

SimdQuickHeap: The QuickHeap Reconsidered

Johannes Breitling  

Karlsruhe Institute of Technology, Germany

Ragnar Groot Koerkamp  

Karlsruhe Institute of Technology, Germany

Marvin Williams  

Karlsruhe Institute of Technology, Germany

Abstract

Motivation. Priority queues are data structures that maintain a dynamic collection of elements and allow inserting new elements and removing the smallest element. The most widely known and used priority queue is likely the implicit binary heap, even though it has frequent cache misses and is hard to optimize using e.g. SIMD instructions.

Contributions. We introduce the SimdQuickHeap, a variant of the QuickHeap that was introduced by Navarro and Paredes in 2010. As suggested by the name, the data structure bears some similarity to QuickSort. We modify the data layout of the original QuickHeap to have all *pivots* adjacent in memory, with elements between consecutive pivots stored in dedicated *buckets*. This allows efficient SIMD implementations for both partitioning of buckets and scanning the list of pivots to find the bucket to append newly inserted elements to.

The SimdQuickHeap has amortized expected complexity $O(\log n)$ per operation, which improves to $O(\frac{1}{W} \log n)$ in non-degenerate cases, where W is the number of words in a SIMD register. In this case, the I/O-complexity is amortized $O(\frac{1}{B})$ per **push** and $O(\frac{1}{B} \log_2 \frac{n}{M})$ per **pop**.

Results. In synthetic benchmarks, the SimdQuickHeap is up to twice as fast as the next-best competitor, including the non-comparison radix heap, and needs around $1.5 \log_2 n$ comparisons and $\log_2 n$ nanoseconds per pair of push and pop operations. On graph benchmarks with Dijkstra's shortest path algorithm and Jarnik-Prim's minimum spanning tree algorithm, the SimdQuickHeap is consistently the fastest.

2012 ACM Subject Classification Theory of computation → Data structures design and analysis

Keywords and phrases Heap; SIMD; Priority Queue

Supplementary Material *Software:* github.com/RagnarGrootKoerkamp/quickheap

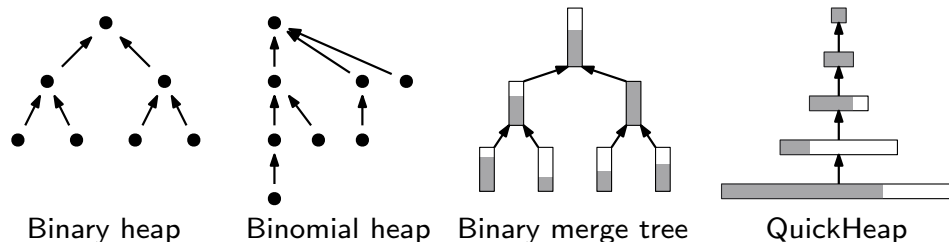
Acknowledgements We thank Peter Sanders for discussions related to this work and feedback on a draft of this paper.

1 Introduction

The Priority Queue (PQ) is a fundamental data structure that manages a collection of elements with associated priorities, allowing the insertion of new elements and the extraction of the element with the highest priority. PQs are central to a wide range of algorithms that rely on dynamically reordering and prioritizing tasks, such as perhaps most famously Dijkstra’s shortest-path algorithm [26], but also job scheduling, discrete event simulation, and branch-and-bound searches. As a result, PQs have been studied extensively, resulting in a large variety of priority queue designs (see Section 3).

Despite being among the oldest designs, (implicit) binary heaps are probably still the most widely used priority queues in practice¹. There are compelling reasons for their prevalence: they are conceptually simple, easy to implement, and have $\Theta(\log_2 n)$ worst-case running time with small constants for both insertions and deletions. While many theoretically superior designs have been proposed, most of them are prohibitively complex for practical use, and surprisingly, only a few priority queue designs that are faster in practice have emerged.

One such design is the *QuickHeap* [43], introduced in 2010 by Navarro and Paredes. It uses a recursive partitioning scheme reminiscent of quicksort [36], hence the name². However, it has received little attention in the literature. Figure 1 shows a schematic overview of the QuickHeap among other priority queue designs. We show that, due to its partitioning scheme, the QuickHeap lends itself to the SIMD (*single instruction, multiple data*) paradigm much more naturally than binary heaps and most other heaps based on merging trees or streams of elements. SIMD instructions are steadily becoming faster, more efficient, more powerful, and are capable of operating on wider registers. In light of these advancements, we argue that the QuickHeap should be reconsidered.



■ **Figure 1** Overview of some heaps. Black dots indicate single elements, and arrows point from larger to smaller elements. Vertical blocks indicate a sorted list of elements, while horizontal blocks indicate an unsorted list of elements. The QuickHeap stores pivots separating the buckets.

Contributions. We present *SimdQuickHeap*, an adaptation of the QuickHeap that achieves significantly improved performance through optimized data layout and SIMD instructions. To our knowledge, it is the first practical SIMD-optimized priority queue implementation.

The *push* and *pop* operations have an amortized worst-case run time bound of $O(\log n)$ and $O(1)$, respectively. In most practical cases, the SIMD implementation speeds up both operations by a factor of the SIMD register width W . Like the QuickHeap, the

¹ They are available in the standard libraries of many programming languages, including C++, Java, C#, Go, Rust, and Python.

² Although not a heap in the tree-like sense, the layers between the pivots (Figure 1) do form a linear tree with unordered sets of elements as nodes.

SimdQuickHeap is cache-oblivious with an amortized I/O-complexity of amortized $O(\frac{1}{B} \log \frac{n}{M})$ per operation for a main memory of M words and block size B , assuming $M = \Omega(\log n)$.

On synthetic benchmarks, SimdQuickHeap is up to $2\times$ faster than state-of-the-art implementations including the radix heap and superscalar sample queue [57]. On real-world instances of Dijkstra’s shortest path algorithm [26] and Jarník-Prim’s minimum-spanning-tree algorithm [38, 47], it is consistently the fastest method. Notably, while the run time of most tree-based and pointer-based priority queues grows relative to the $\Omega(\log_2 n)$ lower bound as n increases, SimdQuickHeap’s run time improves, falling below $\log_2 n$ nanoseconds per push-pop pair.

Notation and Definitions. Given a universe of elements \mathcal{U} and a strict weak ordering $<$ on $\mathcal{U} \cup \{-\infty, \infty\}$, the priority queue (PQ) data structure is a dynamic collection of elements \mathcal{Q} over \mathcal{U} (with duplicates permitted) that supports the following two operations:

push(e) Insert an element $e \in \mathcal{U}$ into \mathcal{Q} .

pop() Remove and return a minimal element $e \in \mathcal{Q}$, i.e., no $e' \in \mathcal{Q}$ with $e' < e$ exists.

For convenience, we assume **top** (pop without removing) and **size** (returns $|\mathcal{Q}|$) to also be available. We deliberately do not consider a **decreaseKey** operation (that replace arbitrary elements in the PQ by smaller ones), since it usually requires stable references to elements, adding significant implementation complexity and runtime overhead. While our definition admits distinct, equivalent elements, we refer to all equivalent elements as *equal* elements for simplicity.

2 The Original QuickHeap

In 2006, Paredes and Navarro presented *incremental quick select* [45], an algorithm that takes a list of n elements and then repeatedly returns the next-smallest element, such that returning the k smallest elements in an online setting takes $\mathcal{O}(n + k \log k)$ time. The authors note that their construction can be generalized to a proper priority queue. Consequently, in follow-up work from 2010, they present the QuickHeap [43]. The QuickHeap organizes the elements in disjoint *buckets* B_1, \dots, B_ℓ that are separated by *pivots* $p_1 > \dots > p_{\ell-1}$ such that

$$\infty \succ B_1 \succeq p_1 \succ B_2 \succeq p_2 \succ \dots \succeq p_{\ell-1} \succ B_\ell \succ -\infty,$$

where \succeq (\succ) indicates that the \geq ($>$) inequality holds for all elements in the sets. For the sake of simplicity, in the following we assume that the PQ contains no two equal elements. See Section 4.1 for details on how equal elements are handled in our design. The **pop** operation picks a new pivot p_ℓ (e.g., uniformly at random) from the *bottom* bucket B_ℓ and moves all element in B_ℓ smaller than p_ℓ into a new bucket $B_{\ell+1}$. If the new bucket is empty, it is discarded and a new pivot is selected. This process is repeated recursively in $B_{\ell+1}$ until the bottom bucket contains only one (the smallest) element. Finally, this element is returned and its bucket is removed. The **push** operation simply inserts the element into the correct bucket according to the pivots.

The original implementation [43] stores all buckets and pivots in an interleaved fashion in a flat array, with an additional array to track the positions of the pivots. Splitting a bucket is implemented analogously to partitioning in quicksort. To make room for a new element in a bucket B_i , all buckets with index smaller than i are shifted one slot to the right by moving two elements per bucket. The data structure is stored in a circular buffer, so that removing

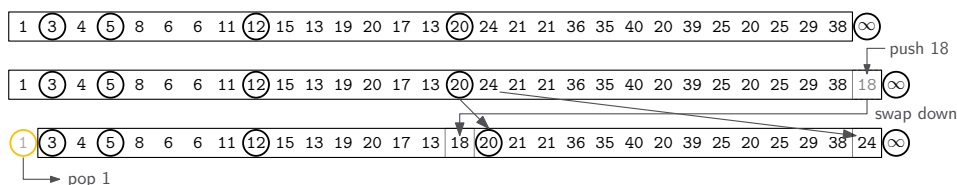


Figure 2 Overview of the QuickHeap data structure [43]. The data is stored in a single buffer. Pivots are circled and partition the data into smaller (left) and larger (right) elements. Their positions are stored in a separate list (not shown). An element is pushed by appending it at the end and then sifting it down until it is less than the corresponding pivot. An element is popped by repeatedly partitioning the bottom layer and then removing the smallest element.

the smallest element only requires incrementing a pointer. Figure 2 shows the data layout and illustrates the **push** and **pop** operations.

The **push** and **pop** operations of the QuickHeap are both in $\mathcal{O}(\log n)$ *in expectation over uniformly distributed input* [43]: the original QuickHeap can degenerate under specific input distributions. For example, when the inserted elements are decreasing, the number of buckets/pivots grows linearly over time, causing a run time of $\mathcal{O}(n)$ for insertions. Follow-up work introduces *Randomized QuickHeaps* [44], which addresses this issue by probabilistic repartitioning. Specifically, each time a bucket of size s is shifted during an insertion, with probability $1/s$, the Randomized QuickHeap unites bucket with all newer buckets by dropping the according pivots. While the Randomized QuickHeap achieves $\mathcal{O}(\log n)$ operations in expectation for *any* input, it adds significant overhead for benign cases. Very recently, Brodal et al. [16] introduced additional rebalancing schemes.

3 Related work

There is a vast body of literature on priority queues that dates back to the 1950s. Here, we highlight work on practical priority queues and focus on RAM-model complexity, I/O-complexity, the number of comparisons, and the running time. For a broader overview, we refer to the survey by Brodal [13] from 2013. A survey on the performance of classic priority queues is given by Larkin et al. [41]. The QuickHeap does not appear in either of these surveys, highlighting the little attention it has received.

Complexity lower bounds. There is a strong connection between sorting and priority queues [54]. Due to the lower bound for comparison based sorting, n insertions followed by n deletions must take at least $n \log_2 n$ comparisons for a comparison-based priority queue. A **pop** followed by a **push** requires at least $\log_2 n$ comparisons in the worst case amortized over many operations [18].

Binary and d -ary heap. A popular priority queue implementation is the binary heap using an implicit array representation [61]. The binary heap organizes the elements as a binary tree, maintaining the *heap invariant* that no element can be larger than its parent in the tree (unless it is the root). Thus, no element can be smaller than the root element. A binary heap can be implemented efficiently using the *Eytzinger* or *heap layout*, where the elements are stored consecutively in a flat array level by level starting from the root. Consequently, the children of node i are stored at positions $2i$ and $2i + 1$ in the array (1-based indexing). Both the **push** and **pop** operations take $\mathcal{O}(\log_2(n))$ time.

The d -ary heap generalizes the binary heap so that each node has up to d children. The main benefits are better cache locality and fewer levels to traverse. This comes at the cost of needing to find the minimum of d elements rather than just two on each level when deleting an element.

Amortized heaps A number of heaps reduce the $\mathcal{O}(\log n)$ cost of binary heap insertions. Binomial heaps (Figure 1) [58, 19] have amortized $\mathcal{O}(1)$ insertions. Fibonacci heaps [33] and pairing heaps [32] reduce this to truly constant time insertions, but have only amortized constant time deletions. The Brodal heap [11] (1996) has $\mathcal{O}(1)$ insertions and non-amortized $\mathcal{O}(\log n)$ deletions, but has a large hidden constant. Other non-amortized structures are the cascade heap [6], which has $\mathcal{O}(1)$ insertions and $\mathcal{O}(\log^* n + \log k)$ for the k 'th delete, and the logarithmic funnel [42], which has $\mathcal{O}(\log^* n)$ insertions and $\mathcal{O}(\log n)$ deletions.

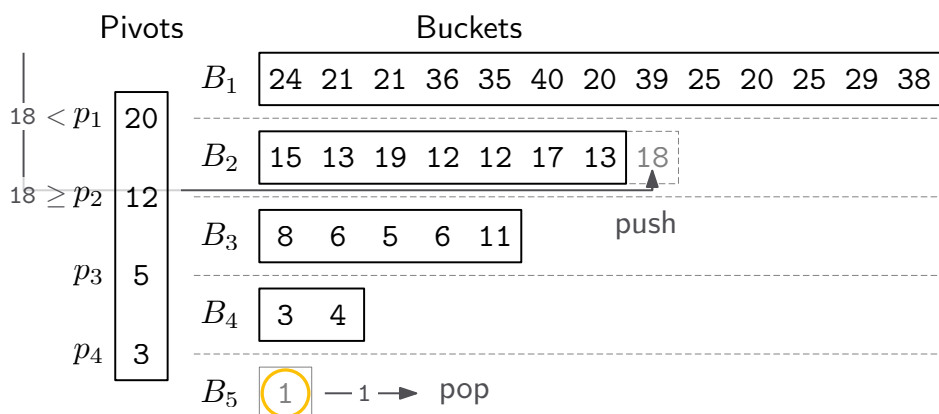
Unfortunately, all these structures are pointer-based and inefficient in practice. The weak heap and weak queue are more practical variants of these ideas that can be implemented using a flat array [27, 20, 28, 29].

decreaseKey. Some of the priority queues mentioned above also support `decreaseKey`, which lowers the key associated with a given element. Typically (but not always [21]), data structures supporting this require *pointer stability* of the elements and thus an additional level of indirection, hurting performance. We do not consider these further, but note that in the pointer-based model there is a lower bound of $\Omega(\log \log n)$ for `decreaseKey` [37] that is attained by slim and smooth heaps [51, 52].

I/O-model. Given that optimal priority queues for the RAM-model were already found early on, further research went into obtaining a good I/O-complexity [1]. In this setting, we are given an M -word internal memory and must minimize the number of transfers of B -word blocks to/from external memory. The lower bound for priority queues matches that of sorting, $\Omega(\frac{1}{B} \log_{M/B} \frac{n}{B})$ amortized transfers per operation [1], and a data structure matching this is given by Jiang and Larsen [17]. Practical implementations of I/O-optimal priority queues are the array heap [10] and sequence-heap [48], which improves the constant factor. The QuickHeap has slightly higher I/O cost of $\mathcal{O}((1/B) \log_2(n/M))$ for the `push` and `pop` operations [43]. With additional support for `decreaseKey`, the lower bound increases to $\Omega(\frac{1}{B} \log B / \log \log n)$ [30], and $\mathcal{O}(\frac{1}{B} \log \frac{n}{B} / \log \log n)$ is achieved by [39]. A *cache-oblivious* [34] priority queue with optimal amortized complexity that does not rely on the value of M and B is given by [3, 4], with improvement to $\mathcal{O}(1/B)$ insertions in [15]. The buffer heap [22, 21] supports `decreasekey` in $\mathcal{O}(\frac{1}{B} \log_2 \frac{n}{M})$ transfers and was discovered in parallel with the bucket heap [14]. The funnel heap [12] is another optimal cache-oblivious queue that relies on *binary mergers* (Figure 1) with strategically chosen buffer sizes.

Other priority queues. There are also data structures that exploit integer keys, that typically are *monotone* and do not allow decreasing keys. See [23] for a survey. The radix heap [2] inserts keys into a bucket based on the largest bit in which they differ from the current minimum, and has amortized complexity $\mathcal{O}(\log C)$ per operation when inserted keys are at most C larger than the current minimum. Like the QuickHeap, bucket sizes grow roughly exponentially. Further examples are bucket queues [25, 24] (similar to bucket sort) and $\mathcal{O}(\log w)$ dynamic predecessor structures for w -bit keys like Van Emde Boas trees [55], Y-fast tries [60], and fusion trees [31].

Practical/engineered priority queues. Given the plethora of theoretical results, there are surprisingly few engineered implementations. See [41] for a 2014 survey of the performance



■ **Figure 3** A schematic overview of the SimdQuickHeap: each bucket is a separate array, and the pivots are stored in their own array in sorted order.

of the classic data structures. Sequence heaps (2000) [48] provide an optimized I/O-efficient implementation based on merging sorted lists, with the drawback that elements are eagerly sorted, even when they are never removed. The QuickHeap (2010) [43] is indeed much faster when only few elements are removed, and competitive on synthetic benchmarks, but in theory is less I/O-efficient than the sequence heap. The superscalar sample queue (S³Q, 2021) [57] is a much more recent engineered data structure based on super scalar sample sort [50], that uses k -way partitioning instead of k -way merging to reach the optimal I/O-complexity. In particular, partitioning tends to be easier to optimize than merging, even though k -way splitting is harder to do evenly than data-independent k -way merging.

One further engineered data structure is the B-heap [40], which implements an implicit binary heap using a recursive layout like the Van Emde Boas tree [55].

Partition-based heaps. Most methods above are based on variously shaped trees with the heap property. Some methods, like the QuickHeap, instead use a path-shaped (linear) heap of buckets, where pivots (implicitly) partition the data into buckets. The priority queue of [3, 4] uses exponentially growing layers where each layer contains an increasing number of sorted buffers of increasing size. The buffer heap [22, 21] is conceptually simpler in that each layer contains an unsorted list of elements. The superscalar sample queue mixes these ideas, and each level contains k unsorted buffers. Brodal et al. [16] recently revisited this concept and provide *rebalancing strategies* to ensure the number of layers is guaranteed to be logarithmic.

4 The new SimdQuickHeap

We start by describing the basic scalar design of the *SimdQuickHeap* based on the original QuickHeap (see Section 2) and analyze its run time. We then discuss practical improvements using SIMD.

4.1 Scalar Design

Our variant of the QuickHeap keeps an explicit decreasing array of all pivots, as shown on the left in Figure 3. Second, the buckets are stored in individual buffers, rather than as a single allocation. The push operation classifies each element using a binary search on the

■ **Algorithm 1** Pseudo code for the `SimdQuickHeap` data structure. The `pop` operation assumes that it is not empty.

```

structure SIMDQUICKHEAP⟨T⟩
  pivots : Vec⟨T⟩, a vector of decreasing pivots
  buckets : Vec⟨Vec⟨T⟩⟩, the vector of bucket vectors
  ℓ: the number of buckets

1: procedure PUSH(x)
2:   i ← CLASSIFY(x, pivots)                                ▷ Binary search the first pivot ≤ x.
3:   buckets[i].PUSHBACK(x)

1: procedure POP
2:   while buckets[ℓ].LEN() > 1 do                            ▷ Partition the bottom layer repeatedly.
3:     i ← SELECTPIVOT(ℓ)
4:     p ← buckets[ℓ][i]
5:     B ← PARTITION(i, ℓ)
6:     if B ≠ ∅ then
7:       buckets.PUSHBACK(B)
8:       pivots.PUSHBACK(p)
9:     ℓ ← ℓ + 1
10:  x ← buckets[ℓ].POPBACK()
11:  ℓ ← ℓ - 1
12:  return x

1: procedure PARTITION(i, b)                                    ▷ Partition bucket b in-place.
2:   p ← buckets[b][i]                                        ▷ Returns elements less than p.
3:   k ← 1
4:   B ← Vec⟨T⟩()
5:   for j ∈ [1, buckets[b].LEN()) do
6:     if buckets[b][j] < p or (j > i and buckets[b][j] ≤ p) then
7:       B.PUSHBACK(buckets[b][j])
8:     else
9:       buckets[b][k] ← buckets[b][j]
10:    k ← k + 1
11:  buckets[b].RESIZE(k)
12:  return B

```

pivot array and then appends it to the respective bucket. To partition the bottom bucket for the `pop` operation, a *pivot* is selected and all elements smaller than the pivot are moved to a newly allocated bucket, while compactifying the remaining elements in-place, as shown in Algorithm 1. In case we selected the smallest element in the bucket as a pivot, no element moves to the new bucket, so we pick a new pivot and start over. However, if all elements in the bucket are equal, this does not terminate. For this reason, we treat elements equal to the pivot specially.

Equal Elements. If all elements in a bucket are equal, we need to ensure that partitioning terminates and does not degenerate to splitting off one element at a time. To achieve this, elements in the bucket that are on the left of the pivot are moved down when they are *strictly* smaller than the pivot, while elements on the right of the pivot are moved down when they are smaller *or equal*. This way, in expectation half the (non-pivot) elements are moved to

the new bucket. Since the pivot itself remains in its bucket, we maintain the invariant that no bucket is ever empty, unless the data structure is completely empty, in which case there is only one empty bucket. Alternatively, one could use dedicated buckets for elements equal to the pivot, but we did not pursue this strategy further. Note that our strategy allows for elements equal to the pivot p_i to be in buckets B_i and B_{i+1} .

Pivot Selection. Similar to quicksort, the performance of the `SimdQuickHeap` depends on the quality of the pivots. The median would be the ideal pivot, as it splits the bucket into two equal parts. Experimentally, the median of three random samples offered good performance across all tested workloads. Stronger strategies, such as the true median or median-of-medians [7], did not justify their additional overhead.

Rebalancing. The `SimdQuickHeap` does not suffer from degenerate inputs as much as the original `QuickHeap` (see Section 2): While the original `QuickHeap` “bubbles down” elements to insert them, the `SimdQuickHeap` uses binary search to classify elements, which takes $\mathcal{O}(\log n)$ time no matter the number of pivots. Thus, we do not employ any rebalancing strategies in our implementation. However, we will consider rebalancing in future work, since it still has many potential benefits: When there are $\mathcal{O}(\log n)$ buckets, the binary search has complexity $\mathcal{O}(\log \log n)$ [16]. Alternatively, this allows for a simple linear scan in $\mathcal{O}(\log n)$.

I/O-complexity. When the number of layers is indeed $\mathcal{O}(\log n)$, the `SimdQuickHeap` has the same amortized I/O-complexity as the `QuickHeap`: $\mathcal{O}(\frac{1}{B} \log_2 \frac{n}{M})$ per operation, assuming that the first size- B block of each bucket fits in the cache, i.e., $M = \Omega(B \log n)$. Specifically, classifying a pushed element then does not require any memory transfers since the pivots are in memory, and appending it to a block incurs amortized $\mathcal{O}(\frac{1}{B})$ transfers by buffering the last block of elements appended to each bucket. The memory access patterns for the partitioning step are similar to the original `QuickHeap`, and the original analysis applies.

4.2 Amortized analysis

For simplicity, we also assume that the true median is used as a pivot (as this can be found in linear time), and that no two elements are equal. Using the median of three random elements yields the same bounds for `push` and, in expectation over the random pivot selection, also for `pop`. This can be shown with a similar, more technical, proof that uses the fact that median of three random elements gives a constant-factor split in expectation.

► **Theorem 1.** *For the `SimdQuickHeap` containing n elements, the `push` operation takes $\mathcal{O}(\log n)$ amortized and worst-case time, the `pop` operation takes $\mathcal{O}(1)$ amortized time.*

Proof. The `push` operation can classify the element to insert in $\mathcal{O}(\log n)$ time via a binary search on the pivots, since there are at most $\ell \leq n - 1$ pivots. Appending the element to the correct bucket takes constant time, resulting in a total worst-case run time of $\mathcal{O}(\log n)$.

For the amortized run time analysis, we define the *potential function*

$$\Phi := c \cdot \sum_{i=1}^{\ell} (|B_i| + 1) \ln(|B_i| + 1)$$

for some sufficiently large constant c . Let $k_i = |B_i| + 1$ be the *weight* of the bucket B_i . Then,

inserting an element into bucket B_i increases the potential by

$$\begin{aligned}\Delta_{\text{push},i}\Phi &= c(k_i + 1) \ln(k_i + 1) - ck_i \ln(k_i) \\ &= ck_i \ln(1 + k_i^{-1}) + c \ln(k_i + 1) \\ &\leq c(1 + \ln(k_i + 1)) && \text{using } \ln(1 + x) \leq x \\ &\leq 2c \ln(n + 2) \in \mathcal{O}(\log n).\end{aligned}$$

Thus, the **push** operation takes $\mathcal{O}(\log n)$ amortized time.

The **pop** operation first finds the median of the bottom bucket of size k_ℓ as pivot and partitions it, resulting in two buckets of equal size $k_\ell/2$. This takes $\mathcal{O}(k_\ell)$ time and changes the potential (ignoring rounding errors) by

$$\Delta_{\text{part}}\Phi = 2 \left(c \frac{k_\ell}{2} \ln \left(\frac{k_\ell}{2} \right) \right) - ck_\ell \ln(k_\ell) = -ck_\ell \ln(2).$$

Repeatedly partitioning until the bottom bucket contains only one element changes the potential (again, ignoring rounding errors) by $\Delta_{\text{pop}}\Phi = -c \ln(2)(k_\ell + \frac{k_\ell}{2} + \frac{k_\ell}{4} + \dots + 2) = -\Omega(ck_\ell)$ and takes $\mathcal{O}(k_\ell)$ time. Thus, with c sufficiently large, the potential can “pay” for the $\mathcal{O}(k_\ell)$ actual cost of the pop operation, giving $\mathcal{O}(1)$ amortized time. ◀

4.3 Practical Optimizations

Here, we describe practical optimizations that take advantage of the data layout of the SimdQuickHeap. Let W be the number of elements that fit into one SIMD register. When the number of layers is $\mathcal{O}(\log n)$ indeed, together, the two optimizations below result in an $\mathcal{O}(\frac{1}{W} \log n)$ algorithm.

Classification. If the number of pivots is small, a linear scan can be faster than a binary search, thanks to better cache efficiency and branch predictions. We therefore employ a SIMD-optimized linear scan unless the number of pivots surpasses some threshold, where we fall back to binary search. In practice, the threshold is large enough such that the linear scan is used for all feasible inputs where the number of buckets is in $\mathcal{O}(\log n)$. We compare W consecutive pivots at once with SIMD, resulting in a run time of $\mathcal{O}(\frac{1}{W} \log n)$ when the number of pivots is in $\mathcal{O}(\log n)$.

Partitioning. Multiple partitioning schemes using SIMD have been proposed [9, 59]. We use a straightforward mechanism: We load W values from the bucket and compare them against the pivot at once in $\mathcal{O}(1)$ time. Using AVX2 `permutevar`³ (~ 3 cycles) or AVX-512 `compress` (~ 6 cycles) instructions, we can then efficiently move the values smaller (or larger) than the pivot to the start of the register. Finally, we append these values to the back of the bucket array. This way, partitioning a bucket of size m takes $\mathcal{O}(m/W)$ time.

In practice, the constant overheads of this partitioning scheme become too large when the bucket is small (of length $\mathcal{O}(W)$). Thus, we stop the iterative partitioning when the bucket size falls below length 16. To be able to delete the minimum in the last bucket efficiently, we keep it sorted in descending order as long as it is smaller than the threshold. Consequently, when a new element is inserted into the last bucket and the bucket size is below the threshold, the element is inserted at the appropriate position using a linear scan.

³ See <https://lemire.me/blog/2017/04/10/removing-duplicates-from-lists-quickly/>.

5 Experiments

The implementation of the `SimdQuickHeap` and the evaluations are written in Rust and are available at [github:ragnargrootkoerkamp/quickheap](https://github.com:ragnargrootkoerkamp/quickheap). We use an AMD Zen 4 EPYC 9684X (92 cores, 64 KiB L1 cache, 1 MiB L2 cache, and 32+64 MiB L3 cache for each package of 8 cores) with a 12-channel DDR5 RAM running Rocky Linux 9.4 for the experiments. All experiments are run using a single thread and repeated three times after one warm-up iteration. We report the median of these runs but note that the measurements had generally very little variance. We enable full link-time optimization to minimize the impact of cross-language function calls. We further reimplemented some of the benchmarks in C++ for the C++ implementations, but measured no difference in performance.

Compared implementations. For `SimdQuickHeap`, we test both 256-bit AVX2 and 512-bit AVX-512 implementations. We compare against a number of external implementations:

- binary heap, the Rust standard library `std::collections::BinaryHeap`.
- 8-ary heap, implementation from the `orx_priority_queue` Rust crate [5]. This is the fastest d -ary heap implementation we found, faster than e.g. the implementation in Boost [8]. The 4-ary variant was slightly slower for large inputs.
- radix heap from the `radix_heap` Rust crate [46].
- weak heap from the `weakheap` Rust crate [53].
- sequence heap, the original C++ implementation [48, 49], via Rust bindings.
- S3Q, the original superscalar sample queue C++ implementation [57, 56], via Rust bindings.

We re-implemented the original `QuickHeap`, since the source code was not available to us. We further implemented a `ScalarQuickHeap` variant of the `SimdQuickHeap` that does not use SIMD instructions and can be instrumented to count internal metrics such as the number of comparisons.

We tested further implementations in preliminary experiments but excluded them because they were consistently slower than the ones listed above. In particular, we found no competitive implementation of theoretically efficient pointer-based structures like Fibonacci heaps and pairing heaps.

5.1 Synthetic benchmarks

Workloads. We tested each implementation on a set of different synthetic workloads, for both 32-bit and 64-bit keys. We vary n , the maximum size of the the heap, from $2^{10} = 1024$ to $2^{25} \approx 32$ million.

The `HeapSort` workload first generates n random values and then measures how long it takes to first push all of them and then pop all of them: $\text{push}^n \circ \text{pop}^n$. In the `Wiggle` workload, we grow the structure to size n via a sequence of $3n$ operations $(\text{push} \circ \text{pop} \circ \text{push})^n$ and then empty it via $(\text{pop} \circ \text{push} \circ \text{pop})^n$. The `ConstantSize` workload is initialized by growing the data structure to contain n elements via $(\text{push} \circ \text{pop} \circ \text{push})^n$ as well. Then, we measure how long it takes to do $10n$ pairs of $\text{pop} \circ \text{push}$ operations.

In `MonotoneWiggle` and `MonotoneConstantSize`, the pushed value is chosen as the last popped value plus a uniform random constant between 0 and n . In the non-monotone `Wiggle`, the pushed value is a uniform random (32 or 64-bit) integer. We note that this non-monotone case is degenerate, in that most values are pushed to the front of the queue.

Metrics. We normalize the total time for each workload by the number of `pop` \circ `push` pairs and by $\log_2 n$, resulting in the number of nanoseconds needed per *fundamental operation*, since there is a lower-bound of $\log_2 n$ comparisons for a `pop` \circ `push` pair when the queue contains n elements. Separately, we also measure the total number of comparisons made by each algorithm, as an indication for how closely they approach the lower bound of $\log_2 n$ comparisons per `pop` \circ `push`.

Results. We see in Figure 4 that the binary heap and d -ary heap slow down significantly as the data grows beyond the size of the CPU caches. The d -ary heap is consistently around $1.4\times$ faster than the binary heap. The weak heap is surprisingly competitive to the binary heap.

The radix heap benefits from high cache-locality and gets faster rather than slower as n grows. It is faster for 32-bit than for 64-bit data. Similarly, as predicted by the theory, the sequence heap is I/O-efficient and does not slow down on larger inputs. The more modern and optimized superscalar-sample-queue is around $2\times$ faster for 64-bit data, very close to the radix heap.

Our reimplementation of the original QuickHeap is on-par with the sequence heap, while the scalar version of our new heap is $2\times$ slower, possibly because the original QuickHeap uses more efficient in-place partitioning. The SimdQuickHeap is around $8\times$ faster than the scalar version. Like the other engineered heaps, it gets faster relative to the lower bound as n increases, indicating that the constant overhead of each `push` and `pop` is relatively large compared to the $\log_2 n$ pivot steps required for each element. The AVX-512 version is up to $1.2 - 1.4\times$ faster than the AVX2 version, and around $2\times$ faster than the superscalar-sample-queue. Surprisingly, it also significantly outperforms the radix heap in all tested cases.

The SimdQuickHeap reaches below $\log_2 n$ nanoseconds per `pop` \circ `push` pair. This means that it requires around 1 nanosecond for each comparison in the lower bound. Since each nanosecond corresponds to 3.7 clock cycles, each of which can contain multiple SIMD instructions on 8 words, the overhead of moving data around is still very large compared to just the comparisons.

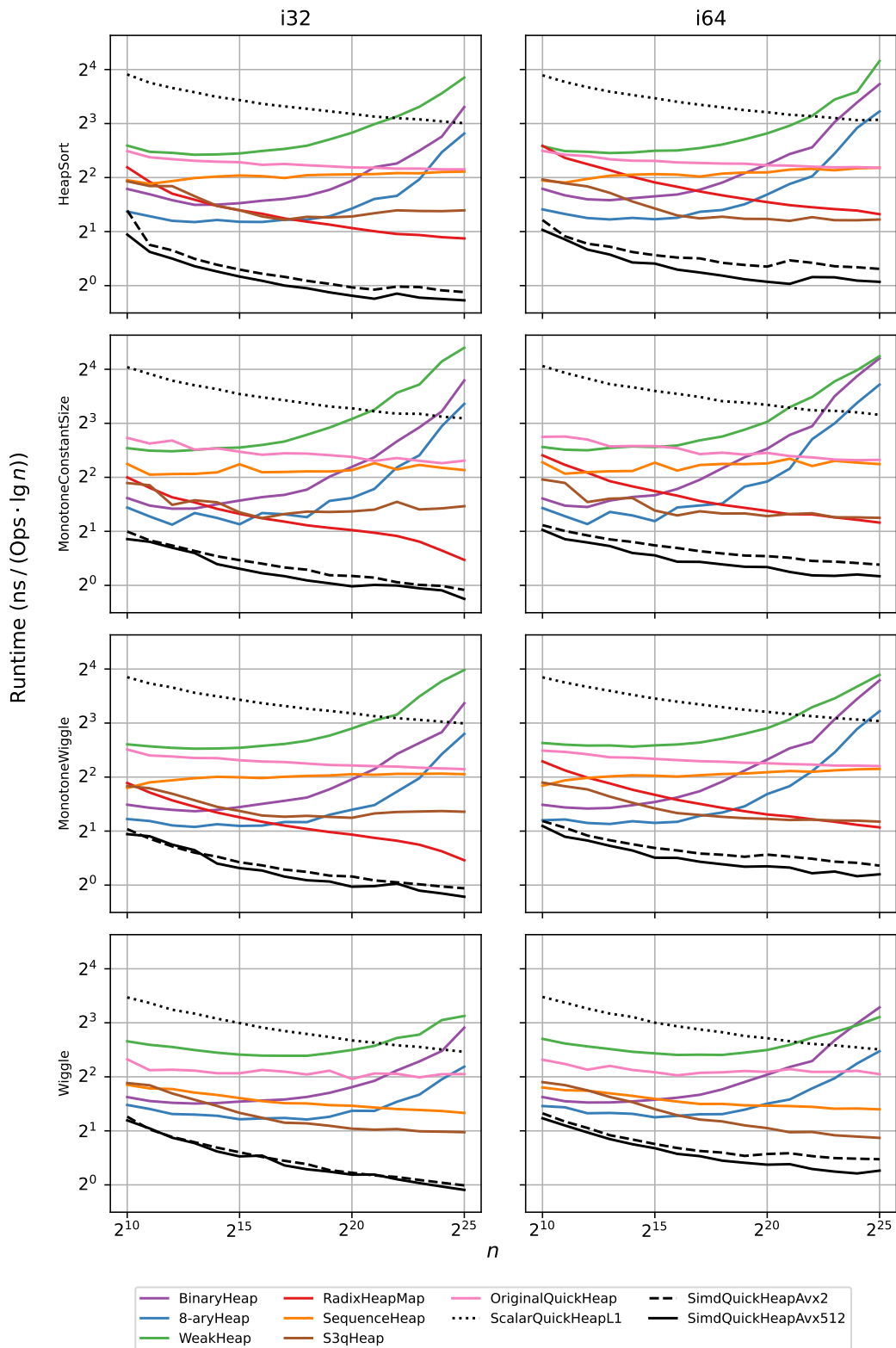
Number of comparisons. On the ConstantSize workload (Figure 6), the binary heap makes less than $1.1 \log_2 n$ comparisons per `pop` \circ `push`, while the 8-ary heap takes around $2.6 \log_2 n$ comparisons for this. The binary heap needs more comparisons on degenerate input, while the weak heap always consistently needs $1.0 \log_2 n$.

The ScalarQuickHeap makes $1.7 \log_2 n$ comparisons when using random pivots in combination with a linear scan over the pivots to insert elements. This decreases to $1.1 \log_2 n$ when using an oracle that gives the exact median for free, showing that there is some room for improvement by using a more accurate partitioning scheme.

Running time distribution. When running SimdQuickHeap with $n = 2^{25}$ on the ConstantSize workload, 21% of the time is spent pushing elements, with 14% (of the total) scanning the list of pivots to find the right layer to append to. The other 78% of time is spent on `pop`: 59% partitioning the input, 11% finding the position of the smallest element in the bottom bucket, and 3% removing and returning this element.

5.2 Graph benchmarks

Two common graph textbook algorithms that require a priority queue are Dijkstra’s algorithm [26] to compute shortest paths on a weighted directed graph with non-negative



■ **Figure 4** Log-log plots showing for a variety of workloads (rows, see Section 5.1), i32 and i64 inputs (columns), and increasing maximum number of elements stored in the data structure (x-axis) the time in nanoseconds, divided by $\text{Ops} \cdot \log_2 n$, i.e., the time relative to the $O(\log_2 n)$ lower bound per operation, where Ops is the number of pairs of push-pop operations

Instance	Description	$ V $ ($\times 10^6$)	$ E $ ($\times 10^6$)	max degree	median weight	max weight
CAL	Rd. California	1	4	8	3115	54k
CTR	Rd. Central USA	14	34	9	3467	54k
GER	Rd. Germany	20	41	9	39	62k
USA	Rd. USA	23	58	9	3473	92k
RHG ₂₀	RHG	1	10	90k	230	999
RHG ₂₂	RHG	4	41	941k	230	999
RHG ₂₄	RHG	16	159	595k	232	999

■ **Table 1** Number of edges and vertices of the graph instances that have been used for benchmarking, as well as the median and maximal edge weight.

edge weights, and Jarnik–Prim’s algorithm [47, 38] to compute minimum-spanning-trees on weighted, undirected graphs. In this section, we compare running times of the two algorithms on multiple graph instances using different priority queues.

Setup. We use 32-bit identifiers for vertices/edges and 32-bit non-negative edge weights. We pack these into a single 64-bit unsigned integer that we use as elements for the priority queues. Since most implementations do not support `decreaseKey` (including the `SimdQuickHeap`, we re-insert a vertex when a new shorter path is found. When popping a vertex, we check if the distance stored in the queue matches the best distance found so far, and if not, discard it.

Graph Types. We used different types of graphs to benchmark the performance of the `SimdQuickHeap` (Table 1). We test on four road networks of varying sizes with edge weights represent travel times, all but the German⁴ network downloaded from 9th DIMACS implementation challenge⁵. We also test on random hyperbolic graphs generated with `KaGen` [35] with 2^{20} , 2^{22} , and 2^{24} nodes, an average degree of 16 and a power law exponent of $\gamma = 2.3$. Edge weights represent the hyperbolic distance.

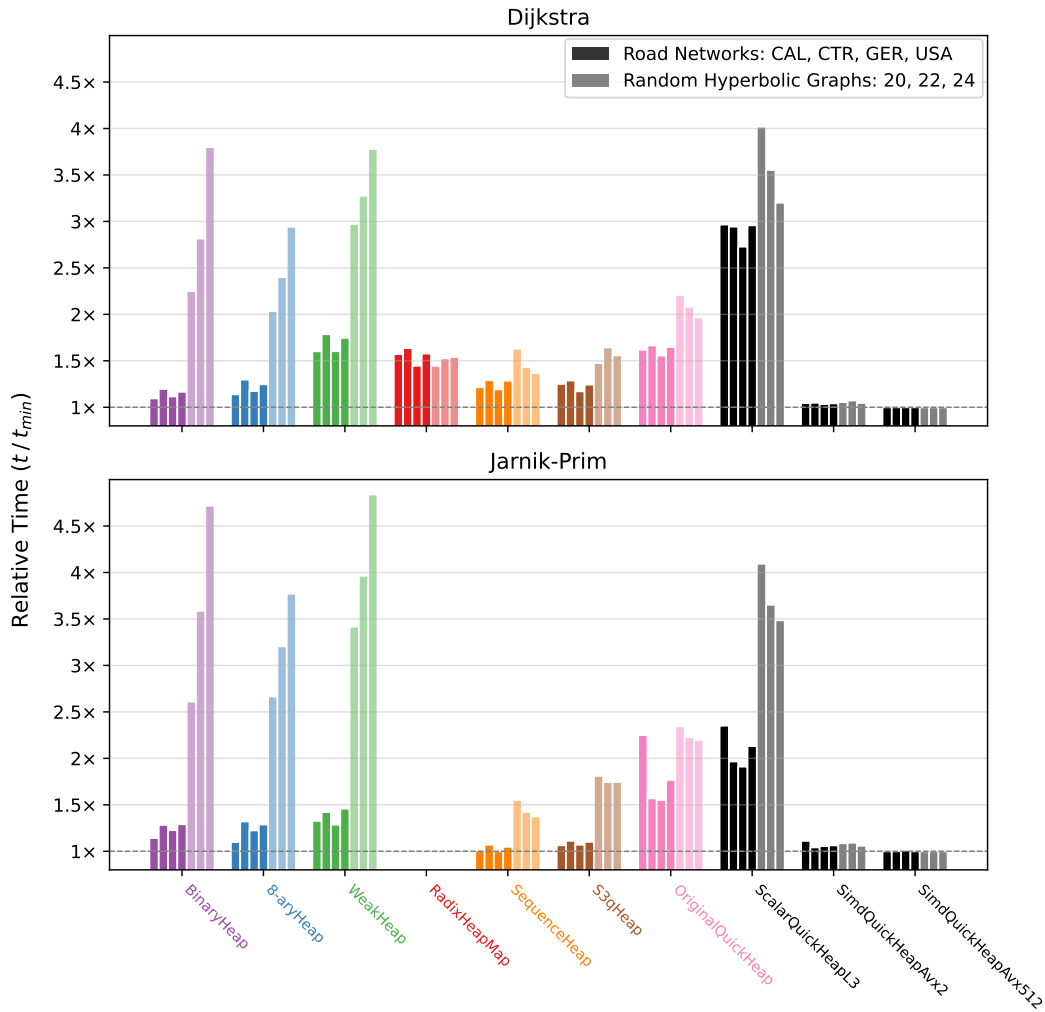
Results. In Figure 5, we see that the `SimdQuickHeap` gives the fastest running time for all graphs for both Dijkstra and Jarnik–Prim. Speedups on road networks are small since the cache misses involved with traversing the graph are likely the bottleneck. On the hyperbolic graphs, we achieve speedup factors of up to $3\times$ compared to the binary heap.

6 Conclusion

Due its conceptually simple design, the `SimdQuickHeap` allows an efficient implementation using SIMD instructions that has both a good theoretical complexity as well as a great performance in practice. As predicted by the improved I/O-complexity, the `SimdQuickHeap` vastly outperforms the binary and d -ary heap, as well as the weak heap. This new design gives a $4\times$ speedup over I/O-efficient structures such as the sequence heap (2000) [48] and QuickHeap (2010) [43], and a $2\times$ speedup over the much more recent superscalar sample queue [57]. Maybe more surprising, the `SimdQuickHeap` is also up to $2\times$ faster than the non-comparison-based radix heap.

⁴ <https://i11www.itk.kit.edu/resources/roadgraphs.php>

⁵ <http://www.diag.uniroma1.it/challenge9/download.shtml>



■ **Figure 5** Relative runtime of the different heaps on multiple graph instances with regard to the minimum runtime on this instance. The **bold** bars represent the road networks, and the light bars represent random hyperbolic graphs of different sizes.

Future Work. Future work will be to implement the recently introduced rebalancing strategies of [16] to limit the number of buckets, and to evaluate if and how much they affect performance in both (currently) degenerate and non-degenerate cases. This will also make the running time and I/O-complexity hold unconditionally. Furthermore, the current $O(\frac{1}{B} \log_2 \frac{n}{M})$ I/O-complexity is missing the more efficient $\log_{M/B}$ that is obtained by multi-way merging/splitting in the sequence heap and superscalar sample queue, and so the question is whether the QuickHeap naturally extends to multi-way splitting.

Another option is to make a SIMD-optimized version of the radix-heap, as it has a similar structure with exponentially growing buckets.

From the engineering side, optimizing the partitioning might provide further gains. A possible drawback of the SimdQuickHeap is the larger memory usage compared to the original QuickHeap. Using a list of reusable blocks rather than single vectors could reduce memory usage. Furthermore, the implementation could be expanded to support key-value pairs that are both 64 bits, bulk-insertions, and multi-threading.

References

- 1 Alok Aggarwal and S. Vitter, Jeffrey. The input/output complexity of sorting and related problems. *Communications of the ACM*, 31(9):1116–1127, September 1988. URL: <http://dx.doi.org/10.1145/48529.48535>, doi:10.1145/48529.48535.
- 2 Ravindra K. Ahuja, Kurt Mehlhorn, James Orlin, and Robert E. Tarjan. Faster algorithms for the shortest path problem. *Journal of the ACM*, 37(2):213–223, April 1990. URL: <http://dx.doi.org/10.1145/77600.77615>, doi:10.1145/77600.77615.
- 3 Lars Arge, Michael A. Bender, Erik D. Demaine, Bryan Holland-Minkley, and J. Ian Munro. Cache-oblivious priority queue and graph algorithm applications. In *Proceedings of the thirty-fourth annual ACM symposium on Theory of computing*, STOC02, page 268–276. ACM, May 2002. URL: <http://dx.doi.org/10.1145/509907.509950>, doi:10.1145/509907.509950.
- 4 Lars Arge, Michael A. Bender, Erik D. Demaine, Bryan Holland-Minkley, and J. Ian Munro. An optimal cache-oblivious priority queue and its application to graph algorithms. *SIAM Journal on Computing*, 36(6):1672–1695, January 2007. URL: <http://dx.doi.org/10.1137/S0097539703428324>, doi:10.1137/s0097539703428324.
- 5 Ugur Arikan. orx-priority-queue: Priority queue traits and efficient d-ary heap implementations. <https://github.com/orxfun/orx-priority-queue/>, 2023.
- 6 Maxim Babenko, Ignat Kolesnichenko, and Ivan Smirnov. Cascade heap: Towards time-optimal extractions. In *Computer Science – Theory and Applications*, page 62–70. Springer International Publishing, 2017. URL: http://dx.doi.org/10.1007/978-3-319-58747-9_8, doi:10.1007/978-3-319-58747-9_8.
- 7 Manuel Blum, Robert W. Floyd, Vaughan Pratt, Ronald L. Rivest, and Robert E. Tarjan. Time bounds for selection. *Journal of Computer and System Sciences*, 7(4):448–461, August 1973. doi:10.1016/S0022-0000(73)80033-9.
- 8 Boost. Boost C++ libraries. boost.org, 2026.
- 9 Berenger Bramas. A Novel Hybrid Quicksort Algorithm Vectorized using AVX-512 on Intel Skylake. *International Journal of Advanced Computer Science and Applications (ijacsa)*, 8(10), October 2017. doi:10.14569/IJACSA.2017.081044.
- 10 Klaus Brengel, Andreas Crauser, Paolo Ferragina, and Ulrich Meyer. An experimental study of priority queues in external memory. In *Algorithm Engineering*, page 345–359. Springer Berlin Heidelberg, 1999. URL: http://dx.doi.org/10.1007/3-540-48318-7_27, doi:10.1007/3-540-48318-7_27.
- 11 Gerth Stølting Brodal. Worst-case efficient priority queues. In *Proceedings of the Seventh Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '96, page 52–58, USA, 1996. Society for Industrial and Applied Mathematics.
- 12 Gerth Stølting Brodal and Rolf Fagerberg. Funnel heap - a cache oblivious priority queue. In *Proceedings of the 13th International Symposium on Algorithms and Computation*, ISAAC '02, page 219–228, Berlin, Heidelberg, 2002. Springer-Verlag.
- 13 Gerth Stølting Brodal. *A Survey on Priority Queues*, page 150–163. Springer Berlin Heidelberg, 2013. URL: http://dx.doi.org/10.1007/978-3-642-40273-9_11, doi:10.1007/978-3-642-40273-9_11.
- 14 Gerth Stølting Brodal, Rolf Fagerberg, Ulrich Meyer, and Norbert Zeh. Cache-oblivious data structures and algorithms for undirected breadth-first search and shortest paths. In *Algorithm Theory - SWAT 2004*, page 480–492. Springer Berlin Heidelberg, 2004. URL: http://dx.doi.org/10.1007/978-3-540-27810-8_41, doi:10.1007/978-3-540-27810-8_41.
- 15 Gerth Stølting Brodal, Michael T. Goodrich, John Iacono, Jared Lo, Ulrich Meyer, Victor Pagan, Nodari Sitchinava, and Rolf Svenning. External-memory priority queues with optimal insertions. volume 351, pages 5:1–5:14. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2025. URL: <https://drops.dagstuhl.de/entities/document/10.4230/LIPICs.ESA.2025.5>, doi:10.4230/LIPICs.ESA.2025.5.

- 16 Gerth Stølting Brodal, John Iacono, Casper Moldrup Rysgaard, and Sebastian Wild. Partition-based simple heaps. *arXiv*, 2026. URL: <https://arxiv.org/abs/2603.01206>, doi:10.48550/ARXIV.2603.01206.
- 17 Gerth Stølting Brodal and Jyrki Katajainen. Worst-case efficient external-memory priority queues. page 107–118, 1998. URL: <http://dx.doi.org/10.1007/BFb0054359>, doi:10.1007/bfb0054359.
- 18 Mark R. Brown. The complexity of priority queue maintenance. In *Proceedings of the ninth annual ACM symposium on Theory of computing - STOC '77*, STOC '77, page 42–48. ACM Press, 1977. URL: <http://dx.doi.org/10.1145/800105.803394>, doi:10.1145/800105.803394.
- 19 Mark R. Brown. Implementation and analysis of binomial queue algorithms. *SIAM Journal on Computing*, 7(3):298–319, August 1978. URL: <http://dx.doi.org/10.1137/0207026>, doi:10.1137/0207026.
- 20 Asger Bruun, Stefan Edelkamp, Jyrki Katajainen, and Jens Rasmussen. Policy-based benchmarking of weak heaps and their relatives,. In *Experimental Algorithms*, page 424–435. Springer Berlin Heidelberg, 2010. URL: http://dx.doi.org/10.1007/978-3-642-13193-6_36, doi:10.1007/978-3-642-13193-6_36.
- 21 Rezaul A. Chowdhury and Vijaya Ramachandran. Cache-oblivious buffer heap and cache-efficient computation of shortest paths in graphs. *ACM Transactions on Algorithms*, 14(1):1–33, January 2018. URL: <http://dx.doi.org/10.1145/3147172>, doi:10.1145/3147172.
- 22 Rezaul Alam Chowdhury and Vijaya Ramachandran. Cache-oblivious shortest paths in graphs using buffer heap. In *Proceedings of the sixteenth annual ACM symposium on Parallelism in algorithms and architectures*, SPAA04, page 245–254. ACM, June 2004. URL: <http://dx.doi.org/10.1145/1007912.1007949>, doi:10.1145/1007912.1007949.
- 23 Jonas Costa, Lucas Castro, and Rosiane de Freitas. Exploring monotone priority queues for dijkstra optimization. *RAIRO - Operations Research*, 59(5):2419–2436, September 2025. URL: <http://dx.doi.org/10.1051/ro/2025082>, doi:10.1051/ro/2025082.
- 24 Eric V. Denardo and Bennett L. Fox. Shortest-route methods: 1. reaching, pruning, and buckets. *Operations Research*, 27(1):161–186, February 1979. URL: <http://dx.doi.org/10.1287/opre.27.1.161>, doi:10.1287/opre.27.1.161.
- 25 Robert B. Dial. Algorithm 360: shortest-path forest with topological ordering [h]. *Communications of the ACM*, 12(11):632–633, November 1969. URL: <http://dx.doi.org/10.1145/363269.363610>, doi:10.1145/363269.363610.
- 26 Edsger W Dijkstra. A note on two problems in connexion with graphs. *Numerische Mathematik*, 1(1):269–271, Dec 1959. doi:10.1007/bf01386390.
- 27 Ronald D. Dutton. Weak-heap sort. *BIT*, 33(3):372–381, September 1993. URL: <http://dx.doi.org/10.1007/BF01990520>, doi:10.1007/bf01990520.
- 28 Stefan Edelkamp, Amr Elmasry, and Jyrki Katajainen. The weak-heap data structure: Variants and applications. *Journal of Discrete Algorithms*, 16:187–205, October 2012. URL: <http://dx.doi.org/10.1016/j.jda.2012.04.010>, doi:10.1016/j.jda.2012.04.010.
- 29 Stefan Edelkamp, Amr Elmasry, and Jyrki Katajainen. Weak heaps engineered. *Journal of Discrete Algorithms*, 23:83–97, November 2013. URL: <http://dx.doi.org/10.1016/j.jda.2013.07.002>, doi:10.1016/j.jda.2013.07.002.
- 30 Kasper Eenberg, Kasper Green Larsen, and Huacheng Yu. Decreasekeys are expensive for external memory priority queues. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, STOC '17, page 1081–1093. ACM, June 2017. URL: <http://dx.doi.org/10.1145/3055399.3055437>, doi:10.1145/3055399.3055437.
- 31 M. L. Fredman and D. E. Willard. Blasting through the information theoretic barrier with fusion trees. In *Proceedings of the twenty-second annual ACM symposium on Theory of computing - STOC '90*, STOC '90, page 1–7. ACM Press, 1990. URL: <http://dx.doi.org/10.1145/100216.100217>, doi:10.1145/100216.100217.

- 32 Michael L. Fredman, Robert Sedgwick, Daniel D. Sleator, and Robert E. Tarjan. The pairing heap: A new form of self-adjusting heap. *Algorithmica*, 1(1–4):111–129, November 1986. URL: <http://dx.doi.org/10.1007/BF01840439>, doi:10.1007/bf01840439.
- 33 Michael L. Fredman and Robert Endre Tarjan. Fibonacci heaps and their uses in improved network optimization algorithms. *Journal of the ACM*, 34(3):596–615, July 1987. URL: <http://dx.doi.org/10.1145/28869.28874>, doi:10.1145/28869.28874.
- 34 Matteo Frigo, Charles E. Leiserson, Harald Prokop, and Sridhar Ramachandran. Cache-oblivious algorithms. *ACM Transactions on Algorithms*, 8(1):1–22, January 2012. URL: <http://dx.doi.org/10.1145/2071379.2071383>, doi:10.1145/2071379.2071383.
- 35 Daniel Funke, Sebastian Lamm, Ulrich Meyer, Manuel Penschuck, Peter Sanders, Christian Schulz, Darren Strash, and Moritz von Looz. Communication-free massively distributed graph generation. *Journal of Parallel and Distributed Computing*, 131:200–217, September 2019. doi:10.1016/j.jpdc.2019.03.011.
- 36 C. A. R. Hoare. Algorithm 64: Quicksort. *Communications of the ACM*, 4(7):321, 1961. URL: <http://dx.doi.org/10.1145/366622.366644>, doi:10.1145/366622.366644.
- 37 John Iacono and Özgür Özkán. Why some heaps support constant-amortized-time decrease-key operations, and others do not. In *Automata, Languages, and Programming*, page 637–649. Springer Berlin Heidelberg, 2014. URL: http://dx.doi.org/10.1007/978-3-662-43948-7_53, doi:10.1007/978-3-662-43948-7_53.
- 38 Vojtěch Jarník. O jistém problému minimálním. (Z dopisu panu O. Borůvkovi). *Práce Moravské Přírodovědecké Společnosti*, 1930.
- 39 Shunhua Jiang and Kasper Green Larsen. A faster external memory priority queue with decreasekeys. *arXiv*, 2018. URL: <https://arxiv.org/abs/1806.07598>, doi:10.48550/ARXIV.1806.07598.
- 40 Poul-Henning Kamp. You’re doing it wrong. *Communications of the ACM*, 53(7):55–59, July 2010. URL: <http://dx.doi.org/10.1145/1785414.1785434>, doi:10.1145/1785414.1785434.
- 41 Daniel H. Larkin, Siddhartha Sen, and Robert E. Tarjan. A back-to-basics empirical study of priority queues. In *2014 Proceedings of the Sixteenth Workshop on Algorithm Engineering and Experiments (ALENEX)*, page 61–72. Society for Industrial and Applied Mathematics, December 2013. URL: <http://dx.doi.org/10.1137/1.9781611973198.7>, doi:10.1137/1.9781611973198.7.
- 42 Christian Loeffeld. The logarithmic funnel heap: An efficient priority queue for extremely large sets, 2023. URL: <https://arxiv.org/abs/1705.10648v4>, doi:10.48550/ARXIV.1705.10648.
- 43 Gonzalo Navarro and Rodrigo Paredes. On sorting, heaps, and minimum spanning trees. *Algorithmica*, 57(4):585–620, March 2010. URL: <http://dx.doi.org/10.1007/s00453-010-9400-6>, doi:10.1007/s00453-010-9400-6.
- 44 Gonzalo Navarro, Rodrigo Paredes, Patricio V. Poblete, and Peter Sanders. Stronger quickheaps. *International Journal of Foundations of Computer Science*, 22(04):945–969, June 2011. URL: <http://dx.doi.org/10.1142/S0129054111008507>, doi:10.1142/s0129054111008507.
- 45 Rodrigo Paredes and Gonzalo Navarro. Optimal incremental sorting. In *Proceedings of the Meeting on Algorithm Engineering & Experiments*, page 171–182, USA, 2006. Society for Industrial and Applied Mathematics.
- 46 Mike Pedersen. radix-heap: Fast monotone priority queues. <https://github.com/mpdn/radix-heap>, 2016.
- 47 R. C. Prim. Shortest connection networks and some generalizations. *The Bell System Technical Journal*, 36(6):1389–1401, 1957. doi:10.1002/j.1538-7305.1957.tb01515.x.
- 48 Peter Sanders. Fast priority queues for cached memory. *ACM Journal of Experimental Algorithmics*, 5:7, December 2000. URL: <http://dx.doi.org/10.1145/351827.384249>, doi:10.1145/351827.384249.
- 49 Peter Sanders. Sequence heap. <https://github.com/raphinesse/SequenceHeap>, 2000.

- 50 Peter Sanders and Sebastian Winkel. Super scalar sample sort. In *Algorithms – ESA 2004*, page 784–796. Springer Berlin Heidelberg, 2004. URL: http://dx.doi.org/10.1007/978-3-540-30140-0_69, doi:10.1007/978-3-540-30140-0_69.
- 51 Corwin Sinnamon and Robert E. Tarjan. A tight analysis of slim heaps and smooth heaps. In *Proceedings of the 2023 Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, page 549–567. Society for Industrial and Applied Mathematics, January 2023. URL: <http://dx.doi.org/10.1137/1.9781611977554.ch24>, doi:10.1137/1.9781611977554.ch24.
- 52 Corwin Sinnamon and Robert E. Tarjan. Efficiency of self-adjusting heaps. *ACM Transactions on Algorithms*, 21(4):1–39, September 2025. URL: <http://dx.doi.org/10.1145/3708989>, doi:10.1145/3708989.
- 53 Egor Starovoitov. weak-heap. <https://github.com/starovoid/weakheap>, 2022.
- 54 Mikkel Thorup. Equivalence between priority queues and sorting. *Journal of the ACM*, 54(6):28, December 2007. URL: <http://dx.doi.org/10.1145/1314690.1314692>, doi:10.1145/1314690.1314692.
- 55 P. van Emde Boas. Preserving order in a forest in less than logarithmic time. In *16th Annual Symposium on Foundations of Computer Science (sfcs 1975)*. IEEE, October 1975. URL: <http://dx.doi.org/10.1109/SFCS.1975.26>, doi:10.1109/sfcs.1975.26.
- 56 Raphael von der Grün. Superscalar sample queue. <https://github.com/raphinesse/s3q>, 2021.
- 57 Raphael von der Grün. Superscalar sample queue: Engineering a distribution-based priority queue. Master’s thesis, Karlsruhe Institute of Technology, 2021. URL: <https://ae.iti.kit.edu/english/4296.php>.
- 58 Jean Vuillemin. A data structure for manipulating priority queues. *Communications of the ACM*, 21(4):309–315, April 1978. URL: <http://dx.doi.org/10.1145/359460.359478>, doi:10.1145/359460.359478.
- 59 Jan Wassenberg, Mark Blacher, Joachim Giesen, and Peter Sanders. Vectorized and performance-portable quicksort. *Software: Practice and Experience*, 52(12):2684–2699, 2022. doi:10.1002/spe.3142.
- 60 Dan E. Willard. Log-logarithmic worst-case range queries are possible in space $\theta(n)$. *Information Processing Letters*, 17(2):81–84, August 1983. URL: [http://dx.doi.org/10.1016/0020-0190\(83\)90075-3](http://dx.doi.org/10.1016/0020-0190(83)90075-3), doi:10.1016/0020-0190(83)90075-3.
- 61 J.W.J. Williams. Algorithm 232: Heapsort. *Communications of the ACM*, 7(6):347–348, June 1964. URL: <http://dx.doi.org/10.1145/512274.3734138>, doi:10.1145/512274.3734138.

A Number of comparisons

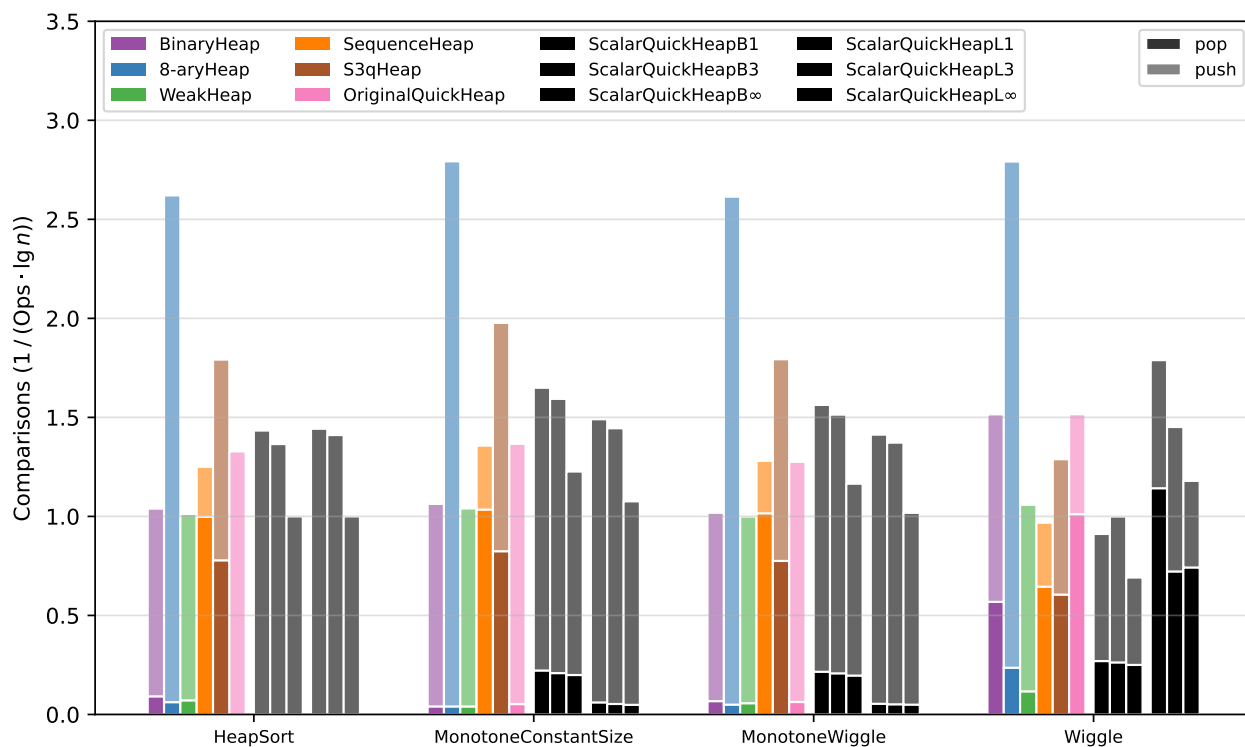


Figure 6 The number of comparisons per operation relative to the $\lg_2 n$ lower bound in different workloads. The number of comparisons during push operations is highlighted on the bottom, with the number of comparisons during pop operations shown on top. The SimdQuickHeap variants differ in using binary search (B) or linear search (L), and the pivot section (random, median of 3, true median oracle).