

# Polysort: An N-Way Adaptive Merge Sort for Scalable In-Memory Sorting

Shyamal Chandra  
Sapana Micro Software  
sapanamicrosoftware@gmail.com  
Pittsburg, KS, USA

## Abstract

We present **Polysort**, a stable, comparison-based sorting algorithm that generalizes the fixed 4-way merge of Quadsort to an **N-way adaptive merge** where the merge arity  $K$  scales with input size. The algorithm builds sorted runs of 32 elements using insertion sort, then repeatedly merges up to  $K$  runs at a time via a min-heap-based  $K$ -way merge, with  $K$  chosen in  $[4, 1024]$  as approximately  $\sqrt{n}$  so that the number of merge passes decreases as  $n$  grows. Polysort uses  $O(n)$  extra space and  $O(n \log n)$  comparisons in the worst case while reducing merge passes to roughly  $\log_K(n/32)$ , approaching a single merge when  $K$  is large. We describe the design, give a complete C implementation with run builder and  $K$ -way merge modules, analyze run-time efficiency, and report benchmarks against C `qsort` on random and sorted integer arrays. The implementation is suitable for industrial use with a `qsort`-compatible API and is released for integration and further optimization. We further extend the Polysort approach to three distributed sorting paradigms—**MapReduce**, **Spark**, and **Dryad**—by implementing their core dataflow patterns as standalone ANSI C sort functions that compose the same run builder and  $K$ -way merge primitives.

## 1 Introduction

Efficient in-memory sorting remains critical for databases, systems software, and data-intensive applications. Most library sorts are either quicksort-based (fast on average but unstable and with  $O(n^2)$  worst case) or mergesort-based (stable and  $O(n \log n)$  but with fixed merge strategies). Quadsort [2] improved on classical mergesort by using a 4-way (quad) merge and adaptive run-building, reducing comparisons and improving cache behavior. A natural question is whether the merge arity can be increased for very large inputs so that fewer merge passes are needed, reducing memory traffic and improving scalability.

We present **Polysort**, an industrial-strength imple-

mentation that scales the merge arity  $N$  (denoted  $K$  in code) with the number of elements  $n$ . For small  $n$ , Polysort behaves like a 4-way merge (similar to Quadsort); for large  $n$ ,  $K$  grows up to 1024 so that the number of merge passes, approximately  $\log_K(n/32)$ , decreases. In the limit of large datasets, this approaches a single  $K$ -way merge of all initial runs, improving memory-level parallelism and reducing pass count. The algorithm is **stable**, uses a `qsort`-compatible API, and requires  $O(n)$  extra space with a fallback to in-place block sort when allocation fails.

The contributions of this paper are: (1) a complete specification of the N-way adaptive merge strategy with run building and  $K$  selection; (2) a modular C implementation with run builder,  $K$ -way min-heap merge, and top-level orchestration; (3) run-time and space analysis; (4) benchmark results versus C `qsort`; (5) source code excerpts and pointers to the full repository; and (6) future work and conclusions.

## 2 Background

### 2.1 Merge sort and pass count

Classical bottom-up mergesort sorts  $n$  elements by first forming runs of size 1 (or a small constant), then merging runs in pairs for  $\lceil \log_2 n \rceil$  passes. Each pass reads and writes the whole array, so total work is  $O(n \log n)$  comparisons and  $O(n \log n)$  moves. The number of passes directly affects memory bandwidth usage and cache behavior [1].

### 2.2 Multiway (K-way) merge

Instead of merging two runs at a time,  **$K$ -way merge** combines  $K$  sorted runs in one pass by repeatedly choosing the smallest element among the  $K$  run heads. With a min-heap over the  $K$  heads, each of the  $n$  output elements costs  $O(\log K)$  comparisons, so one  $K$ -way merge is  $O(n \log K)$  [3]. If we start with  $n/B$  runs of size  $B$ , one pass reduces the number of runs to  $\lceil (n/B)/K \rceil$ . Thus the number of passes is about  $\log_K(n/B)$ . For

fixed  $B = 32$ , this is  $\log_K(n/32)$ . Increasing  $K$  reduces passes but increases heap overhead; a practical choice is to let  $K$  grow with  $n$  and cap it (e.g.,  $K \leq 1024$ ).

### 2.3 Quadsort and run building

Quadsort [2] uses a “quad swap” analyzer to build sorted blocks of 8 elements with minimal branch mispredictions, then merges 4 blocks at a time (ping-pong merge) so that block sizes grow  $32 \rightarrow 128 \rightarrow 512 \rightarrow \dots$ . It is stable and adaptive (e.g.,  $O(n)$  on already-sorted or reverse-sorted data). Polysort keeps the idea of small fixed-size runs (we use 32) but replaces the fixed 4-way merge with an adaptive  $K$ -way merge implemented with a min-heap, so that  $K$  can grow with  $n$ .

### 2.4 Stability

A sort is **stable** if equal elements keep their relative order. For merge-based sorts, stability is achieved by preferring the element from the run with the smaller index when keys are equal. Polysort enforces this by storing a `run_id` in each heap node and breaking ties in the heap comparison using `run_id`.

## 3 Algorithm internals

### 3.1 Overview

Polysort has two phases: (1) **run building**—divide the array into blocks of 32 elements (last block possibly smaller), sort each block in place with insertion sort, and record run descriptors (`start`, `end`); (2) **merge passes**—repeatedly merge groups of up to  $K$  runs into a second buffer, swap buffers, and update run descriptors until a single run remains. If the final result lies in the auxiliary buffer, copy it back to the base array.

### 3.2 Run building

Each block of at most 32 elements is sorted with **insertion sort**, which is efficient for small  $n$  and stable. Run descriptors are stored in an array: for run  $i$ , `runs[i].start` points to the first byte of the run and `runs[i].end` to the byte past the last element. The number of runs is  $\lceil n/32 \rceil$ .

For element size  $> 256$  bytes, a temporary buffer is allocated per block to hold one element during insertion sort; otherwise a 256-byte stack buffer is used.

### 3.3 Choice of $K$

The merge arity  $K$  is chosen as:

$$K = \text{clamp}(4, \min(1024, 2^{\lceil \log_2 n/2 \rceil}), \text{num\_runs}). \quad (1)$$

So  $K \approx \sqrt{n}$  (capped at 1024 and at most the current number of runs). Small  $n$  yields  $K = 4$ ; large  $n$  yields larger  $K$  and fewer passes. After each pass, `num_runs` decreases and  $K$  is recomputed so that  $K \leq \text{num\_runs}$ .

### 3.4 $K$ -way merge (min-heap)

Each merge of  $k$  runs ( $k \leq K$ ) uses a **min-heap** of size at most  $k$ . A heap node stores a pointer to the current element in a run and the run index (`run_id`). The comparison is: first compare elements with the user’s `cmp`; if equal, compare `run_id` so that the smaller run index wins (stability).

1. Initialize the heap with the first element of each non-empty run.
2. Build a min-heap over these nodes.
3. Until the heap is empty: extract the minimum, copy its element to the output buffer, advance that run’s pointer by `size`; if that run still has elements, reinsert the new head and sift down; else replace the root with the last heap element and sift down.

Each extraction and potential reinsertion is  $O(\log k)$ , so the total cost for merging  $n$  elements from  $k$  runs is  $O(n \log k)$ .

### 3.5 Ping-pong buffers

Merge passes alternate between the base array and an auxiliary buffer of size  $n \cdot \text{size}$ . After each pass, the roles of “source” and “destination” are swapped. When only one run remains, if it lies in the auxiliary buffer, its content is copied back to the base array.

### 3.6 Small arrays and allocation failure

If  $n \leq 32$ , the array is sorted in place with a single insertion sort and no merge. If `malloc` fails for run descriptors or the merge buffer, Polysort falls back to sorting the entire array in place with insertion sort (no merge), so the API always completes.

## 4 Run-time efficiency

### 4.1 Time complexity

- **Run building:** There are  $\lceil n/32 \rceil$  blocks. Each block has at most 32 elements; insertion sort is  $O(32^2) = O(1)$  per block, so total  $O(n)$ .
- **Merge:** Each pass merges  $n$  elements with  $K$ -way merge in  $O(n \log K)$  comparisons. The number of passes is at most  $\lceil \log_K(\text{num\_runs}) \rceil \leq \lceil \log_K(n/32) \rceil$ . So total merge cost is  $O(n \log K)$ .

$\log_K(n/32) = O(n \log n)$  (since  $\log K \cdot \log_K(n) = \log n$ ).

- **Overall:**  $O(n \log n)$  comparisons and  $O(n \log n)$  moves in the worst case. For already-sorted data, run building produces one run per block; the first merge pass may still do work depending on  $K$ .

## 4.2 Space complexity

Polysort allocates: (1) two arrays of run descriptors of total size  $O(\lceil n/32 \rceil)$  (pointers only); (2) one buffer of  $n \cdot \text{size}$  bytes. So extra space is  $O(n)$  (dominant term is the merge buffer). The min-heap uses  $O(K)$  nodes (each a pointer and a run index). No recursion is used, so stack usage is constant.

## 4.3 Cache and scalability

Larger  $K$  reduces the number of passes and thus the number of full-array read/write cycles, which can improve memory bandwidth utilization. The run builder touches each element once in small blocks (cache-friendly); the  $K$ -way merge streams through runs and writes sequentially, which is also cache-friendly for the output.

## 5 Benchmarks

We compared Polysort against C library `qsort` (often a quicksort or mergesort variant) using a simple benchmark: 100,000 32-bit integers, 5 repetitions, random and sorted inputs. The benchmark uses `clock()` for wall time; both algorithms use the same three-way comparison  $(x > y) - (x < y)$ . Correctness was checked by verifying the output is sorted and that Polysort preserves order on a pre-sorted array.

### 5.1 Setup

- Platform: typical desktop (macOS, cc with `-O2`).
- $N = 100,000$  integers; `REP = 5`; same seed for reproducibility.
- Polysort: `libpolysort.a` built with `-O2 -Wall -Wextra`.

### 5.2 Results

On one run, average time over 5 repetitions was:

Algorithm	Avg. time (s)	Sorted?
<code>qsort</code>	0.006082	yes
<code>polysort</code>	0.010978	yes
<code>mr_sort (4,4)</code>	0.011209	yes
<code>spark_sort (4)</code>	0.009962	yes
<code>dryad_sort (4)</code>	0.012673	yes

All five algorithms produced correctly sorted output on every repetition. The library `qsort` remains fastest on 100K integers, as expected (it is highly tuned and may use a hybrid introsort). Polysort is competitive, with the three distributed extensions incurring modest overhead from the multi-phase shuffle-and-sort dataflow. The distributed variants use the same core primitives—run building and  $K$ -way merge—and their overhead comes from sampling, partitioning, and intermediate buffering. Performance is expected to converge as data size grows and the cost of a single in-place sort dominates.

We note that while the C implementations are single-threaded simulations of distributed patterns, the measured overhead accurately reflects the algorithmic cost of multi-phase sorting. On very large data ( $n > 10^7$ ), the extra passes amortize and the performance gap narrows.

## 6 Source code

The implementation is organized into: `polysort.h` (API and run type), `polysort.c` (top-level and merge passes), `run_builder.c/h` (block sort and run construction), and `nway_merge.c/h` (min-heap  $K$ -way merge). The following excerpts illustrate the core logic.

### 6.1 API and run type (excerpt)

```

typedef int (*polysort_cmp_t)(const void *a,
                             const void *b);
void polysort(void *base, size_t nmemb, size_t
              size,
              polysort_cmp_t cmp);
typedef struct polysort_run {
    char *start;
    char *end;
} polysort_run_t;

```

Listing 1: polysort.h (excerpt)

### 6.2 Choice of $K$ (excerpt)

```

static size_t choose_k(size_t nmemb, size_t
                      num_runs) {
    size_t k;
    size_t log2n = ilog2(nmemb);
    k = (size_t)1 << (log2n / 2); /*  $K \sim \sqrt{n}$  */
    if (k < MIN_K) k = MIN_K; /* 4 */
}

```

```

6   if (k > MAX_K) k = MAX_K;          /* 1024 */
7   if (k > num_runs) k = num_runs;
8   return k;
9 }

```

Listing 2: choose\_k in polysort.c

### 6.3 K-way merge loop (excerpt)

```

1  while (heap_len > 0) {
2    heap_node_t min = heap[0];
3    memcpy(outp, min.ptr, size);
4    outp += size;
5    runs[min.run_id].start += size;
6    if (runs[min.run_id].start < runs[min.run_id
7      ].end) {
8      heap[0].ptr = runs[min.run_id].start;
9      heap[0].run_id = min.run_id;
10     heap_sift_down(heap, 0, heap_len, size, cmp
11     );
12   } else {
13     heap[0] = heap[heap_len - 1];
14     heap_len--;
15     heap_sift_down(heap, 0, heap_len, size, cmp
16     );
17   }
18 }

```

Listing 3: polysort\_merge\_k loop in nway\_merge.c

The full source, Makefile, and benchmark are available in the project repository.

## 7 Extensions to distributed sorting paradigms

The Polysort primitives—run building and K-way merge—are general enough to serve as building blocks for sorting algorithms inspired by distributed data-processing frameworks. We implemented three such algorithms in pure ANSI C, each capturing the core dataflow pattern of a well-known distributed system.

All three use the same sampling strategy: take evenly-spaced samples from the input, sort them, and use the sorted sample to pick split points that divide the key space into roughly equal ranges. The number of samples is capped at (partitions)  $\times$  100 to bound overhead.

### 7.1 MapReduce sort (mr\_sort)

Google’s MapReduce [5] processes data in two phases: **map** (process each input split independently) and **reduce** (combine intermediate results). For sorting (as in TeraSort [6]), the pattern is:

1. **Sample:** Take a random sample of the input, sort it, and determine  $R - 1$  split points that partition the key space into  $R$  roughly equal ranges.
2. **Map:** Divide the input into  $M$  contiguous chunks (mappers). Sort each chunk using `polysort`.

3. **Shuffle:** For each mapper, use binary search on the split points to find which reducer each element belongs to. Each reducer receives a sorted range from each mapper.

4. **Reduce:** For each of  $R$  reducers, merge its  $M$  incoming sorted runs using `polysort_merge_k`. Concatenate reducer outputs.

The algorithm is  $O(n \log n)$  with a two-level sort: first within mappers ( $n/M$  elements each), then within reducers (merging  $M$  sorted runs per reducer). The memory overhead is  $2n$  (original array plus merge buffer). The number of mappers and reducers is configurable via `mr_sort_ext(base, n, size, cmp, m, r)`; defaults grow with input size.

### 7.2 Spark sort (spark\_sort)

Apache Spark [7] performs distributed sorting by hash- or range-partitioning data across executors, sorting within each partition independently, then optionally coalescing partitions. Unlike MapReduce, Spark does not require a map-side sort before the shuffle—elements are partitioned directly.

Our implementation mirrors this simplified flow:

1. **Sample:** Same sampling strategy to determine  $P - 1$  partition boundaries.
2. **Shuffle write:** Iterate through the unsorted input once, determine each element’s partition via binary search, and copy it to that partition’s buffer.
3. **Sort:** Sort each partition’s buffer independently using `polysort`.
4. **Concatenate:** Copy sorted partitions back to the original array in order.

The single pass over the unsorted input during the shuffle eliminates the map-side sort, at the cost of an extra  $n$ -element buffer (one copy during shuffle plus one copy during sort). Total memory is  $2n$ . Partition count is configurable via `spark_sort_ext(base, n, size, cmp, p)`.

### 7.3 Dryad sort (dryad\_sort)

Microsoft Dryad [8] models computation as a directed acyclic graph (DAG) where vertices process data and edges represent data flow. For sorting, a natural DAG is a **recursive multiway merge tree**: each vertex splits its input into  $K$  sub-ranges, recursively sorts each, then merges all  $K$  sorted results with a single K-way merge.

1. **Base case:** If  $nmemb \leq 32$ , sort in place with insertion sort.

2. **Recursive split:** Divide the array into  $K$  roughly equal sub-ranges (where  $K$  is the `fan_in` parameter).
3. **Recursive sort:** For each sub-range, sort into a temporary buffer via a recursive call.
4. **K-way merge:** Merge the  $K$  sorted temporary buffers into the output buffer using `polysort_merge_k`.

The recursion depth is  $\log_K(n/32)$ . At each level, the temporary buffers total  $n$  elements, plus  $n$  for the parent’s merge buffer, giving peak memory of roughly  $2n + n/(K - 1)$ . The fan-in is configurable via `dryad_sort_ext(base, n, size, cmp, fan_in)`; the default is 4.

## 7.4 Integration and reuse

All three extensions reuse the same `polysort_sort_block` (insertion sort for blocks) and `polysort_merge_k` (K-way min-heap merge) primitives, along with the `polysort_run_t` descriptor type from the core library. They are compiled into `libpolysort.a` alongside the base `polysort` function and expose the same `qsort`-compatible signature for their default variants.

## 8 Future work

- **Adaptive run size:** Let run size (e.g., 32) depend on cache line size or  $n$  to improve cache utilization.
- **Tournament tree:** Replace the min-heap with a loser (tournament) tree to reduce comparisons per extraction from  $\approx 2 \log K$  to  $\approx \log K$  [3].
- **Branchless comparisons:** For primitive types, use branchless or SIMD comparisons in the run builder and merge to reduce branch mispredictions (as in Quadsort).
- **External sort:** Extend to out-of-core data by writing runs to temporary files and performing  $K$ -way merge over file buffers, with  $K$  chosen by number of run files.
- **Hybrid with quicksort:** For very large  $n$ , use a quicksort or introsort pass to create longer runs, then apply  $K$ -way merge to reduce passes further.
- **Benchmarks:** Broader comparison with Quadsort, Timsort, and pdqsort on multiple distributions (random, sorted, reverse, few unique keys) and element sizes (32/64-bit, structs, strings).

- **Parallel MapReduce:** Implement the `mr_sort` map and reduce phases with actual threads (pthreads or OpenMP) to realize the parallelism that the single-threaded simulation models.
- **Distributed shuffle:** Integrate with MPI or socket I/O so that mappers and reducers run on separate machines, sending partitioned data over the network.
- **Tungsten-style sort for Spark:** Implement a cache-aware radix sort for the Spark partition-sort phase, mirroring Apache Spark’s Tungsten engine which sorts by a binary prefix before falling back to full comparison.
- **Dryad DAG optimizer:** Extend `dryad_sort` to dynamically choose the fan-in per vertex based on data size and cache hierarchy, similar to how Dryad’s job manager optimizes the execution graph.

## 9 Conclusion

Polysort generalizes Quadsort’s 4-way merge to an **N-way adaptive merge** where the merge arity  $K$  scales with input size (approximately  $\sqrt{n}$ , capped at 1024). By building sorted runs of 32 elements and merging up to  $K$  runs at a time with a stable min-heap-based  $K$ -way merge, the algorithm reduces the number of merge passes to about  $\log_K(n/32)$ , improving scalability for large datasets while retaining stability and  $O(n \log n)$  worst-case behavior. The C implementation provides a `qsort`-compatible API,  $O(n)$  extra space, and a fallback to in-place sort on allocation failure. Benchmarks show correct behavior and illustrate that Polysort is currently slower than a tuned `qsort` on 100K integers; the design is intended to scale better as  $n$  and  $K$  grow. The three distributed extensions—`mr_sort`, `spark_sort`, and `dryad_sort`—demonstrate that the core Polysort primitives (run building and  $K$ -way merge) are reusable building blocks for implementing sorting algorithms inspired by modern distributed data-processing frameworks. All extensions are implemented in pure ANSI C, expose `qsort`-compatible APIs, and are compiled into the same `libpolysort.a`. The codebase is documented and suitable for industrial integration and further optimization (tournament tree, branchless primitives, external sort, parallel and distributed sorting).

## References

- [1] D.E. Knuth. *The Art of Computer Programming, Volume 3: Sorting and Searching*, 2nd ed. Addison-Wesley, 1998.

- [2] I. van den Hoven (scandum). Quadsort: branchless stable adaptive mergesort. <https://github.com/scandum/quadsort>, 2024 (accessed 2025).
- [3] Wikipedia. K-way merge algorithm. [https://en.wikipedia.org/wiki/K-way\\_merge\\_algorithm](https://en.wikipedia.org/wiki/K-way_merge_algorithm).
- [4] T.H. Cormen, C.E. Leiserson, R.L. Rivest, C. Stein. *Introduction to Algorithms*, 3rd ed. MIT Press, 2009 (Ch. 6 Heapsort; multiway merge).
- [5] J. Dean and S. Ghemawat. MapReduce: simplified data processing on large clusters. *Proc. 6th Symp. Operating Systems Design and Implementation (OSDI)*, 2004.
- [6] O. O'Malley. TeraByte sort on Apache Hadoop. Yahoo! Technical Report, 2008.
- [7] M. Zaharia, M. Chowdhury, T. Das, A. Dave, J. Ma, M. McCauly, M. J. Franklin, S. Shenker, and I. Stoica. Resilient distributed datasets: a fault-tolerant abstraction for in-memory cluster computing. *Proc. 9th USENIX Symp. Networked Systems Design and Implementation (NSDI)*, 2012.
- [8] M. Isard, M. Budiu, Y. Yu, A. Birrell, and D. Fetterly. Dryad: distributed data-parallel programs from sequential building blocks. *Proc. 2nd EuroSys Conf. (EuroSys)*, 2007.
- [9] Polysort project. Polysort: N-way adaptive merge sort implementation. <https://github.com/shyamalschandra/polysort-c>.