like graph edges, key-value pairs of integers and/or pointers, etc. The second case is when the input is irregular, like strings, texts, images, or any structure that is either large or the sizes vary among different elements. In this case, moving the data is expensive, so we will sort and output a list of indices pointing to the actual data fields. Both cases are widely used in practice.

Performance on sorting float-point numbers. We now show the experiments that perform data movement. We sort 10 million random i.i.d. double-precision float-point numbers and in this case the block size $B = 8$. The numbers of comparisons and read and write transfers of four algorithms are summarized in Table 8. The I/O costs of the algorithms for two typical values of ω (10 and 100) are visualized in Figure 5.

Both quicksort and mergesort require about $\log_2(n/M)$ rounds, which is 8 to 16, to fully fit the sub-problems into the small-memory. Both algorithms require approximately the same number of writes, since quicksort is in-place while mergesort is not and has double memory footprint, but the partition of quicksort is not exactly even. These factors offset each other. However, the read transfer are doubled in mergesort, and in practice causing its less efficient practical performance. Mergesort requires fewer comparisons, but this does not lead to a significant difference in modern architecture, compared to the I/O cost to the main memory.

BST sort, as we analyzed theoretically, requires about $\log_2(nB/M)$ reads and two random writes plus some small cost on initializations. BST sort uses less writes compared to quicksort or mergesort when $2B < \log_2(n/M)$, which is hard to be satisfied given the parameters of the hardware. BST sort thus will still be less efficient to sort small and fixed-structure entries in future memories.

The experiment result of samplesort follows our analysis as well: each round of sample sort approximately requires three sequential reads and one sequential write, and the algorithm uses roughly $\lceil \log_{M/B}(n/B) \rceil$ rounds to reach the base case. Other costs (sampling and sorting pivots, transposing the offsets, etc.) are negligible.

Cache Size	Quicksort			Mergesort			BST sort			Samplesort		
	RT	WT	#comp	RT	WT	#comp	RT	WT	#comp	RT	WT	#comp
100	2.06	2.04	39.49	4.01	2.00	24.11	24.49	2.04	43.76	1.60	0.82	63.05
1000	1.58	1.57	39.49	3.11	1.56	24.11	19.39	2.04	43.76	1.01	0.51	56.63
10000	1.11	1.10	39.49	2.25	1.13	24.11	14.17	2.03	43.76	0.50	0.25	52.47

Table 8: Number of read and write transfers of different sorting algorithms on 10^7 random i.i.d. double-precision float-point numbers. Values in the table are numbers of read and write transfers and comparisons per input element. Cache size is measured by the number of 64-byte cache-lines.

Cache Size	Quicksort			Mergesort			BST sort			Samplesort		
	RT	WT	#comp	RT	WT	#comp	RT	WT	#comp	RT	WT	#comp
100	19.83	1.14	37.68	17.04	1.02	23.30	46.24	1.93	40.93	6.19	0.79	35.97
1000	15.57	0.91	37.68	13.42	0.81	23.30	35.30	1.92	40.93	3.93	0.45	34.34
10000	11.12	0.67	37.68	8.96	0.56	23.30	24.65	1.90	40.93	2.38	0.38	34.43

Table 9: Number of read and write transfers of different sorting algorithms on 2×10^6 indices pointing to structures with average size of 64 bytes. Other setup is the same as that in Table 8.

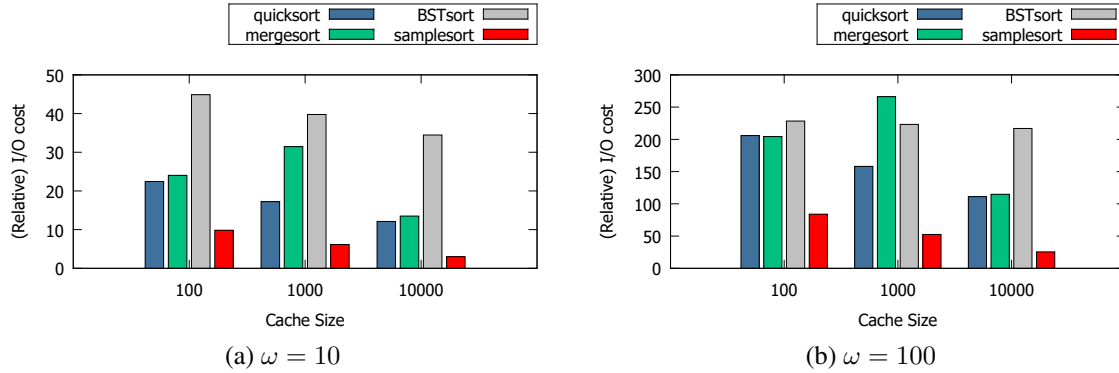


Figure 4: I/O costs of different sorting algorithms on 10^7 **doubles** with various read-write asymmetry ω . Numbers are from Table 8.

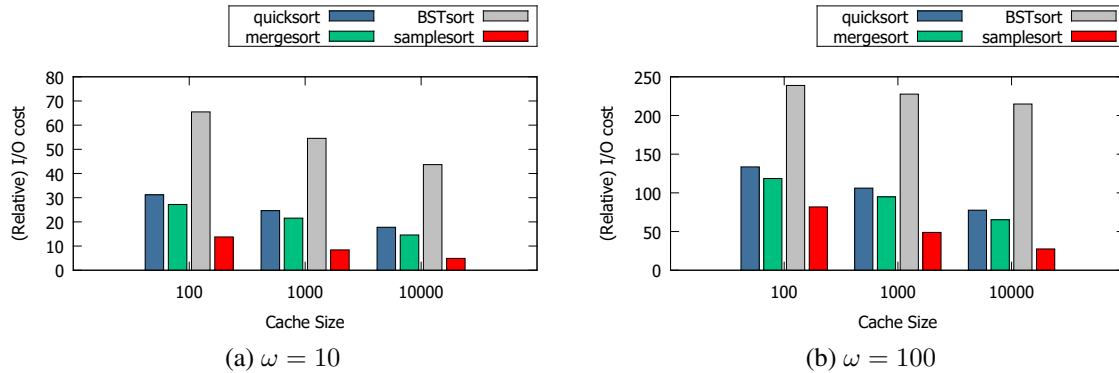


Figure 5: I/O costs of different sorting algorithms on 2×10^6 **pointers** with various read-write asymmetry ω . Numbers are from Table 9.

Based on our previous experience, after highly optimizing the performance in the sequential setting, samplesort is slightly slower than quicksort (but much faster in multicore parallelism) even though the I/O cost of samplesort is lower. We believe the reason to be that, samplesort does not explicitly utilize L1 and L2 caches (the binary search, offset transpose, etc.) so most small-memory accesses are toward L3 cache. The memory access pattern for quicksort, however, is mostly sequential scanning, which is highly optimized (like by prefetching) in the current architecture. Since the write costs will significantly increase for future non-volatile main memory, we project that samplesort will be more efficient in this setting, since the I/O cost will very likely be the bottleneck of all these sorting algorithms (similar to the existing multicore parallel setting).

Sorting larger entries. As shown previously, although BST sort requires only at most random two writes per insertion, this cost is still significant compared to other sorting algorithms that have better spatial locality in memory access. However, when the size of each input entry becomes larger, each cache-line holds fewer entries and the number of write transfers per entry increases.

We run the experiment with a fixed number of inputs (10^7) and cache-lines (1000), and varies the size of each input entry from 8 bytes to 64 bytes. The results are shown in Table 10. Based on the analysis in Section 5.1, for quicksort, mergesort, and samplesort, the I/O cost is proportional to the entry size. However, BST sort is mostly not affected. The number of reads just marginally increases since the cache can hold less top-tree levels, and slightly more writes are required since the overall footprint of the search tree increases as the growth of the entry size. From the results, we can see that when the entry size becomes at least 64 bytes, BST sort outperforms all other algorithms in the asymmetric setting.

Performance on sorting via pointers. In this experiment, we sort 2 million random strings that the characters are stored contiguously in the memory with an average size of 64 bytes. The input also contains 2 million 8-byte pointers to the strings, and we only sort the pointers.

The implementations of quicksort, mergesort, and BST sort are identical except that each comparison requires two indirect addresses to locate the data field. Samplesort, since it is cache-aware, requires two modifications: first, no oversampling for the pivots since now each pivot requires one or two cache-lines for its data; second, instead binary-search the bucket of each element twice, we now store the bucket label after the first search to avoid the second search (roughly double the writes while half the reads).

The numbers of read and write transfers are reported in Table 9 and visualized for two typical values of ω (10 and 100) in Figure 5. Writes of quicksort, mergesort, and BST sort stay approximately the same as the previous experiment (the decrease in values is due to fewer input elements). Reads are now increased by about $\log_2(nB/M)$, the cost of locating the data fields.

For each round in sample sort, each element now requires one random read, three sequential reads and two sequential writes (i.e. 1.375 read transfers and 0.25 write transfer). Also, one extra random read per element is in the base case to load the data into the cache. The algorithm requires $O(\log_{M/B}(n))$ rounds to reach the base case, and in the experiment, this round number is approximately 3, 2 and 1 for three cache sizes. Finally, output needs to be moved back when reaching the base case in an odd number of rounds.

As what we understand, the wall-clock performance of samplesort when sorting pointers on current architecture is already faster than the other sorts. This is because samplesort requires fewer I/Os and comparisons, and all algorithms randomly access the memory so techniques like prefetching cannot help. This gap will be even enlarged in future NVMs when the write costs being exaggerated.

Entry Size	Quicksort		Mergesort		BSTsort		Samplesort	
	RT	WT	RT	WT	RT	WT	RT	WT
8	1.58	1.57	3.25	1.63	19.39	2.04	1.01	0.51
16	3.26	3.17	7.00	3.50	21.11	2.29	2.09	1.08
32	6.58	6.13	15.00	7.50	23.04	2.79	4.25	2.24
64	13.18	11.58	32.00	16.00	25.35	3.79	8.91	4.87

Table 10: Numbers of read and write transfers of different sorting algorithms with various data sizes (bytes). The input size is 10^7 and the cache contains 1000 cache-lines.

5.3 Conclusions

Based on our implementations and experiments, the following conclusions and the projection of the techniques for future non-volatile main-memories can be drawn.

Samplesort. Samplesort generally requires fewer I/Os than other sorting algorithms. On existing hardware, since samplesort does not explicitly optimize for L1/L2 cache, its sequential performance is slightly slower than quicksort. However, in the multicore parallel setting, samplesort is always the fastest due to its efficiency on I/O cost. We predict that samplesort will play a more significant role with the future hardware even in the sequential setting, because of the lower number of accessing to the large asymmetric memory and the asymmetry of bandwidth and energy on the new memories.

Sorting indirect pointers. The advantage of samplesort will be enlarged when sorting via pointers. This is because samplesort requires fewer rounds to finish, and therefore it uses fewer reads, fewer writes, and fewer comparisons compared to other algorithms. This also matches the observation on the existing symmetric setting that samplesort is more efficient on wall-clock performance.

BST sort. Although BST sort only requiring constant writes on large-memory per update, it is still inefficient compared to other approaches when sorting integers or float-point numbers because of its lacking on utilizing the cache-lines. However, when sorting entries with at least 64 bytes (i.e. a cache-line size), the I/O cost of BST sort outperforms the rest three algorithms.

6 Graph Traversal Algorithms

We now provide the full version of the graph-traversal algorithms in this paper. We discuss two of the most commonly-used graph-traversing algorithms: Breadth-first search (BFS), and Dijkstra’s Algorithm. We show that with appropriate implementations, these algorithms can write much fewer to the main memory, comparing to the classic implementations.

Throughout this section we assume the input graph $G = (V, E)$ contains $n = |V|$ vertices and $m = |E|$ edges.

6.1 Breadth-First Search

Breadth-first search (BFS) is an algorithm for graph traversing or searching. It starts at the source node, which is an arbitrary node of a graph, and explores the vertices with respect to the distances (the number of hops) from the source node. BFS is commonly-used in computing single-source shortest paths on unweighted graphs, as a subroutine for graph radii estimation, eccentricity estimation and betweenness centrality, and as

a basic building block for other graph algorithms like graph connectivity, reachability, bridges, biconnected components, and strongly connected components. In this paper we focus on the most basic application: shortest paths on undirected graphs, and the techniques discussed here can apply to many other applications.

6.1.1 The classic implementation

Given a graph $G = (V, E)$ with n vertices and m edges, the classic implementation of BFS keeps a queue with size n , and an array of boolean flags with size n indicating that each vertex is visited or not during the search. This implementation requires at most 2 writes for each vertex: one sequential write for adding it to the queue and one random access for changing the flag of this vertex. Meanwhile, searches along an edge to an already visited vertex require no writes. The overall I/O cost of BFS $Q(n, m) = O(\omega n + m)$ [BFG⁺16].

This bound is asymptotic optimal for arbitrary graphs, since the output of BFS, the shortest-path tree, has size $O(n)$. However, a number applications (e.g. s-t shortest-path or connectivity, graph radii estimation or eccentricity estimation) have output size $O(1)$, which allow techniques utilizing the small-memory and reducing the number of writes.

6.1.2 BFS implementation using rotating arrays

The rotating arrays are used in many algorithms to reduce the space requirement. For example, the diamond DAG computation in the longest common subsequence (LCS) problem only requires storing two consecutive columns (or rows) at any time during the dynamic programming process, since each DP value only depends on three other nearby vertices in the diamond DAG. The algorithm maintains two arrays each with the size of either input string: one to hold the results in odd columns (rows) while the other for the even ones. We call this structure the *rotating arrays* since the two arrays are rotating and holding computed and uncomputed values. They only require temporal (cache) space of one dimension although the nodes in the entire DAG have two dimensions.

Here we observe that BFS on undirected graphs can apply a similar approach: instead of keeping a queue and a global array of boolean flags with size n , we maintain separate queues, each corresponding to a specific frontier (i.e. the set of vertices with the same distance to the source). This is because, during the BFS on undirected graphs, the outgoing edges from the i -th frontier can only reach the vertices in either the $(i - 1)$ -th frontier or the $(i + 1)$ -th frontier. Otherwise, assume an edge reaches a vertex in the $j (< i - 1)$ -th frontier, then the vertex in j -th frontier will visit this vertex in $(j + 1)$ -th frontier. As a result, when processing the i -th frontier, we only need to keep three frontiers with distance $i - 1$, i and $i + 1$.

Since each frontier is separately kept and all vertices in one frontier have the same distances, we no longer need to keep the relative orders of the elements within each frontier. We hence directly use an unordered set to maintain each frontier. Since only three frontiers are useful during the BFS process, we also only allocate and keep three unordered sets and their roles rotate among the previous, the current and the next frontier.

The 2-level hash table introduced in Section 4.1.1 works perfectly well here since we have no control of the frontier sizes, which can be either greater or smaller than previous ones. For each frontier, we loop over all vertices in the set, and for each outgoing edge, if the other endpoint is not in either of the three unordered sets then it is added to the next frontier. Therefore, the operations to the sets include lookups and insertions. We do not explicitly delete the sets. Instead, when the searching of one frontier is finished, the new next frontier will reuse this space of the previous frontier.

6.1.3 Bidirectional search

This previous implementation based on rotating arrays can greatly reduce the number of write transfer when the frontier size is much smaller than the number of vertices, like grids, meshes, roadmaps, etc. For other graphs with larger frontier size and smaller diameter (e.g. real-world networks that follow power law), I/O cost is not improved significantly. However for these graphs, the average distance between every pair of vertices is usually very small, and this is called the “small-world phenomenon” [WS98].

To utilize this property, we employ bidirectional breadth-first search on computing the s-t shortest path. Assuming the distance between this pair of query vertices is d , then the search from each direction will only visit the vertices within distance $\lceil d/2 \rceil$. On these power-law graphs, the frontier size grows exponentially on the first several steps until most of the vertices are reached. The number of vertices on the top-half levels in the shortest-path tree is therefore far less than n . This difference however, is negligible on graphs with bounded frontier size, since the distance between a random pair of vertices is large expectedly, and the bidirectional search eventually reaches about the same amount of vertices unless d is much smaller than the graph diameter.

To implement bidirectional search efficiently, we also use three rotating arrays for the search process. The extra detail here is that, the source, destination, and all intermediate vertices are put together in the same sets, and use one bit (the sign bit) to identify if the vertex is by the search from the source or from the destination. We only keep three instead of six sets, which reduces the number of read transfers approximately by a half.

6.1.4 Experiment

In the experiment, we examine whether and how the new algorithms we just discussed improve the I/O cost. We implement four versions of breadth-first search: classic and bidirectional search with and without using rotating arrays. The experiment is run on 8 graph instances with various cache size.

Graph instances. We use graphs of various types from the SNAP datasets [LK14a] as the input of the algorithms. The graphs include the road networks in Pennsylvania and Texas (real-world planar graphs), the web graphs of the University of Notre Dame and Stanford University, the DBLP collaboration network, and the Youtube online social network (4 real-world networks). In the case of web graphs, each edge represents a hyperlink between two web pages. We also use synthetic graphs of 2D and 3D grids. For each of the graphs used, the numbers of vertices and edges are given in the table below. If a graph does not come equipped with weights, we assign to every edge a random integer between 1 and 10,000. The graphs we used are relatively small due to the overhead in our software and hardware simulator, but we adjust the cache size accordingly.

Type	Graph	# Vertices	# Edges
Sparse Graphs	Grid: 2D	1M	3.96M
	Grid: 3D	1M	5.94M
	Roadmap: Pennsylvania	1.09M	3.08M
	Roadmap: Texas	1.39M	3.84M
Social Networks	Webgraph: Notre Dame	325K	2.20M
	Webgraph: Stanford	281K	3.98M
	DBLP collab. network	317K	2.10M
	Youtube network	1.13M	5.98M

Overall performance. The numbers of read and write transfers on 8 graph instances with various cache size are given in Table 11. Their weighted I/O costs are given in Table 12 considering two typical values of ω , and visualized (merged into two categories) in Figure 6.

Algorithm	Classic BFS						BFS based on RA					
	500		2000		10000		500		2000		10000	
	RT	WT	RT	WT	RT	WT	RT	WT	RT	WT	RT	WT
2D Grid	18.47	6.17	17.60	6.01	10.88	4.59	8.75	0.84	6.52	0.16	6.20	0.01
3D Grid	19.38	5.45	16.42	5.31	13.51	4.51	88.24	4.26	31.03	2.01	6.39	0.38
PA Roadmap	8.73	2.78	8.07	2.54	2.43	1.05	29.61	2.95	5.37	0.65	1.26	0.02
TX Roadmap	6.79	2.16	6.01	1.91	1.95	0.85	23.45	2.29	4.32	0.50	1.00	0.02
Stan Webgraph	39.75	3.97	27.11	3.75	9.63	1.48	232.09	5.99	150.70	4.74	35.19	2.10
NDU Webgraph	3.47	1.04	3.02	0.97	2.52	0.77	84.67	4.80	58.91	3.91	18.33	1.73
DBLP Network	19.70	5.75	15.43	5.07	6.85	1.68	140.01	7.34	104.12	6.27	37.67	3.20
Youtube Network	13.23	2.47	9.64	2.29	5.63	1.84	91.70	6.77	71.98	6.39	39.95	5.08

Algorithm	Classic Bidirectional BFS						Bidirectional BFS based on RA					
	500		2000		10000		500		2000		10000	
	RT	WT	RT	WT	RT	WT	RT	WT	RT	WT	RT	WT
2D Grid	18.14	4.97	17.61	4.94	9.39	4.61	7.60	0.55	4.94	0.13	4.69	0.01
3D Grid	11.66	3.08	10.52	3.06	9.47	2.98	38.24	1.97	8.80	0.69	2.90	0.09
PA Roadmap	7.87	3.19	7.58	3.09	3.37	1.60	27.00	2.62	4.25	0.43	1.06	0.02
TX Roadmap	5.31	2.30	5.03	2.19	1.79	1.06	11.12	1.18	2.30	0.21	0.60	0.01
Stan Webgraph	4.70	1.38	4.24	1.37	2.74	1.29	7.32	0.84	3.00	0.47	0.63	0.13
NDU Webgraph	0.80	0.70	0.77	0.70	0.74	0.68	2.63	0.43	1.12	0.28	0.16	0.06
DBLP network	1.40	1.01	1.29	0.99	1.05	0.89	1.33	0.25	0.37	0.12	0.12	0.04
Youtube network	0.74	0.68	0.71	0.68	0.69	0.67	0.25	0.09	0.08	0.04	0.02	0.01

Table 11: Numbers of read and write transfers of BFS implementations on different cache sizes. Numbers are r/w transfers per vertex per **10** queries. We pick the number 10 to fit the numbers in one table.

To better understand the performance of different implementations, we also count various statistics of the searching including the depths and frontier sizes, and the results are generalized in Table 13.

Unidirectional search. Indicated in Table 11, classic BFS requires no more than one sequential write (pushed into the queue) and one random write (marking the distance) per vertex. The random write can be avoided if the associated cache line is cached, so better locality of vertex indices of some graphs (roadmaps and NDU webgraph) and larger cache size reduce writes per vertex. The algorithm also reads the edges of each vertex and checks whether the other endpoint is visited or not, which is also affected by the edges per vertex and the locality of vertex indices.

For the implementation using rotating arrays, the key factor is whether the frontier fits into the cache. Shown in Table 13, the sparse graphs (grids and roadmaps) have smaller frontier sizes, so as long as the cache can hold each of them, the writes are largely minimized. This is also true for reads since checking is always in the hash table. However, once the frontier is larger than the cache size, then each insertion to the hash table now becomes a random write and can hardly be cached because of the hash function. The reads are increased even more, since checking whether a vertex is visited can lead to at most 6 cache misses: three rotating arrays each with 2 levels. Summarized in Figure 6, the I/O cost is largely improved using the rotating arrays when frontiers fit into the cache, but deteriorates when not.

Bidirectional search. The I/O performance of classic bidirectional BFS is similar to the classic unidirectional BFS since they essentially search in a similar pattern. The bidirectional search requires fewer reads and

$\omega = 10$												
Algorithm	Classic BFS			BFS based on RA			Classic Bidir. BFS			Bidir. BFS based on RA		
Cache Size	500	2000	10000	500	2000	10000	500	2000	10000	500	2000	10000
2D Grid	80.18	77.72	56.74	17.13	8.15	6.25	67.80	67.04	55.53	13.14	6.27	4.76
3D Grid	73.85	69.49	58.60	130.86	51.16	10.18	<i>42.44</i>	41.16	39.29	57.98	15.65	3.82
PA Roadmap	<i>36.50</i>	33.42	12.89	59.10	11.84	1.44	<i>39.76</i>	38.46	19.39	53.24	8.52	1.24
TX Roadmap	<i>28.42</i>	25.14	10.44	46.36	9.36	1.18	<i>28.27</i>	26.95	12.37	22.89	4.37	0.69
Stan Webgraph	<i>79.41</i>	<i>64.64</i>	<i>24.39</i>	291.99	198.08	56.18	18.48	17.93	15.62	15.75	7.71	1.93
NDU Webgraph	<i>13.84</i>	<i>12.70</i>	<i>10.20</i>	132.71	98.04	35.58	7.80	7.73	7.54	6.90	3.90	0.75
DBLP network	<i>77.19</i>	<i>66.16</i>	<i>23.67</i>	213.43	166.84	69.63	11.53	11.22	9.91	3.84	1.61	0.56
Youtube network	<i>37.92</i>	<i>32.59</i>	<i>24.06</i>	159.43	135.92	90.72	7.58	7.50	7.38	1.17	0.51	0.17

$\omega = 100$												
Algorithm	Classic BFS			BFS based on RA			Classic Bidir. BFS			Bidir. BFS based on RA		
Cache Size	500	2000	10000	500	2000	10000	500	2000	10000	500	2000	10000
2D Grid	635.5	618.7	469.4	92.5	22.7	6.7	514.7	511.8	470.7	62.9	18.2	5.3
3D Grid	564.0	547.1	464.3	514.4	232.3	44.2	319.4	316.9	307.6	235.6	77.4	12.0
PA Roadmap	<i>286.4</i>	261.5	107.0	324.4	70.0	3.0	<i>326.7</i>	316.4	163.6	289.4	47.0	2.8
TX Roadmap	<i>223.1</i>	197.3	86.8	252.5	54.7	2.8	<i>234.8</i>	224.2	107.6	128.8	23.0	1.5
Stan Webgraph	<i>436.3</i>	<i>402.4</i>	157.1	831.1	624.5	245.1	142.5	141.2	131.6	91.6	50.1	13.6
NDU Webgraph	<i>107.1</i>	<i>99.8</i>	79.2	565.0	450.2	190.8	70.8	70.3	68.8	45.3	28.9	6.0
DBLP network	<i>594.6</i>	<i>522.6</i>	175.1	874.2	731.3	357.3	102.6	100.5	89.7	26.4	12.8	4.5
Youtube network	<i>260.2</i>	<i>239.1</i>	190.0	769.0	711.4	547.6	69.1	68.5	67.6	9.4	4.3	1.5

Table 12: I/O costs of BFS implementations on different cache sizes. Numbers of r/w transfers are from Table 11. Italic-font number indicates that the classic implementation has a lower I/O cost in that setting.

	Unidirectional search			Bidirectional search		
	depth	maximum frontier size	average frontier size	depth	maximum frontier size	average frontier size
2D Grid	1571.4	1460	636.4	370.2	2091	1102.8
3D Grid	229.9	13277	4349.7	49.3	15058	4628.9
PA Roadmap	574.5	5032	1893.1	163.9	7204	2718.3
TX Roadmap	748.6	6058	1804.9	174.7	5752	2102.1
Stan Webgraph	107.1	87614	2383.4	3.5	50094	3782.8
NDU Webgraph	29.2	124519	11155.1	3.8	51562	2053.5
DBLP Network	16.1	129283	19694.4	3.9	36373	1420.1
Youtube Network	15.8	607534	71828.5	2.7	1841	107.4

Table 13: The depths (number of levels in the shortest-path tree) and frontier sizes (number of vertices) during the search processes. Note that the numbers are averaged from multiple searches.

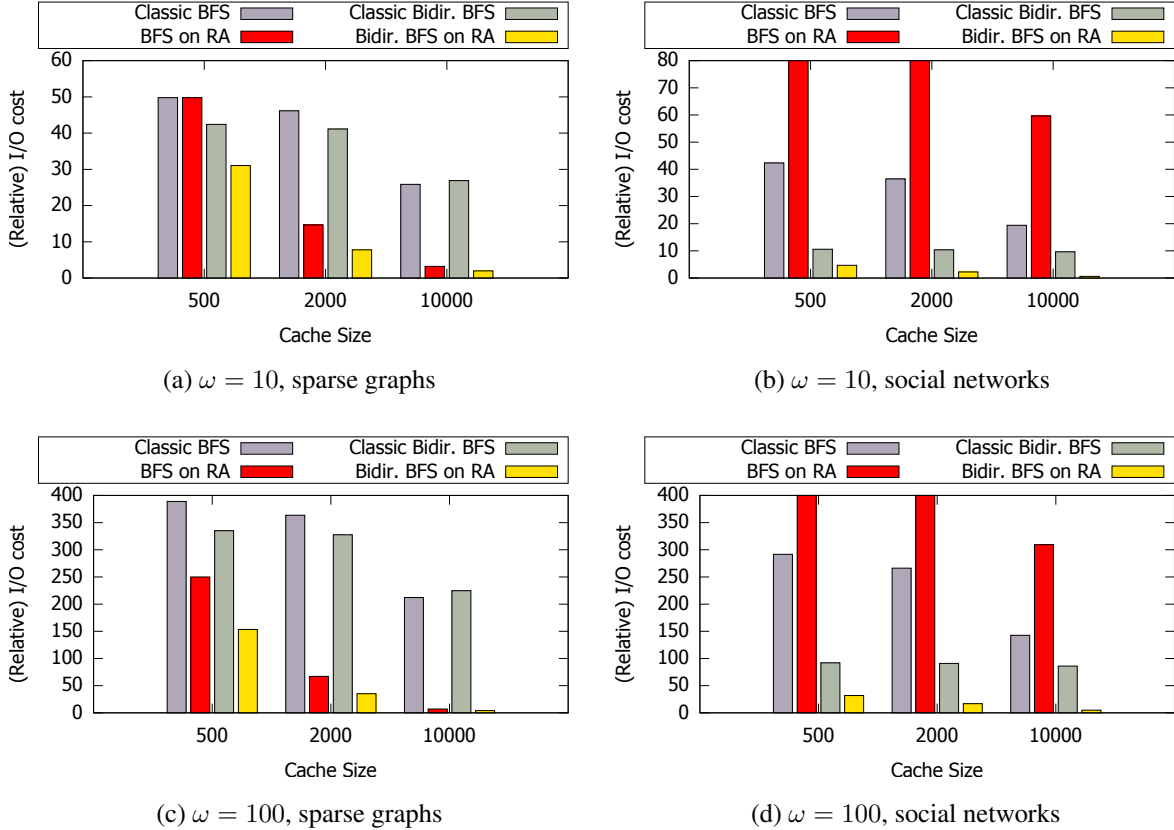


Figure 6: The trends of the I/O costs of four different implementations of BFS. The new implementations shown in this paper are the BFS and bidirectional BFS based on rotating arrays (red and orange bars). Graphs used in the experiment are shown in Section 6.1.4 and categorized into sparse (almost planar) graphs and social networks. We show the relative I/O cost based on varied cache sizes, and each number is geometric mean of the four graphs in that category (the exact numbers are given in Table 12). We can see the consistent advantages of the new BFS implementation on sparse graphs, and the improvement of the new bidirectional version in all cases. Notice that in (b) and (d) some values exceed the ranges of vertical axis.

writes since it strictly searches fewer vertices, especially in social networks since these graphs have smaller diameters and the two searches usually meet earlier before the majority of the vertices are visited.

This special property of social network largely helps our implementation using rotating arrays. Shown in Table 13, the two bidirectional searches use 2.7-3.9 rounds on average to meet, which preserves the frontier in the cache during the search even for very small cache sizes. As a result, the bidirectional BFS using rotating arrays has a constantly good performance on all combinations of graphs and cache sizes, which is shown in Table 12 and Figure 6.

6.1.5 Conclusions

We discuss how to efficiently implement BFS in the asymmetric setting and experiment the I/O performance for four implementations on a variety of graphs. We show that if the query is s-t (pairwise) distance, our bidirectional BFS using rotating arrays shows an overwhelming advantage in all cases we tested. If all

distances to the source are required, the unidirectional BFS using rotating arrays has a better performance if the cache can hold each frontier.

6.2 Dijkstra’s Algorithm

Dijkstra’s Algorithm [Dij59] is a well-known algorithm to compute single-source shortest paths on a non-negative weighted graph $G = (V, E)$. Due to the rapid growth of the data size, real-world graphs nowadays can easily go beyond the size of the order of gigabytes, and they need to be stored on the large non-volatile memory. Running the classic implementation of Dijkstra’s algorithm can be costly in this setting. We show that with an appropriate implementation, the algorithm can write much fewer to the non-volatile memory, which further leads to much lower I/O cost.

Throughout this section we assume the input graph $G = (V, E)$ contains $n = |V|$ vertices and $m = |E|$ edges.

Dijkstra’s algorithm maintains a set of visited vertices associated with their shortest distances to the source (denoted by d_v for $v \in V$), and the unvisited neighbors of these vertices form the “frontier” (d_v for unvisited vertex v is $+\infty$). Each vertex u in the frontier stores a tentative distance to be $\min_{v \in N(u)} \{d_v + e_{u,v}\}$ where $N(u)$ is the incoming neighbor set of vertex u . Initially the visited set contains only one vertex: the source node. The invariant of this algorithm is that, the minimum tentative distance in the frontier set is indeed the shortest distance of this vertex, and thus Dijkstra’s algorithm iteratively move this vertex from the frontier to the visited set and update the frontier accordingly. The correctness can easily be proven by induction on the number of visited nodes.

Due to the widespread use of Dijkstra’s algorithm, there are plenty of works on the efficient implementations of Dijkstra’s algorithm and we refer the readers to some recent works [MS03, BGST16] for summaries of the work bounds of different approaches. Different implementations have different costs on the two operations in Dijkstra’s algorithm: EXTRACT-MIN that finds and removes the vertex with minimum distance in the frontier set, and DECREASE-KEY that updates the tentative distance of the other endpoint of an edge of this vertex removed from the frontier in this iteration. The data type that supports the queries is abstracted as a priority queue. Specifically when the priority queue is implemented using Fibonacci Heap [FT87] to maintain the frontier set, the overall time complexity is $O(m + n \log n)$.

6.2.1 Classic implementation using a binary heap

Although there are many advanced implementations of the priority queue with lower time complexities, the constant hidden in the asymptotic bound is large. In practice the classic implementation using a binary heap works reasonably well on general sparse real-world graphs, and its wall-clock performance is competitive or even better on modern computer architecture. We hence implement it as a baseline and measure the number of read and write transfers of this algorithm as a comparison to our write-efficient version.

For each vertex, we maintain the shortest distance in a global array with size n . For practical purpose, instead of initializing the priority queue of size n with infinite distances, we insert a vertex when it is first added to the frontier (so the priority queue needs to support INSERT, which most implementations do). The binary heap only stores the indices of the vertices which optimizes the memory footprint and number of write transfers. To perform DECREASE-KEY in a binary heap efficiently, we keep another global array, an auxiliary structure that maps each vertex to its position in the heap, and is maintained up-to-date as the priority queue changes. This implementation has worst-case time complexity to be $O(m \log n)$. Since this implementation does not take any special optimization on caching, the I/O cost is therefore $O(\omega m \log(nB/M))$, assuming the cache always keeps the first M/B vertices in the priority queue.

6.2.2 Phased Dijkstra

Phased Dijkstra [BFG⁺16] is a specific implementation of Dijkstra’s algorithm that the goal is to fully maintain the priority queue in small-memory. It only requires $O(n)$ writes to large-memory, the shortest distances to all vertices. The idea is to partition the computation into phases such that for a parameter M' each phase keeps a priority queue of size at most $(1 + \epsilon)M'$ and visits at least M' vertices. By selecting $M' = M/c$ for an appropriate constant c , the priority queue fits in small-memory, and the only writes to large-memory are the final distances.

Algorithm 3: Phased Dijkstra

Input: A connected weighted graph $G = (V, E)$ and a source vertex s
Output: The shortest distances $\delta = \{\delta_1, \dots, \delta_n\}$ from source s

```

1 Priority Queue  $P \leftarrow \emptyset$ 
2 Mark vertex  $s$  as visited and set  $\delta(s) \leftarrow 0$ 
3 while there exists unvisited vertices do
4   if  $P = \emptyset$  then
5     Scan over all edges in  $E$  and store at most  $M'$  closest unvisited vertices in  $P$ 
6     if  $|P| = M'$  then
7       Set  $d_{max}$  as the distance to the farthest vertex in  $P$ 
8     else
9        $d_{max} \leftarrow +\infty$ 
10     $u = P.EXTRACT-MIN()$ 
11    Set  $\delta_u$  as the distance from  $s$  to  $u$ , and mark  $u$  as visited
12    foreach  $(u, v, dis_{u,v}) \in E$  do
13      if  $\delta_u + dis_{u,v} < d_{max}$  then
14        if  $v \in P$  then
15           $P.DECREASE-KEY(v, \delta_u + dis_{u,v})$ 
16        else
17           $P.INSERT(v, \delta_u + dis_{u,v})$ 
18          if  $|P| = (1 + \epsilon)M'$  then
19            Remove  $\epsilon M'$  vertices with larger distances in  $P$ 
20            Set  $d_{max}$  as the farthest distance in  $P$ 
21 return  $\delta(\cdot)$ 

```

The pseudocode of the algorithm is provided in Algorithm 3. Technically the priority queue P can be implemented using an arbitrary heap since it is fully in the small-memory and will not affect the I/O cost. In our experiment we implement it using a binary heap.

In phased Dijkstra, each phase starts and ends with an empty priority queue P . The priority queue is kept small by discarding the $\epsilon M'$ largest elements (vertex distances) whenever $|P| = (1 + \epsilon)M'$. To achieve this, P is flattened into an array, the M' -th smallest element d_{max} is found by selection, and the priority queue is reconstructed from the elements no greater than d_{max} , all taking linear time. After such a truncation, all further insertions in a given phase are not added to P if they have a value greater than d_{max} .

The processing of each phase in phased Dijkstra consists of two parts. The first part (line 5–9) of each phase loops over all edges in the graph and relaxes any that go from a visited to an unvisited vertex (possibly inserting or decreasing a key in P). The second part (repeatedly loop over line 10–20) then runs standard Dijkstra’s algorithm repeatedly visiting the vertex with minimum distance and relaxing its neighbors until P is empty. Similar to other implementations, to implement relax, the algorithm needs to know whether a vertex is already in P , and if so where in P so that it can do a decrease-key on it. However in phased Dijkstra it is too costly to store this information using a global array. Instead, we use an unordered map introduced in Section 4.1.1 for this mapping. Theoretically the hash table can be set with fixed size, but in practice we use

a 2-level hash table since it leads to better performance when the frontier size is consistently small, and equal performance otherwise. This map is referred as *vertex map* later.

The I/O cost of phased Dijkstra is $Q(n, m) = O\left(n\left(\frac{m}{M} + \omega\right)\right)$. More details on the correctness and complexity analysis can be found in [BFG⁺16].

We make a special optimization that once all outgoing edges of a vertex are visited, we remove this node and all associated edges in the scan in Line 5. This is done by using the sign-bit of the output distances, such that it is more likely being cached. We call the *active set* that contains visited but not removed vertices.

For an efficient implementation that optimizes memory footprint and I/O efficiency, each heap element contains the vertex index, the pointer to the vertex map, and the tentative distance. In total it takes 16 bytes. Each element in the vertex map only stores a pointer, and the vertex index can be check via the corresponding heap element. Throughout our experiment we set the maximum occupancy rate of the 2-level hash table to be 0.8, and truncation ratio $\epsilon = 0.25$. M' is chosen such that the hash table and priority queue occupy about 40% of the small-memory, while various values of M' are also discussed.

6.2.3 Experiments

The experiments are run on two Dijkstra implementations with different parameter combinations on cache size, cache policy, and the priority queue size. The experiment is run on eight graph instances with various cache size.

Since in phased Dijkstra the number of reads is significantly more than that of writes, to keep the priority queue in the cache, the special cache strategy in Section 3.2 is required. In Table 17 we show that different cache policies only cause minor differences, so in the majority of this section we use the **Static** policy.

Overall performance. In Table 14 we show the number of read and write transfers of two implementations on different graphs with various cache size. Cache size varies from 2,000 to 50,000 cache lines each with 64 bytes. In Table 15 the overall I/O costs with different values of ω are listed based on the numbers in Table 14.

For the binary-heap implementation, the actual reads and writes mostly match the theoretical bound $O(m \log(nB/M))$. Reads are about slightly less than twice as writes: each edge is read once during Dijkstra (linear scan and requires no modification), and an update in the binary heap always requires one more read than writes. The only exception is on the roadmaps when the frontier size is consistently small while nearby vertices share contiguous indices. In this case the cache efficiently holds the entire heap and leads to fewer reads and writes to the large-memory.

For phased Dijkstra, we first observed that the number of writes is always no more than 1.3 per vertex: one write per vertex when the distance is finalized, plus some other cost to maintain the active set. The number of read transfers is mainly decided the number of phases, and the size of the active set (the edge lists of active vertices at the beginning of each phase needs to be scanned).

The overall I/O performance (shown in Table 15) indicates that phased Dijkstra is consistently better than the binary-heap version except for the only case that both ω and cache size are small and the active set is large. This case can hardly happen in practice since the cache size is even much smaller than the current L3 cache. The improvement on I/O cost in all cases is up to 3 and 7.6 when ω is 10 and 100. This peak occurs when the small-memory sizes is large and frontier size is larger, since at this time the cost to maintain the binary heap in the classic implementation is costly, while the extra cost to run phased Dijkstra is insignificant.

Different cache maintenance policy. We show the number of read and write transfer of phased Dijkstra on two different cache policies in 17. Different policies actually give a very similar performance in the graph instances. More analysis is shown in the full version of this paper.

Cache Size	Classic Dijkstra using binary heap						Phased Dijkstra								
	2k		10k		50k		2k		10k		50k				
	RT	WT	RT	WT	RT	WT	RT	WT	RT	WT	RT	WT			
2D Grid	10.37	4.87	4.61	2.34	3.53	1.89	5.93	1.12	(1)	3.62	1.11	(1)	2.74	1.04	(1)
3D Grid	25.87	14.53	17.44	8.54	8.20	3.35	126.56	1.12	(383)	28.16	1.12	(44)	5.90	1.07	(1)
PA Roadmap	7.82	4.14	2.16	0.81	1.12	0.40	30.40	0.99	(202)	1.89	0.52	(1)	1.04	0.28	(1)
TX Roadmap	7.69	4.10	2.16	0.82	1.09	0.40	37.93	0.97	(262)	1.88	0.53	(1)	1.01	0.27	(1)
Stan Webgraph	37.31	18.13	26.73	12.06	10.60	3.91	404.61	1.03	(100)	81.02	1.02	(16)	13.83	0.88	(2)
NDU Webgraph	24.97	15.34	15.14	8.21	5.19	1.65	99.86	1.10	(132)	22.57	1.01	(24)	5.81	0.67	(3)
DBLP Network	32.14	19.66	23.17	13.02	10.08	4.58	341.08	1.12	(131)	64.75	1.10	(24)	11.53	0.94	(4)
Youtube Network	35.19	22.86	26.80	16.53	21.18	12.90	836.64	1.13	(477)	161.36	1.11	(93)	33.08	1.05	(18)

Table 14: Number of read and write transfers of different Dijkstra implementations on different cache size. Numbers are normalized to read or write transfers per vertex. We ran SSSP queries on 10 different randomly-chosen source nodes. Numbers in the parentheses are the average phases.

Cache Size	$\omega = 10$						$\omega = 100$					
	Classic			Phased Dijkstra			Classic			Phased Dijkstra		
	2k	10k	50k	2k	10k	50k	2k	10k	50k	2k	10k	50k
2D Grid	59.1	28.0	22.4	17.1	14.7	13.1	497.8	238.9	192.2	118.3	114.3	106.8
3D Grid	171.2	102.8	41.7	137.8	39.3	16.6	1478.9	871.2	343.5	239.0	140.1	113.1
PA Roadmap	49.2	10.2	5.1	40.3	7.1	3.8	421.4	83.1	41.2	129.0	53.9	28.7
TX Roadmap	48.7	10.4	5.1	47.6	7.1	3.7	417.5	84.3	40.9	134.5	54.5	28.3
Stan Webgraph	218.6	147.3	49.7	414.9	91.2	22.6	1850.2	1232.2	401.8	507.7	182.8	101.9
NDU Webgraph	178.4	97.2	21.6	110.9	32.7	12.5	1559.4	836.3	169.7	209.9	123.9	72.8
DBLP Network	228.7	153.4	55.9	352.3	75.7	20.9	1997.6	1325.2	468.4	453.1	174.7	105.9
Youtube Network	263.8	192.1	150.1	847.9	172.5	43.6	2321.0	1679.4	1310.8	949.3	272.8	138.3

Table 15: The I/O costs on different Dijkstra implementations on different cache size. The write-read ratio ω are selected to be typical projected values 10 (latency, bandwidth) and 100 (energy). Results are based on the numbers in Table 14.

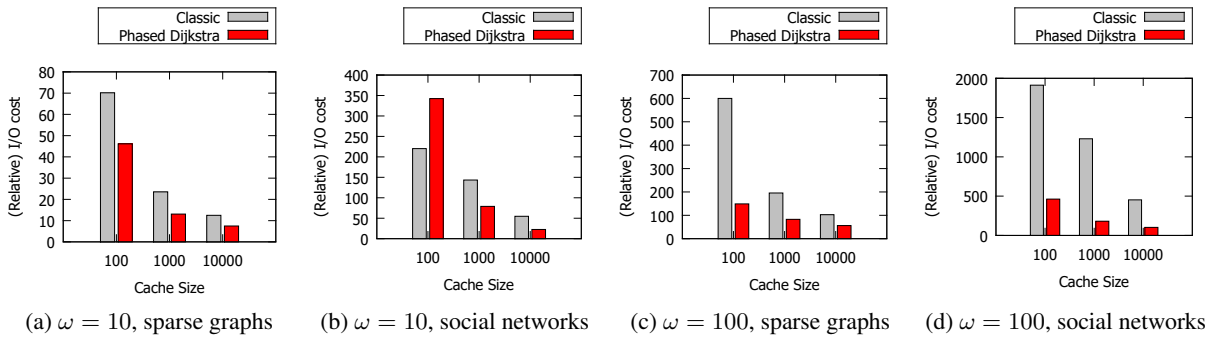


Figure 7: I/O costs of classic Dijkstra and phased Dijkstra. Graphs are categorized into sparse (almost planar) graphs and social networks and the I/O cost is geometric mean of the four. Numbers are from Table 15.

Cache Policy	SplitPool		Static	
	RT	WT	RT	WT
2D Grid	3.74	1.11	3.62	1.11
3D Grid	27.43	1.12	28.16	1.12
PA Roadmap	1.84	0.88	1.89	0.52
TX Roadmap	1.83	0.87	1.88	0.53
Stan Webgraph	81.40	1.02	81.02	1.02
NDU Webgraph	21.96	1.07	22.57	1.01
DBLP Network	65.20	1.11	64.75	1.10
Youtube Network	163.38	1.12	161.36	1.11

Table 16: Number of read and write transfer of phased Dijkstra on two different cache policies: the **SplitPool** policy and the **Static** policy. The cache contains 10,000 cache-lines. Different policies give very similar performance except for the writes on roadmaps.

Picking appropriate priority queue’s size M' . In previous experiment we empirically set the overall size of the priority queue and vertex map to be about 40% of the overall cache size. However, this percentage can affect the performance of phased Dijkstra. Larger percentage leads the heap to contain more elements so that the overall number of phases is decreased. Meanwhile the size left for the rest cache is smaller, which decreases the cache performance. Hence, the number of phases and specific graph property lead to different optimality point of the the priority queue’s size and no easy conclusions can be drawn. In Table 18 we show a snapshot on the cache size of 10,000 and the percentage varies from 30% to 60%. The trend is that, when the priority queue’s size is smaller than the average frontier size then larger priority queue helps, and vice visa. In general different priority queue’s sizes only make a minor difference on I/O cost and will not affect the relatively lower cost comparing to the binary-heap implementation.

6.2.4 Conclusions

We discussed phased Dijkstra and experiment its performance on a variety of graphs. The high-level idea is to fit the computation within the small-memory (i.e. the cache) and thus requires no intermediate writes to the large asymmetry memory. The experimental results show that phased Dijkstra consistently outperforms the binary-heap version on I/O cost except for the combination of small ω ($= 10$), small cache size, and on social networks. Although phased Dijkstra contains some parameters, we also show that they do not affect the efficiency of phased Dijkstra when they are within a reasonable large range. A similar case also holds for different cache policies.

Notice that the idea here that fits the computation in the small-memory can also be applied to computing minimum spanning tree, sorting, and many other problems.

References

- [ABCR12] Manos Athanassoulis, Bishwaranjan Bhattacharjee, Mustafa Canim, and Kenneth A. Ross. Path processing using solid state storage. In *Proc. International Workshop on Accelerating Data Management Systems Using Modern Processor and Storage Architectures (ADMS)*, 2012.

Cache Policy	SplitPool		Static	
	RT	WT	RT	WT
2D Grid	3.74	1.11	3.62	1.11
3D Grid	27.43	1.12	28.16	1.12
PA Roadmap	1.84	0.88	1.89	0.52
TX Roadmap	1.83	0.87	1.88	0.53
Stan Webgraph	81.40	1.02	81.02	1.02
NDU Webgraph	21.96	1.07	22.57	1.01
DBLP Network	65.20	1.11	64.75	1.10
Youtube Network	163.38	1.12	161.36	1.11

Table 17: Number of read and write transfer of phased Dijkstra on two different cache policies: the **SplitPool** policy and the **Static** policy. The cache contains 10,000 cache-lines. Different policies give very similar performance except for the writes on roadmaps.

Priority Queue Size	30%		40%		50%		60%	
	RT	WT	RT	WT	RT	WT	RT	WT
2D Grid	8.4	1.10	8.6	1.11	8.9	1.11	9.3	1.11
3D Grid	33.7	1.12	28.9	1.12	25.6	1.12	23.3	1.12
PA Roadmap	2.7	0.48	2.8	0.52	3.0	0.58	3.2	0.66
TX Roadmap	2.6	0.48	2.7	0.53	2.9	0.59	3.1	0.66
Stan Webgraph	94.4	1.02	82.9	1.02	72.1	1.02	67.1	1.02
NDU Webgraph	26.1	1.01	22.8	1.01	20.2	1.03	18.5	1.04
DBLP Network	75.1	1.10	64.9	1.10	57.3	1.10	51.5	1.11
Youtube Network	183.8	1.11	161.5	1.11	144.4	1.12	134.9	1.12

Table 18: Number of read and write transfers of phased Dijkstra on different priority queue’s size. The overall percentage of priority queue and vertex map varies from 30% to 60% comparing to the cache size, which is 10,000 cache-lines.

- [ACM⁺11] Ameen Akel, Adrian M. Caulfield, Todor I. Mollov, Rajech K. Gupta, and Steven Swanson. Onyx: A prototype phase change memory storage array. In *Proc. USENIX Workshop on Hot Topics in Storage and File Systems (HotStorage)*, 2011.
- [Ada92] Stephen Adams. Implementing sets efficiently in a functional language. Technical Report CSTR 92-10, University of Southampton, 1992.
- [Ada93] Stephen Adams. Efficient sets—a balancing act. *Journal of functional programming*, 3(04):553–561, 1993.
- [AL63] M AdelsonVelskii and Evgenii Mikhailovich Landis. An algorithm for the organization of information. Technical report, DTIC Document, 1963.
- [ALT16] Mahdi Amani, Kevin A. Lai, and Robert E. Tarjan. Amortized rotation cost in AVL trees. *Information Processing Letters*, 116(5):327–330, 2016.

- [AS89] C.R. Aragon and R.G. Seidel. Randomized search trees. In *Foundations of Computer Science*, pages 540–545, 1989.
- [AV88] Alok Aggarwal and Jeffrey S. Vitter. The Input/Output complexity of sorting and related problems. *Communications of the ACM*, 31(9), 1988.
- [BAT06] Avraham Ben-Aroya and Sivan Toledo. Competitive analysis of flash-memory algorithms. In *Proc. European Symposium on Algorithms (ESA)*, 2006.
- [Bay72] Rudolf Bayer. Symmetric binary B-trees: Data structure and maintenance algorithms. *Acta informatica*, 1(4):290–306, 1972.
- [BDBF⁺16] Naama Ben-David, Guy E. Blelloch, Jeremy T. Fineman, Phillip B. Gibbons, Yan Gu, Charles McGuffey, and Julian Shun. Parallel algorithms for asymmetric read-write costs. In *Proc. 28th ACM Symposium on Parallelism in Algorithms and Architectures (SPAA)*, 2016.
- [BDBF⁺18] Naama Ben-David, Guy E. Blelloch, Jeremy T. Fineman, Phillip B. Gibbons, Yan Gu, Charles McGuffey, and Julian Shun. Implicit decomposition for write-efficient connectivity algorithms. In *Proc. IEEE International Parallel & Distributed Processing Symposium (IPDPS)*, 2018.
- [BFG⁺15] Guy E. Blelloch, Jeremy T. Fineman, Phillip B. Gibbons, Yan Gu, and Julian Shun. Sorting with asymmetric read and write costs. In *Proc. 27th ACM Symposium on Parallelism in Algorithms and Architectures (SPAA)*, 2015.
- [BFG⁺16] Guy E. Blelloch, Jeremy T. Fineman, Phillip B. Gibbons, Yan Gu, and Julian Shun. Efficient algorithms with asymmetric read and write costs. In *Proc. European Symposium on Algorithms (ESA)*, pages 14:1–14:18, 2016.
- [BFS16] Guy E. Blelloch, Daniel Ferizovic, and Yihan Sun. Just join for parallel ordered sets. In *Proc. 28th ACM Symposium on Parallelism in Algorithms and Architectures*, pages 253–264. ACM, 2016.
- [BGSS18] Guy E. Blelloch, Yan Gu, Julian Shun, and Yihan Sun. Parallel write-efficient geometric algorithms. In *ACM Symposium on Parallelism in Algorithms and Architectures (SPAA)*. ACM, 2018.
- [BGST16] Guy E. Blelloch, Yan Gu, Yihan Sun, and Kanat Tangwongsan. Parallel shortest paths using radius stepping. In *Proc. 28th ACM Symposium on Parallelism in Algorithms and Architectures (SPAA)*, SPAA '16, pages 443–454, 2016.
- [CDG⁺16] Erin Carson, James Demmel, Laura Grigori, Nicholas Knight, Penporn Koanantakool, Oded Schwartz, and Harsha Vardhan Simhadri. Write-avoiding algorithms. In *Proc. IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, pages 648–658, 2016.
- [CGN11] Shimin Chen, Phillip B. Gibbons, and Suman Nath. Rethinking database algorithms for phase change memory. In *Proc. Conference on Innovative Data Systems Research (CIDR)*, 2011.
- [CL09] Sangyeun Cho and Hyunjin Lee. Flip-N-Write: A simple deterministic technique to improve PRAM write performance, energy and endurance. In *Proc. IEEE/ACM International Symposium on Microarchitecture (MICRO)*, 2009.

- [CNF⁺09] Jeremy Condit, Edmund B Nightingale, Christopher Frost, Engin Ipek, Benjamin Lee, Doug Burger, and Derrick Coetzee. Better i/o through byte-addressable, persistent memory. In *Proc. ACM SIGOPS 22nd symposium on Operating systems principles*, pages 133–146. ACM, 2009.
- [Dij59] Edsger W. Dijkstra. A note on two problems in connexion with graphs. *Numerische mathematik*, 1(1), 1959.
- [DJX09] Xiangyu Dong, Norman P. Jouupi, and Yuan Xie. PCRAMsim: System-level performance, energy, and area modeling for phase-change RAM. In *Proc. ACM International Conference on Computer-Aided Design (ICCAD)*, 2009.
- [DWS⁺08] Xiangyu Dong, Xiaoxia Wu, Guangyu Sun, Yuan Xie, Hai H. Li, and Yiran Chen. Circuit and microarchitecture evaluation of 3D stacking magnetic RAM (MRAM) as a universal memory replacement. In *Proc. ACM Design Automation Conference (DAC)*, 2008.
- [EGMP14] David Eppstein, Michael T. Goodrich, Michael Mitzenmacher, and Pawel Pszozna. Wear minimization for cuckoo hashing: How not to throw a lot of eggs into one basket. In *Proc. ACM International Symposium on Experimental Algorithms (SEA)*, 2014.
- [FT87] Michael L. Fredman and Robert E. Tarjan. Fibonacci heaps and their uses in improved network optimization algorithms. *Journal of the ACM*, 34(3), 1987.
- [GS78] Leo J Guibas and Robert Sedgwick. A dichromatic framework for balanced trees. 1978.
- [GT05] Eran Gal and Sivan Toledo. Algorithms and data structures for flash memories. *ACM Computing Surveys*, 37(2), 2005.
- [HP15] HP, SanDisk partner on memristor, ReRAM technology. <http://www.bit-tech.net/news/hardware/2015/10/09/hp-sandisk-reram-memristor>, October 2015.
- [HZX⁺14] Jingtong Hu, Qingfeng Zhuge, Chun Jason Xue, Wei-Che Tseng, Shouzhen Gu, and Edwin Sha. Scheduling to optimize cache utilization for non-volatile main memories. *IEEE Transactions on Computers*, 63(8), 2014.
- [IBM14] www.slideshare.net/IBMZRL/theseus-pss-nvmw2014, 2014.
- [Int15] Intel and Micron produce breakthrough memory technology. http://newsroom.intel.com/community/intel_newsroom/blog/2015/07/28/intel-and-micron-produce-breakthrough-memory-technology, July 2015.
- [JS17] Riko Jacob and Nodari Sitchinava. Lower bounds in the asymmetric external memory model. In *Proc. 29th ACM Symposium on Parallelism in Algorithms and Architectures (SPAA)*, pages 247–254, 2017.
- [KSDC14] Hyojun Kim, Sangeetha Seshadri, Clement L. Dickey, and Lawrence Chu. Evaluating phase change memory for enterprise storage systems: A study of caching and tiering approaches. In *Proc. USENIX Conference on File and Storage Technologies (FAST)*, 2014.

- [LIMB09a] Benjamin C. Lee, Engin Ipek, Onur Mutlu, and Doug Burger. Architecting phase change memory as a scalable DRAM alternative. In *Proc. ACM International Symposium on Computer Architecture (ISCA)*, 2009.
- [LIMB09b] Benjamin C Lee, Engin Ipek, Onur Mutlu, and Doug Burger. Architecting phase change memory as a scalable dram alternative. In *ACM SIGARCH Computer Architecture News*, volume 37, pages 2–13. ACM, 2009.
- [LK14a] Jure Leskovec and Andrej Krevl. SNAP Datasets: Stanford large network dataset collection. 2014.
- [LK14b] Jacob Leverich and Christos Kozyrakis. Reconciling high server utilization and sub-millisecond quality-of-service. In *Proc. Ninth European Conference on Computer Systems*, page 4. ACM, 2014.
- [LRR⁺03] Charles Lefurgy, Karthick Rajamani, Freeman Rawson, Wes Felter, Michael Kistler, and Tom W Keller. Energy management for commercial servers. *Computer*, 36(12):39–48, 2003.
- [MDS⁺15] Jasmina Malicevic, Subramanya Dulloor, Narayanan Sundaram, Nadathur Satish, Jeff Jackson, and Willy Zwaenepoel. Exploiting nvm in large-scale graph analytics. In *Proc. 3rd Workshop on Interactions of NVM/FLASH with Operating Systems and Workloads*, page 2. ACM, 2015.
- [Mic] MARSSx86. <http://marss86.org>.
- [MS03] Ulrich Meyer and Peter Sanders. Δ -stepping: a parallelizable shortest path algorithm. *J. Algorithms*, 49(1), 2003.
- [MSCT14] Jagan S. Meena, Simon M. Sze, Umesh Chand, and Tseung-Yuan Tseng. Overview of emerging nonvolatile memory technologies. *Nanoscale Research Letters*, 9, 2014.
- [MSG⁺12] Krishna T Malladi, Ian Shaeffer, Liji Gopalakrishnan, David Lo, Benjamin C Lee, and Mark Horowitz. Rethinking dram power modes for energy proportionality. In *Proc. 45th Annual IEEE/ACM International Symposium on Microarchitecture*, pages 131–142. IEEE Computer Society, 2012.
- [NR73] Jürg Nievergelt and Edward M Reingold. Binary search trees of bounded balance. *SIAM journal on Computing*, 2(1):33–43, 1973.
- [PS09] Hyoungmin Park and Kyuseok Shim. FAST: Flash-aware external sorting for mobile database systems. *Journal of Systems and Software*, 82(8), 2009.
- [PTL] PTLsim. <http://www.ptlsim.org>.
- [QGR11] Moinuddin K. Qureshi, Sudhanva Gurusurthi, and Bipin Rajendran. *Phase Change Memory: From Devices to Systems*. Morgan & Claypool, 2011.
- [QSR09] Moinuddin K Qureshi, Vijayalakshmi Srinivasan, and Jude A Rivers. Scalable high performance main memory system using phase-change memory technology. *ACM SIGARCH Computer Architecture News*, 37(3):24–33, 2009.

- [SA96] Raimund Seidel and Cecilia R Aragon. Randomized search trees. *Algorithmica*, 16(4-5):464–497, 1996.
- [SFB18] Yihan Sun, Daniel Ferizovic, and Guy E Belloch. PAM: parallel augmented maps. In *ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP)*, pages 290–304, 2018.
- [SK13] Daniel Sanchez and Christos Kozyrakis. Zsim: fast and accurate microarchitectural simulation of thousand-core systems. In *ACM SIGARCH Computer Architecture News*, volume 41, pages 475–486. ACM, 2013.
- [ST85] Daniel Dominic Sleator and Robert Endre Tarjan. Self-adjusting binary search trees. *Journal of the ACM (JACM)*, 32(3):652–686, 1985.
- [Tar83] Robert Endre Tarjan. *Data Structures and Network Algorithms*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 1983.
- [Vig12] Stratis D. Viglas. Adapting the B⁺-tree for asymmetric I/O. In *Proc. East European Conference on Advances in Databases and Information Systems (ADBIS)*, 2012.
- [Vig14] Stratis D. Viglas. Write-limited sorts and joins for persistent memory. *Proc. VLDB Endowment*, 7(5), 2014.
- [WS98] Duncan J. Watts and Steven H. Strogatz. Collective dynamics of “small-world” networks. *Nature*, 393(6684):440–442, 1998.
- [XDJX11] Cong Xu, Xiangyu Dong, Norman P. Jouppi, and Yuan Xie. Design implications of memristor-based RRAM cross-point structures. In *Proc. IEEE Design, Automation and Test in Europe (DATE)*, 2011.
- [YLK⁺07] Byung-Do Yang, Jae-Eun Lee, Jang-Su Kim, Junghyun Cho, Seung-Yun Lee, and Byoung-Gon Yu. A low power phase-change random access memory using a data-comparison write scheme. In *Proc. IEEE International Symposium on Circuits and Systems (ISCAS)*, 2007.
- [Yol13] Yole Developpement. Emerging non-volatile memory technologies, 2013.
- [ZWT13] Omer Zilberberg, Shlomo Weiss, and Sivan Toledo. Phase-change memory: An architectural perspective. *ACM Computing Surveys*, 45(3), 2013.
- [ZZYZ09] Ping Zhou, Bo Zhao, Jun Yang, and Youtao Zhang. A durable and energy efficient main memory using phase change memory technology. In *Proc. ACM International Symposium on Computer Architecture (ISCA)*, 2009.