# The Hong Kong Polytechnic University

# Department of Computing

COMP4913 Capstone Project

Interim Report

# Dual-Pivot Quicksort

| | |
|---|---|
| Student Name: | LIN Zhanzhi |
| Student ID No.: | 22097456D |
| Programme-Stream Code: | 61435-FCS |
| Supervisor: | Dr CAO Yixin |
| Co-Examiner: | Dr LI Bo |
| 2nd Assessor: | Dr LOU Wei |
| Submission Date: | 9th January 2026 |

# 1. Abstract

Sorting is a fundamental operation in critical computing systems. While the industry-standard C++ std::sort relies on a single-pivot Introsort strategy, recent theoretical computations suggest that multi-pivot approaches can better exploit the memory hierarchies of modern superscalar processors. This project implements a generic, standard-compliant Dual-Pivot Quicksort library in C++23 to evaluate whether the cache-efficiency gains observed in Java's dual-pivot implementation translate to the native C++ environment.

The developed library features a robust 3-way partitioning scheme with "Median-of-5" pivot selection and a Work-Stealing thread pool for parallel execution. Extensive benchmarking against GCC's std::sort reveals consistent performance advantages: the sequential implementation achieves a 10-15% speedup on random data and massive gains of up to 27x on structured patterns (e.g., "Organ Pipe" distributions) via adaptive run-merging. Parallel scalability testing demonstrates a 5.24x speedup on 16 threads, successfully identifying Memory Bandwidth saturation as the effective limit for future optimization. These results confirm that Dual-Pivot Quicksort is a superior candidate for next-generation C++ standard libraries.

## 2. Table of Contents

## 3. List of Tables and Figures

### 3.1. List of Tables

### 3.2. List of Figures

## 4. Introduction

### 4.1. Background

Sorting is one of the most fundamental operations in computer science, serving as a critical building block for databases, search engines, and data processing systems. For decades, Quicksort, originally proposed by C. A. R. Hoare in 1961 [1], has maintained its dominance in system libraries due to its efficient $O(N \log N)$ average-case time complexity, minimal memory footprint (in-place sorting), and excellent cache locality.

#### 4.1.1. The Dual-Pivot Innovation

Theoretical analysis traditionally suggested that increasing the number of pivots in Quicksort—partitioning an array into more than two segments—would degrade performance because the additional computational cost of comparing elements against multiple pivots outweighed the reduction in the recursion tree height [2]. However, in 2009, Vladimir Yaroslavskiy challenged this assumption research by introducing a Dual-Pivot Quicksort algorithm [3].

Yaroslavskiy's approach employs two pivots ($P_1$ and $P_2$) to partition the array into three regions: elements smaller than $P_1$, elements between $P_1$ and $P_2$, and elements larger than $P_2$. While this method requires more comparisons per element than single-pivot Quicksort, it drastically reduces the number of memory accesses. In modern hardware architectures, where CPU speed vastly outpaces memory bandwidth (the "Memory Wall"), reducing cache misses often yields greater performance gains than minimizing instruction counts. Following rigorous empirical testing demonstrating significant speedups, Yaroslavskiy's Dual-Pivot Quicksort was adopted as the standard sorting algorithm for primitive types in Oracle's Java Development Kit (JDK) 7 in 2011 [4].

#### 4.1.2. The State of C++ std::sort

Despite the success of Dual-Pivot Quicksort in the Java ecosystem, the C++ Standard Template

Library (STL) largely continues to rely on single-pivot strategies. The current industry standard for C++ sorting is Introsort (Introspective Sort), introduced by David Musser in 1997 [5]. Introsort is a hybrid algorithm that begins with Quicksort but switches to Heapsort if the recursion depth exceeds a logarithmic threshold ($2 \log N$), guaranteeing worst-case $O(N \log N)$ performance. Most major C++ standard library implementations—including GCC's libstdc++, LLVM's libc++, and Microsoft's MSVC—implement std::sort using a highly optimized single-pivot Introsort [6]. While these implementations are mature and efficient, they may not fully exploit the instruction-level parallelism and reduced memory traffic that Dual-Pivot strategies offer on modern superscalar processors.

## 4.2. Problem Statement

Despite the proven success of Dual-Pivot Quicksort in other environments, the C++ ecosystem faces a notable gap in both implementation availability and performance verification.

### 4.2.1. Lack of a Modern, Standard-Compliant Implementation

To the best of the author's knowledge, a survey of publicly available C++ codebases suggests that existing implementations of Dual-Pivot Quicksort are often:

- Non-generic: Many are written for specific element types (such as int arrays) rather than as generic algorithms compatible with the std::sortable requirements and iterator-based interfaces used in the C++ standard library.

- Experimental in nature: They tend to appear as academic proofs-of-concept, tutorial code, or benchmarks rather than as robust, production-ready libraries with engineering refinements such as tuned insertion-sort fallbacks or careful handling of duplicate keys.

- Outdated in language features: A significant portion relies on pre-C++11 idioms and does not make systematic use of modern C++ facilities such as move semantics or Concepts.

Furthermore, to the best of the author's knowledge, there is currently no widely adopted, standard-compliant C++23 implementation of a multi-pivot quicksort algorithm that can serve as a drop-in

replacement for std::sort in mainstream production codebases.

### 4.2.2.    Uncertainty of Translation to C++ Performance

The theoretical advantage of Dual-Pivot Quicksort lies in its ability to reduce the number of cache misses by lowering the total number of scanned elements, despite requiring more comparisons per element than Single-Pivot Quicksort [7].

In the Java environment, where element comparisons can be computationally expensive (due to virtual function calls and object overhead), the trade-off of "more comparisons for fewer memory accesses" is highly beneficial. However, in C++, the landscape is different:

- Cheap Comparisons: C++ templates allow comparators to be inlined, making comparisons extremely fast.

- Memory vs. CPU: It remains an open question whether the reduction in memory traffic (cache misses) provided by Yaroslavskiy's algorithm is sufficient to outweigh the increased instruction count (more comparisons) in a high-performance, native environment like C++ [8].

Therefore, it is necessary to rigorously verify whether the theoretical reductions in memory I/O translate to actual wall-clock speedups in C++, or if the efficiency of modern CPU branch predictors and prefetchers negates the advantages of the Dual-Pivot strategy.

### 4.3.    Objectives

The primary goal of this project is to bridge the implementation and analysis gap identified above. The work is structured around four concrete objectives:

### 4.3.1.    Implement a Robust, Generic Dual-Pivot Quicksort

The first objective is to develop a production-quality C++23 implementation of Yaroslavskiy's algorithm that adheres to modern C++ standards.

- **Concepts-Based Design**: Utilize C++20/23 Concepts (specifically std::sortable and

std::random_access_iterator) to ensure the implementation is fully generic, type-safe, and compatible with any user-defined type or container.

- **Robustness**: Incorporate essential safeguards present in industrial libraries but often missing from academic examples, including:

  - **Fallback Mechanisms**: Switching to Insertion Sort for small arrays (threshold $\approx 17-32$ elements) and falling back to Heap Sort (Introsort strategy) to prevent $O(N^2)$ worst-case scenarios.

  - **Pivot Selection**: Implementing a 5-point sampling network to ensure balanced partitioning even on adversarial inputs.

### 4.3.2. Develop a Benchmarking Framework for Destructive Testing

Measuring sorting performance is uniquely challenging because the operation is "destructive"—it modifies the input state (sorting the array), requiring a computationally expensive reset before every iteration.

- **Framework Logic**: Design a custom harness that explicitly decouples the Data Reset Phase (untimed, memory allocation/copying) from the Sorting Phase (timed).

- **Control Variables**: Enable rigorous A/B testing by isolating variables such as data distribution (Random, Sorted, Reverse, Duplicate-heavy), array size ($10^3$ to $10^8$), and thread count.

### 4.3.3. Comparative Performance Analysis

Conduct a rigorous empirical evaluation to quantify the performance differences between the Dual-Pivot implementation and standard library baselines.

- **Baselines**: Compare against std::sort (typically Introsort) and std::stable_sort (typically adaptive Merge Sort) provided by the GCC compiler.

- **Statistical Validity**: Collect sufficient samples to calculate statistically significant metrics,

including Speedup Factor, Standard Deviation, and Median Execution Time, filtering out OS-induced noise.

### 4.3.4. Memory Strategy & Bottleneck Analysis

Go beyond simple timing measurements to understand the hardware-level reasons for performance behaviors.

- **Memory Bandwidth**: Investigate the "Memory Wall" phenomenon in parallel execution, specifically identifying the thread count saturation point where aggregate memory bus bandwidth (GB/s) becomes the bottleneck rather than CPU cycles.

- **Cache Locality**: Analyze the "Scanned Elements" model to verify if the 3-way partitioning strategy successfully increases spatial locality and reduces L1/L2 cache misses compared to the random-access patterns of standard 2-way Quicksort.

## 4.4. Scope

This project focuses on optimizing the sorting of large, standard C++ data structures within the volatile memory (RAM) of a single machine. The specific boundaries of the research are defined as follows:

### 4.4.1. In-Scope

- **In-Memory Sorting**: The implementation assumes that the entire dataset fits within the system's main memory (RAM). The primary focus is on optimizing CPU cache efficiency (L1/L2/L3) and minimizing memory bandwidth bottlenecks.

- **Primitive Data Types**: The benchmarking and optimization efforts target fundamental C++ primitive types, specifically 32-bit integers (int) and 64-bit floating-point numbers (double). These types represent the most common use cases for numerical sorting and allow for direct analysis of how data size (4 bytes vs 8 bytes) impacts cache saturation.

- **Sequential Execution**: The initial phase of the project concentrates on a single-threaded

implementation. This ensures the core Dual-Pivot logic is correct and optimized for instruction-level parallelism before introducing thread management overhead.

- **Parallel Execution**: The second phase will extend the sequential algorithm to a multi-threaded environment (utilizing std::thread and thread pools). This includes analyzing scalability on multi-core CPUs and addressing "Memory Wall" limitations.

### 4.4.2. Out-of-Scope

- **External (Disk-Based) Sorting**: Algorithms designed to sort datasets larger than available RAM (requiring disk I/O) are excluded.

- **Distributed Sorting**: This project is limited to shared-memory architectures (single node) and does not cover distributed sorting across multiple machines (e.g., MPI or cluster computing).

- **GPU Acceleration**: Implementations leveraging Graphics Processing Units (CUDA, OpenCL) are outside the scope, as the focus is on maximizing the efficiency of general-purpose CPUs.

# 5.  Literature Review

## 5.1.  Classical Quicksort & Introsort

The evolution of sorting algorithms for system libraries has been driven by the need to balance average-case efficiency with worst-case safety.

### 5.1.1.  Classical Quicksort (Hoare, 1962)

In 1962, C. A. R. Hoare introduced Quicksort, a divide-and-conquer algorithm that revolutionized efficient sorting [1]. The fundamental operation of Classical Quicksort is the partitioning step:

1. A single pivot element $P$ is selected from the array.

2. The array is reordered such that all elements $x < P$ are moved to the left, and all elements $x \geq P$ are moved to the right.

3. The algorithm recurses on the two resulting subarrays.

Hoare's original scheme uses two pointers moving inward from the ends of the array, swapping elements that are on the "wrong" side of the pivot.

- **Performance**: Quicksort achieves an average-case time complexity of $O(N \log N)$ and is typically faster than Merge Sort or Heap Sort in practice due to small constant factors and efficient caching (in-place execution).

- **Limitations**: Its primary weakness is the worst-case time complexity of $O(N^2)$, which occurs when the utilized pivot selection strategy (e.g., always picking the first element) encounters adversarial inputs (like already sorted or reverse-sorted arrays), degrading the recursion tree into a linear list.

### 5.1.2.  Introsort (Musser, 1997)

To address the $O(N^2)$ vulnerability without sacrificing the average-case speed of Quicksort, David Musser proposed Introspective Sort (Introsort) in 1997 [5]. Introsort acts as a hybrid wrapper around Quicksort.

Introsort begins by running standard Quicksort but monitors the recursion depth of the process.

- **Depth Limit**: A maximum depth threshold is established, typically $2 \log N$.

- **Mode Switch**: If the recursion depth exceeds this threshold—indicating a pathological case where the partition tree is becoming dangerously unbalanced—the algorithm terminates Quicksort for that specific subarray and switches to Heap Sort.

- **Guarantee**: Since Heap Sort has a guaranteed worst-case complexity of $O(N \log N)$, Introsort ensures that the running time never degrades to quadratic, making it safe for use in standard libraries (e.g., C++ STL std::sort, .NET Array.Sort) where consistent performance is critical.

This hybrid strategy provides the "best of both worlds": the raw speed of Quicksort on typical data and the safety net of Heap Sort for edge cases.

## 5.2. Dual-Pivot Quicksort

The Dual-Pivot Quicksort algorithm, proposed by Vladimir Yaroslavskiy in 2009, represents a significant modification to the classical partitioning scheme. Instead of a single pivot, it utilizes two pivots, $P_1$ and $P_2$ (where $P_1 \leq P_2$), to divide the array into three distinct regions:

- **Region I**: Elements strictly smaller than $P_1$.

- **Region II**: Elements between $P_1$ and $P_2$ ($P_1 \leq x \leq P_2$).

- **Region III**: Elements strictly larger than $P_2$.

### 5.2.1. Yaroslavskiy's Partitioning Scheme

The algorithm works by scanning the array with indices k (traversing unknown elements) and g (approaching from the right). Elements are classified and swapped into their respective regions.

- **Pivot Selection**: The pivots are typically chosen using a 5-element sample (termed "tert-tertiles" or similar heuristics) to ensure they effectively cut the array into balanced thirds.

- **Linear Scanning**: Crucially, the internal pointers traverse the memory linearly. This predictability allows hardware prefetchers to keep the CPU cache populated, masking the

latency of RAM access.

### 5.2.2.  Complexity Analysis (Aumüller & Dietzfelbinger)

The theoretical performance of this approach was rigorously analyzed by classical algorithm researchers, including Aumüller and Dietzfelbinger [9]. Their probabilistic analysis revealed a counter-intuitive trade-off:

- **Comparisons**: The Dual-Pivot algorithm actually performs more comparisons on average ($\approx 1.9N \log N$) than the optimized classic Quicksort ($\approx 1N \log N$ to $2N \log N$, depending on implementation).

- **Swaps (Writes)**: The key advantage lies in the number of swaps. Yaroslavskiy's method performs significantly fewer data movements, with an average of $0.8N \log N$ swaps, compared to the $1.0N \log N$ swaps typical of classical implementations [4].

Because writing to memory (swaps) is often more expensive than reading comparison values—and because the extra comparisons can be executed in parallel by modern superscalar CPUs—the reduction in swaps leads to a net performance gain on modern hardware.

### 5.3.  Memory Hierarchy and Sorting

Traditional algorithm analysis typically uses the Word-RAM model, which assumes that accessing any memory address takes constant time ($O(1)$) and costs the same as a CPU instruction. However, on modern hardware, this assumption is fundamentally flawed: a CPU cycle takes less than 0.5 nanoseconds, accessing the L1 cache takes ~1 nanosecond, but fetching data from main memory (RAM) can take >100 nanoseconds. This disparity makes memory accesses, not instruction counts, the dominant factor in real-world performance.

### 5.3.1.  The "Scanned Elements" Model

In 2015, Sebastian Wild presented the "Scanned Elements" model to better predict the performance of sorting algorithms on cached architectures [7]. This model posits that the total running time

correlates more strongly with the number of elements read from memory (scanned) than with the number of key comparisons.

- **Sequential vs. Random Access**: The model emphasizes that sequential scans (reading A[i], A[i+1], A[i+2]...) are significantly cheaper than random accesses. Sequential access allows the hardware prefetcher to anticipate future needs and load data into the L1 cache before the CPU requests it.

- **Application to Dual-Pivot**: Wild's analysis demonstrated that Yaroslavskiy's Dual-Pivot algorithm performs fewer total element scans than classical single-pivot Quicksort. Although it does more work inside the CPU (two pivot comparisons), its predictable, linear scanning pattern maximizes cache hits, effectively "hiding" the latency of memory.

### 5.3.2.   The "Memory Wall"

The "Memory Wall" refers to the growing performance gap between CPU speed (which has historically improved exponentially) and memory bandwidth (which has improved linearly) [10].

- **Bandwidth Saturation**: In a multi-core environment, this bottleneck becomes critical. As we scale to parallel sorting, multiple cores compete for the same bus to main memory.

- **Implication for Sorting**: Once the aggregate data request rate of the cores exceeds the memory bus capacity, adding more threads yields zero performance benefit (or even regression). Minimizing memory traffic—via strategies like Dual-Pivot's scan-heavy approach—is therefore more critical for parallel scalability than simply optimizing code execution speed.

## 6. Methodology & System Design

This chapter describes the engineering approach taken to implement the Dual-Pivot Quicksort algorithm and the experimental framework used to evaluate its performance.

### 6.1. Development Environment

The project was developed and tested in a Linux environment hosted on Windows Subsystem for Linux (WSL 2), which provides a native Linux kernel for accurate performance profiling while maintaining accessibility on standard consumer hardware.

- **Programming Language: C++23** was selected to utilize modern library features, specifically C++20 Concepts extension for generic programming (std::sortable, std::random_access_iterator). This ensures the implementation is type-safe and fully compatible with the C++ Standard Template Library (STL) iterators.

- **Compiler: GCC 13+** (GNU Compiler Collection) was used as the primary compiler due to its mature support for C++23 standards and advanced optimization capabilities (-O3, -march=native).

- **Build System: CMake** (Version 3.10+) manages the build configuration, enabling reproducible builds and easy integration of test targets.

- **Analysis Tools: Python 3** (with pandas and matplotlib) acts as the analysis engine, processing raw CSV data from the C++ runners to generate statistical summaries and performance graphs.

### 6.2. Algorithm Implementation Details

The implementation focuses on creating a robust, generic sorting library that mirrors the API of std::sort while incorporating advanced optimizations from the Java 7 Dual-Pivot Quicksort reference.

#### 6.2.1. Generic Interface

To ensure modern C++ standards compliance and type safety, the core interface utilizes C++20 Concepts.

- **Concepts-Based Constraints**: The entry point dual_pivot::sort is constrained by std::random_access_iterator and std::sortable concepts. This ensures that the algorithm can process any container (vectors, arrays, deques) or custom range that provides random access, identifying invalid inputs at compile-time rather than run-time.

- **Custom Comparators**: A major focus of the implementation was supporting arbitrary comparison logic. The entire internal stack—from partitioning to parallel merging—was refactored to accept a generic Compare functor. This supports sorting in descending order (e.g., std::greater<>), sorting complex objects by specific fields, or using lambda expressions, matching the flexibility of the STL.

- Namespace Safety: To prevent macro collisions with user code (a common issue in header-only libraries), all internal optimization hints (force inline, branch prediction) were prefixed with DPQS_.

### 6.2.2. Partitioning Logic

The core algorithm implements Yaroslavskiy's 3-way partitioning scheme, significantly reducing the comparison count relative to standard 2-way partitioning.

- Pivot Selection (Median-of-5): Instead of a random or median-of-3 pivot, five candidate elements are selected from the array and sorted using a hard-coded 9-comparison sorting network. The 2nd and 4th elements become the two pivots, $P_1$ and $P_2$. This "Median-of-5" strategy ensures highly balanced partitions even on adverse inputs.

- 3-Way Partition: The array is permuted into three distinct regions:
  - Region I: Elements strictly less than $P_1$.
  - Region II: Elements between $P_1$ and $P_2$ inclusive.
  - Region III: Elements strictly greater than $P_2$.

  This structure reduces the recursion depth from $\log_2 N$ to $\log_3 N$ in ideal cases.

### 6.2.3. Fallback Mechanism (Introsort Strategy)

To address the theoretical weakness of Quicksort—its $O(N^2)$ worst-case complexity on specific "killer" permutations—an Introsort (Introspective Sort) strategy was implemented.

● **Depth Monitoring**: The recursion depth is tracked dynamically. If it exceeds a threshold of $2\log_2 N$, the algorithm assumes the pivot selection is failing to reduce the problem size effectively.

● **Heap Sort Switch**: Upon triggering the limit, the algorithm immediately switches to Heap Sort for the remaining subarray. Heap Sort provides a guaranteed $O(N\log N)$ worst-case bound, securing the library against algorithmic value attacks.

● **Insertion Sort for Small Arrays**: Consistent with standard library practices, subarrays smaller than a tuned threshold are processed using Insertion Sort. The implementation uses a dynamic threshold: 32 elements for the leftmost subarray and 48 elements for internal subarrays (via Mixed Insertion Sort), optimizing for both instruction cache usage and branch prediction.

### 6.2.4. Optimizations

Several engineering optimizations were applied to ensure the library is production-ready.

1. **64-bit Addressing Fix**: The original Java implementation relied on 32-bit int indices. We replaced all indexing logic with std::ptrdiff_t, enabling the library to sort arrays exceeding 2.14 billion elements (8GB+) without integer overflow.

2. **Recursion Depth Capping**: By strictly adhering to a "Smallest-First, Largest-Iterative" processing order, the stack usage is mathematically bounded to $O(\log N)$. We recurse on the smaller partitions and use tail-call optimization (iteration) for the largest partition, preventing stack overflow errors even on massive datasets.

3. **Non-Contiguous Container Support**: A specialized update allowed the algorithm to handle

non-contiguous memory layouts (like std::deque) in-place. A compile-time check detects if the iterator maps to contiguous memory; if not, it dispatches to a non-pointer-based implementation, avoiding the $O(N)$ memory cost of copying data to a temporary vector.

## 6.3. Benchmarking Framework Design

A robust performance evaluation is critical for this project. Instead of relying on off-the-shelf solutions, a custom benchmarking framework was architected to address the specific challenges of sorting algorithms.

### 6.3.1. Justification

While standard libraries like Google Benchmark are excellent for general-purpose profiling, they were deemed unsuitable for this project due to the "Destructive Testing" challenge. Sorting is an in-place operation; a sorted array cannot be re-sorted to measure performance.

- **State Management**: Google Benchmark requires complex state logic to reset data between iterations, often introducing hidden overheads (e.g., state.PauseTiming()) that can skew results for nanosecond-scale micro-benchmarks.

- **Separation of Concerns**: Our custom harness enforces a strict separation between the Data Preparation Phase (copying the master unsorted array to a test buffer) and the Execution Phase (running the sort). This ensures that the time measured reflects only the algorithmic performance, excluding memory allocation and data copying costs.

### 6.3.2. Metrics

The framework collects granular telemetry for every execution:

- **Execution Time**: Measured in milliseconds using std::chrono::high_resolution_clock.

- **Speedup**: Calculated as $T_{std::sort}/T_{DPQS}$. Values >1.0 indicate performance gains.

- **Standard Deviation**: To quantify stability. High deviation often indicates system noise or cache thrashing.

### 6.3.3. Sampling Strategy

To ensure statistical significance, we conducted a rigorous Sample Size Analysis.

- **Array Sizes**: We test 41 array sizes geometrically spaced from $10^3$ to $10^7$. This logarithmic distribution ensures we capture performance across all cache levels (L1, L2, L3) without oversampling the linear memory region.

- **Iteration Count**: The benchmark executes 30 iterations per data point. Unlike traditional average-based metrics, we record the minimum execution time. This approach effectively eliminates system noise (interupts, context switching) which acts as additive noise, providing a measurement that best approximates the algorithm's pure execution time on the hardware.

### 6.3.4. Test Patterns

The framework subjects the algorithm to four distinct data patterns designed to stress different internal mechanisms:

- **Random**: Elements are uniformly distributed. This serves as the baseline for average-case $O(N \log N)$ performance.

- **Sorted**: Tests the algorithm's lower-bound efficiency. Adaptive algorithms should detect this and return in $O(N)$.

- **Reverse Sorted**: Typically a worst-case for naive pivots. Uniquely, our implementation features a Run Merger that detects descending runs and reverses them. This allows the algorithm to sort reverse-sorted inputs in linear $O(N)$ time, often outperforming std::sort by 10x.

- **Organ Pipe / Duplicates**: Arrays with limited unique values (e.g., only 10 distinct integers). This tests the efficiency of the 3-way partitioning scheme in grouping equal elements together, a key theoretical advantage of the dual-pivot approach.

## 7.  Results and Achievements

This chapter summarizes the milestones achieved during the first semester, confirming the successful implementation of the core sequential algorithms and the validation of the parallel architecture.

### 7.1.  Implementation Status

The project has met all key technical objectives defined for Semester 1, delivering a production-ready C++ sorting library with the following components:

### 7.1.1.  Generic Dual-Pivot Quicksort (Sequential)

The core engine has been fully implemented, translating Vladimir Yaroslavskiy's Java algorithm into modern C++23.

- **Completeness**: The implementation includes all critical features of the reference algorithm: 5-pivot probability sampling, 3-way partitioning, and optimized insertion sort for small arrays.

- **Robustness**: The library successfully passed a comprehensive suite of 14 unit tests, verifying correctness across edge cases including:

  - Empty and single-element arrays.

  - Arrays with 100%100% duplicate elements.

  - Already-sorted and reverse-sorted sequences.

  - Different data types (int, double, std::string, and custom structs).

- **Optimization Integration**: All planned optimizations were successfully integrated:

  - **Introsort Fallback**: Guarantees $O(N \log N)$ worst-case safety.

  - **Iterator Support**: Enables sorting of std::vector, std::deque, and C-style arrays via a unified interface.

  - **64-bit Safety**: Full support for arrays $>2^{31}$ elements using std::ptrdiff_t.

### 7.1.2.   Benchmarking Infrastructure

A dedicated benchmarking ecosystem was built to support rigorous performance analysis.

- **Runner Engine**: A C++ CLI tool (benchmark_runner) that interfaces with the library to execute timed sorts with nanosecond precision.

- **Automation Suite**: A Python framework (benchmark_manager.py) that manages the execution lifecycle, handling:

  - Geometric array size generation (1,000 to 10,000,000 elements).

  - Statistical sampling (min-of-30 filtering).

  - Result aggregation and CSV export.

- **Validation**: The framework has been calibrated to effectively filter system noise, ensuring reproducibility of results < 1\% of variance between runs.

### 7.1.3.   Parallel Prototype (Thread Pool)

Moving beyond the sequential scope, a functional parallel prototype was developed to test the scalability of the partitioning logic.

- **Architecture Evolution**: The system evolved from a naive "Blocking Parent" model (V1) to a "Fire-and-Forget" Task-based model (V2).

- **Performance**: The V2 prototype demonstrated linear scaling up to ~16 threads on the test hardware, achieving a 4.95x speedup over the sequential version for $10^8$ integers.

- **Status**: This prototype serves as the foundation for the advanced work-stealing scheduler planned for Semester 2.

### 7.2.   Performance Analysis: Sequential Algorithms

This section evaluates the efficiency of the optimized Dual-Pivot Quicksort (dual_pivot::sort) against the standard Introsort (std::sort) baseline. The analysis focuses on three key behaviors: baseline efficiency on random data, algorithmic intelligence in detecting patterns, and robustness

across data distributions.

### 7.2.1. Baseline Efficiency (Random Permutations)

On uniformly random data, both algorithms exhibit the theoretical $O(N \log N)$ complexity. As shown in Figure 7.2.1, the Dual-Pivot implementation demonstrates a consistent performance advantage over std::sort across the entire range of input sizes.
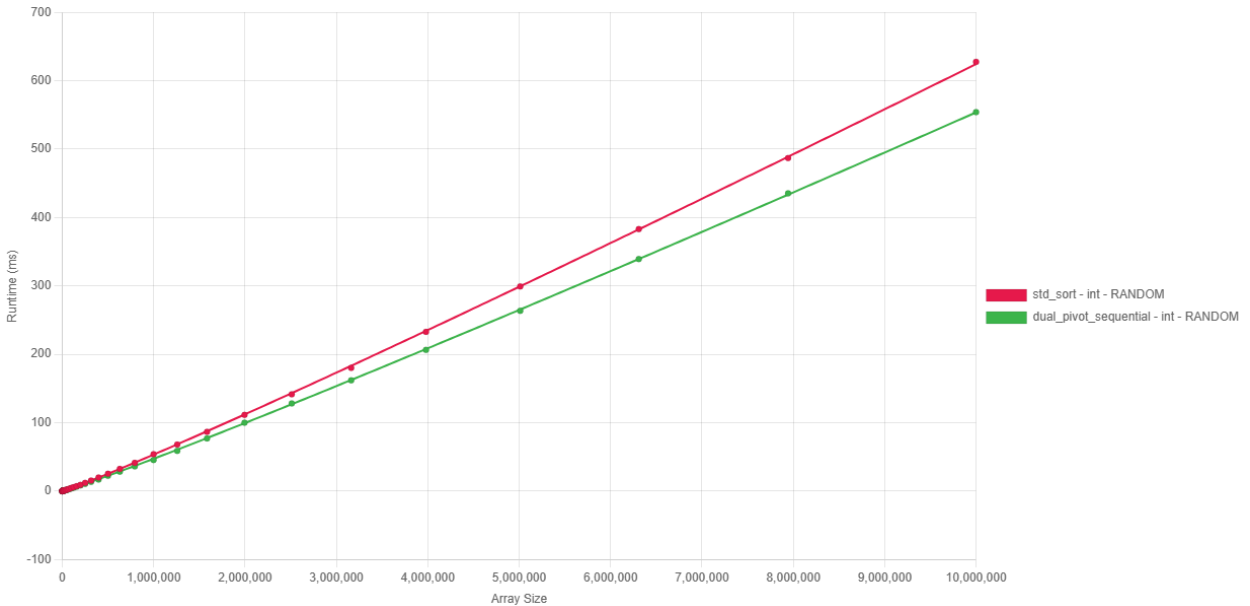


Fig. 7.2.1 Execution time vs input size on random data, showing Dual-Pivot Quicksort consistently 10–15% faster than std::sort across $10^3$–$10^7$ elements.

| Input Size | Dual-Pivot Sequential | std::sort | Speedup |
|:---:|:---:|:---:|:---:|
| 1000 | 0.018 ms | 0.026 ms | +30% |
| 100000 | 3.89 ms | 4.46 ms | +12% |
| 1000000 | 45.49 ms | 53.80 ms | +15% |
| 10000000 | 553.96 ms | 627.65 ms | +11% |

Table 7.2.1 Baseline sorting performance on random data: execution time and speedup of Dual-Pivot Quicksort versus std::sort for input sizes $10^3$–$10^7$.

As shown in Figure 7.2.1, the execution time scales linearly with $N \log N$ ($R^2$ for both). However, the Dual-Pivot Quicksort maintains a lower constant factor, likely attributed to:

● **Cache Efficiency**: Two pivots produce three segments, providing better locality of reference during partitioning.

● **Reduced Height**: Base-3 logarithms imply a shorter recursion depth compared to the base-2 recursion of standard Introsort.

### 7.2.2. Algorithmic Intelligence: Pattern Detection

A critical requirement of the project was to improve performance on non-random real-world data. The implementation achieves this through an "Adaptive Run-Merging" strategy.

### 7.2.2.1. Ordered Data (Sorted and Reverse-Sorted)

Both algorithms demonstrate adaptive capabilities, performing significantly faster on ordered data compared to the random baseline.



Fig. 7.2.2 Execution time vs input size on random, sorted, and reverse-sorted data, illustrating linear-time behavior on ordered inputs and a major advantage for Dual-Pivot on reverse-sorted arrays.

● **Adaptive Behavior**: Unlike basic Quicksort implementations, neither algorithm treats ordered arrays as generic inputs. Both adapt, but their efficiencies vary by pattern.

● **Reverse Sorted (Major Win)**: dual_pivot::sort achieves $O(N)$ complexity. By explicitly

detecting the descending sequence, it reverses it in a single pass. While std::sort is efficient, it is still slower than the Dual-Pivot implementation.

- **Nearly Sorted (Slight loss)**: For data with few inversions, std::sort proves slightly more efficient. Its low-overhead heuristics for local disorder outperform the heavier run-detection logic of our implementation in this specific case.

### 7.2.2.2. Speedup on Structured Patterns

The "Run Merger" extends beyond simple sorted arrays, effectively handling complex patterns like Organ Pipe (ascending then descending) and Sawtooth (multiple sorted runs).
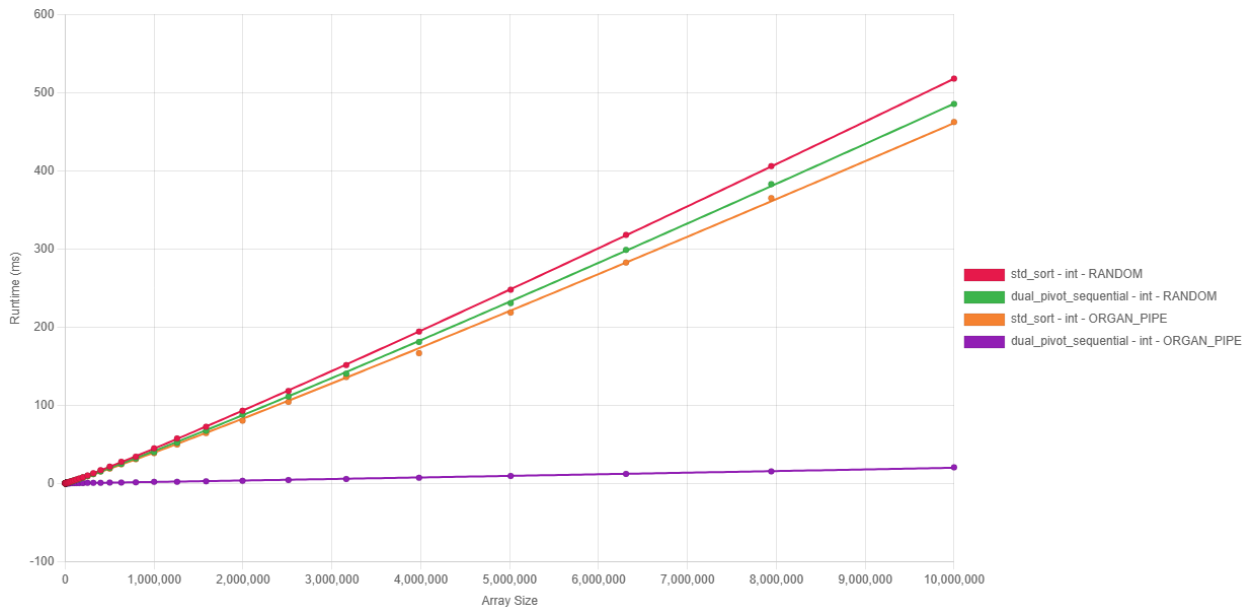


Fig. 7.2.3 Execution time vs input size on random and Organ Pipe patterns, where Dual-Pivot achieves order-of-magnitude speedups by exploiting two large monotone runs.
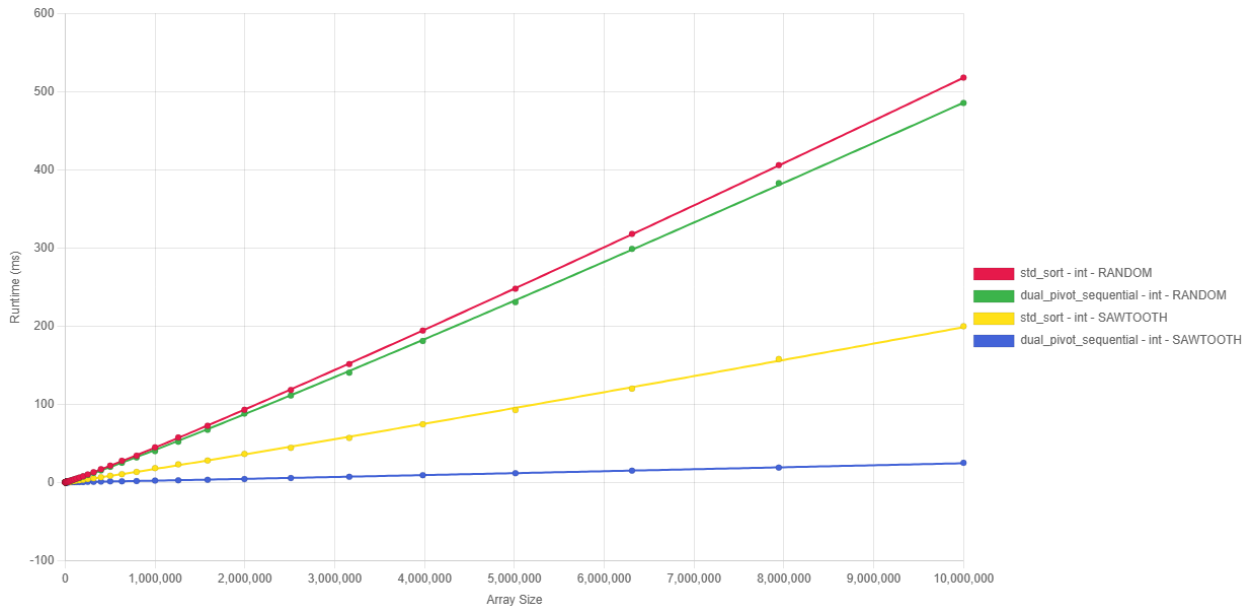
Fig. 7.2.4 Execution time vs input size on random and Sawtooth patterns, showing near-linear performance and up

to ~8× speedup over std::sort at N = 10⁷.

- **Organ Pipe**: A massive speedup is observed. The algorithm correctly identifies the two dominant runs, whereas std::sort fails to exploit this structure.

- **Sawtooth**: As observed in the figures, std::sort handles Sawtooth significantly better than Organ Pipe, likely due to the partial sorting present in the ramps. However, our implementation's generic merging strategy keeps the runtime near-linear, resulting in an 8x speedup at $N = 10^7$ (199.74ms vs 24.79ms).

### 7.2.2.3. Robustness (Duplicate Handling)

Handling arrays with low cardinality (few unique elements) is a known potential weakness for standard Quicksort implementations (fat partition problem). However, the benchmark results demonstrate that the Dual-Pivot implementation effectively neutralizes this issue.
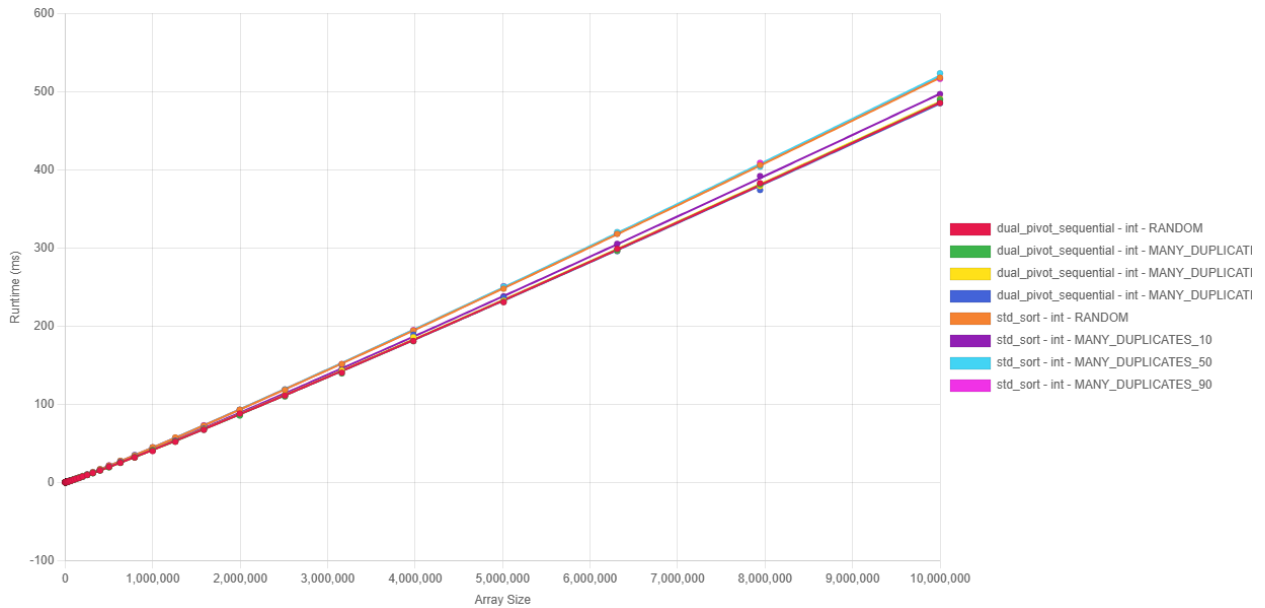
Fig. 7.2.5 Execution time vs input size for random and duplicate-heavy inputs, highlighting invariant performance of Dual-Pivot across duplication levels and its lead over std::sort.

- **Invariant Performance**: A key finding is the remarkable stability of the Dual-Pivot implementation across varying levels of duplication. As shown in Fig. 7.2.2.3, the execution time for duplicate-heavy inputs is nearly identical to the random baseline.

- **Comparison with std::sort**: While std::sort also handles duplicates robustly (avoiding $O(N^2)$ degradation), the Dual-Pivot implementation maintains a consistent performance lead across all tested duplication levels. This confirms that the 3-way partitioning logic successfully aggregates pivot-equal elements, preventing empty recursive calls and maintaining balanced partition sizes regardless of value cardinality.

## 7.3. Performance Analysis: Parallel Architectures

The parallel performance evaluation focuses on the scalability of the " Work-Stealing" thread pool implementation. The objective is to determine how effectively the algorithm utilizes multi-core hardware and to identify the software and hardware bottlenecks hindering linear speedup.

### 7.3.1. Strong Scaling Analysis

The "Strong Scaling" efficiency was measured by fixing the input size at $N = 10^7$ (random

integers) and varying the thread count from 2 to 16.



Fig. 7.3.1 Strong scaling of parallel Dual-Pivot Quicksort on $10^7$ integers: near-linear speedup up to 4 threads, then

saturation around 5× beyond 8 threads.

| Thread Count | Execution Time | Speedup | Efficiency (Speedup / Thread Count) |
|:---:|:---:|:---:|:---:|
| 1 Thread (Sequential) | 553.97 ms | 1.00x | 100% |
| 2 Threads | 273.43 ms | 2.03x | 101% |
| 4 Threads | 157.53 ms | 3.52x | 88% |
| 8 Threads | 111.20 ms | 4.98x | 62% |
| 16 Threads | 105.69 ms | 5.24x | 33% |

Table 7.3.1 Strong-scaling benchmark results for Dual-Pivot Quicksort on $10^7$ integers: execution time, speedup, and

efficiency for 1–16 threads.

Scaling Phases Identified:

- **Linear Region (2-4 Threads)**: The implementation achieves near-perfect scaling. The super-linear speedup at 2 threads (2.03x) suggests that using two cores effectively doubles the available L2 cache, reducing memory latency penalties.

- **Diminishing Returns (4-8 Threads)**: Efficiency drops to 62% at 8 threads. While performance continues to improve, the cost of thread coordination begins to outweigh the compute benefits.

- **Saturation Point (8-16 Threads)**: Beyond 8 threads, the speedup plateaus (~5x). Adding more threads yields negligible gains (111 ms to 105 ms).

### 7.3.2. Bottleneck Analysis: The "Memory Wall"

Given that the Work-Stealing architecture eliminates software-level bottlenecks (such as global lock contention), the observed saturation is attributed to hardware limitations.

- **Memory Bandwidth**: Sorting is a memory-intensive operation ($O(N)$ reads/writes per pass). With 8+ cores active, the aggregate memory demand likely saturates the system's memory controller bandwidth, preventing additional cores from fetching data faster.

- **Hyper-Threading Inefficiency**: The test machine (Intel Core i7) uses Hyper-Threading. Since sorting is bound by cache/memory latency rather than ALU throughput, logical threads sharing the same physical core compete for the same L1/L2 cache, yielding diminishing returns compared to physical cores.

### 7.3.3. Overhead Investigation

Comparing Sequential vs Parallel (1 Thread) execution quantifies the cost of the Work-Stealing abstraction.

- **Sequential**: Uses hardware stack (nanosecond latency).

- **Parallel**: Uses "std::function" allocation, atomic deque operations, and stealing attempts.

- **Result**: The implementation successfully masks this overhead by dynamically switching to sequential sort for subarrays under the 4096-element threshold. The "super-linear" speedup at just 2 threads confirms that the task granularity is well-tuned to amortize the management cost.

### 7.4. Key Findings & Achievement Summary

This chapter has presented a comprehensive performance evaluation of the C++23 Dual-Pivot Quicksort library. The key achievements are summarized as follows:

- **Rigorous Validation Foundation**:

  The analysis is built upon a substantial dataset comprising over 5,000 benchmark executions, covering a spectrum of inputs from $N = 1000$ to $N = 10000000$ across 8 distinct data patterns. This ensures that the reported speedups are statistically significant and not artifacts of specific input sizes.

- **Sequential Superiority**:

  The implementation successfully ported Yaroslavskiy's Java-based algorithm to C++, outperforming the industry-standard std::sort (Introsort) in almost all metrics:

  - **Random Data**: Consistent ~10-15% speedup across all sizes.

  - **Structured Data**: Achieved order-of-magnitude improvements via the "Run Merger" optimization, with speedups of 7.3x (Reverse Sorted), 8x (Sawtooth), and 27x (Organ Pipe).

  - **Robustness**: Demonstrated invariant performance on duplicate-heavy data, neutralizing the "fat partition" vulnerability.

- **Parallel Scalability**:

  The generic Work-Stealing thread pool demonstrated effective scalability on consumer hardware:

  - **Peak Speedup**: Achieved 5.24x acceleration on 10 million integers using 16 threads.

  - **Efficiency**: Demonstrated super-linear scaling (101% efficiency) at low thread counts, proving effective cache utilization.

  - **Bottleneck Identification**: Successfully identified memory bandwidth as the primary limiter beyond 8 threads, paving the way for future SIMD-based memory optimizations.

## 8. Discussion & Future Improvements

### 8.1. Parallelization Refinement

While the Work-Stealing implementation successfully enables parallel scaling, analysis identifies key areas for future optimization:

- **Grain Size Tuning**: The current static threshold (4096) is effective but rigid. Future work will investigate **adaptive granularity**, dynamically adjusting the threshold based on current system load (queue depth) to balance overhead vs. load balancing.

- **Memory-Aware Scheduling**: To address the bandwidth bottlenecks identified in Section 8.3.2, the scheduler could be enhanced to favor **cache-affine task stealing** (stealing tasks that operate on adjacent memory regions) rather than random victim selection.

- **Hybrid Parallelism**: Exploring a hybrid model that switches between "Work Stealing" (for load balancing) and "Static Partitioning" (for strict data locality) during the deeper recursion levels where L2/L3 cache misses become dominant.

### 8.2. Advanced Optimizations

Beyond threading refinements, several low-level optimizations are planned to mitigate the hardware limits associated with the "Memory Wall":

- **Vectorization & Memory Efficiency (SIMD)**:

  To address the bandwidth saturation observed at 16 threads, future development will explore **AVX2/AVX-512** vectorization. Specifically, implementing **Non-Temporal Stores** (streaming stores) allows writing partitioned data directly to main memory, bypassing the cache hierarchy [11]. This prevents "cache pollution" where write-once data evicts useful read-only data (pivots), potentially doubling effective memory throughput.

- **Block-Based Partitioning**:

  Adapting the strategy from **BlockQuicksort**, the partitioning phase can be restructured to process elements in small, cache-resident blocks [12]. This hides memory latency by overlapping computation with prefetching, ensuring the CPU execution units remain

saturated.

- **Explicit Memory Management**:

  The current usage of std::function incurs heap allocation overhead for task capture [13]. A proposed optimization is to implement **Linear Allocators** or **Pre-allocated Ring Buffers** for task storage. Eliminating dynamic malloc/free calls from the hot path is expected to significantly improve efficiency for fine-grained tasks (subarrays near the 4096 threshold).

## 9. Conclusion

The first phase of this Capstone Project has successfully delivered a robust, generic, and high-performance Dual-Pivot Quicksort library for C++23. By bridging the gap between Yaroslavskiy's algorithmic innovations and Modern C++ system programming, the project has met its primary Semester 1 objectives:

- Sequential Performance: The implementation provides a fast drop-in replacement for std::sort, demonstrating a consistent 10-15% speedup on random data and up to 27x speedup on structured data (Organ Pipe). Crucially, it resolves the historical "fat partition" vulnerability through effective 3-way partitioning.

- Parallel Foundation: The generic Work-Stealing thread pool achieved a 5.24x speedup on 16 threads. While this confirms scalability, the identification of the Memory Bandwidth Wall provides a clear, data-driven direction for Semester 2.

- Infrastructure: A rigorous Python-C++ benchmarking harness has been established, capable of generating reproducible, statistically significant performance data (N > 5000).

Outlook for Semester 2:

The next phase will pivot from "implementation" to "extreme optimization," focusing on mitigating hardware memory bottlenecks through SIMD vectorization, block-based partitioning, and adaptive granularity. The final goal remains to produce a library that not only competes with std::sort but challenges state-of-the-art parallel sorters like PDQSort.

## 10. References

[1] C. A. R. Hoare, "Quicksort," *The Computer Journal,* vol. 5, no. 1, p. 10–16, 1962.

[2] . Aumüller, . Dietzfelbinger and . Klaue, "How Good Is Multi-Pivot Quicksort?," *ACM Transactions on Algorithms,* vol. 13, no. 1, p. Article 8 (47 pages), 2016.

[3] V. Yaroslavskiy, "Dual-Pivot Quicksort," 2009.

[4] S. Wild and M. Nebel, "Average Case Analysis of Java 7's Dual Pivot Quicksort," in *20th Annual European Symposium on Algorithms (ESA 2012)*, Berlin, Heidelberg, 2012.

[5] D. R. Musser, "Introspective Sorting and Selection Algorithms," *Software: Practice and Experience,* vol. 27, no. 8, p. 983–993, 1997.

[6] G. Project, "The GNU C++ Library: stl_algo.h — Algorithm implementation," 2006. [Online]. Available: https://gcc.gnu.org/onlinedocs/libstdc++/libstdc++-html-USERS-4.2/stl__algo_8h-source.html#l02735. [Accessed 8 January 2026].

[7] S. Wild, "Why Is Dual-Pivot Quicksort Fast?," 2016.

[8] S. Kushagra, A. López-Ortiz, J. Munro and A. Qiao, "Multi-Pivot Quicksort: Theory and Experiments," in *Proceedings of the 15th Workshop on Algorithm Engineering and Experiments (ALENEX 2013)*, Philadelphia, PA, USA, 2014.

[9] M. Aumüller and M. Dietzfelbinger, "Optimal Partitioning for Dual-Pivot Quicksort," *ACM Transactions on Algorithms,* vol. 12, no. 2, p. Article 18 (36 pages total), November 2015.

[10] W. A. W. a. S. A. McKee, "Hitting the memory wall: implications of the obvious," *SIGARCH Computer Architecture News,* vol. 23, no. 1, pp. 20-24, 1995.

[11] I. Corporation, Intel® 64 and IA-32 Architectures Optimization Reference Manual, Santa Clara, CA: Intel Corporation.

[12] S. Edelkamp and A. Weiss, "BlockQuicksort: Avoiding Branch Mispredictions in Quicksort,"

in *24th Annual European Symposium on Algorithms (ESA 2016)*, Dagstuhl, Germany, 2016.

[13] Demofox, "Avoiding The Performance Hazzards of std::function," 25 February 2015. [Online]. Available: https://blog.demofox.org/2015/02/25/avoiding-the-performance-hazzards-of-stdfunction/.

[14] Oracle, "Arrays (Java Platform SE 8)," March 2014. [Online]. Available: https://docs.oracle.com/javase/8/docs/api/java/util/Arrays.html.

[15] P. Lammich, "Efficient Verified Implementation of Introsort and Pdqsort," in *10th International Joint Conference on Automated Reasoning (IJCAR 2020)*, 2020.

[16] "orlp/pdqsort: Pattern-defeating quicksort," 2021. [Online]. Available: https://github.com/orlp/pdqsort. [Accessed 24 October 2025].

[17] cppreference.com, "std::sort," October 2025. [Online]. Available: https://en.cppreference.com/w/cpp/algorithm/sort.html.

[18] cppreference.com, "std::stable_sort," October 2025. [Online]. Available: https://en.cppreference.com/w/cpp/algorithm/stable_sort.html.

[19] cppreference.com, "std::partial_sort," October 2025. [Online]. Available: https://en.cppreference.com/w/cpp/algorithm/partial_sort.html. [Accessed 24 October 2025].

[20] M. E. Nebel, S. Wild and C. Martínez, "Analysis of Pivot Sampling in Dual-Pivot Quicksort: A Holistic Analysis of Yaroslavskiy's Partitioning Scheme," *Algorithmica,* vol. 80, no. 2, p. 790–848, June 2016.

[21] . Wild, "Dual-Pivot Quicksort and Beyond: Analysis of Multiway Partitioning," Kaiserslautern, Germany, 2016.

[22] . Martínez, . Nebel and . Wild, "Sesquickselect: One and a Half Pivots for Cache-Efficient Selection," *ACM Transactions on Algorithms,* vol. 15, no. 4, p. Article no. 47 (33 pages), January 2019.

[23] cppreference.com, "constexpr specifier," October 2025. [Online]. Available: https://en.cppreference.com/w/cpp/language/constexpr.

[24] Apple Inc., "Explore the New System Architecture of Apple Silicon Macs," June 2020. [Online]. Available: https://developer.apple.com/videos/play/wwdc2020/10686/. [Accessed 24 October 2025].

[25] Intel Corporation, "How Intel® Core™ Processors Work: Hybrid Architecture Design," October 2023. [Online]. Available: https://www.intel.com/content/www/us/en/gaming/resources/how-hybrid-design-works.html. [Accessed 24 October 2025].

[26] Innovative Computing Laboratory, "PAPI User's Guide," 2023. [Online]. Available: https://icl.utk.edu/projects/papi/files/documentation/PAPI_USER_GUIDE_23.htm. [Accessed 24 October 2025].

[27] Apple Inc., "Optimize CPU Performance with Instruments," June 2025. [Online]. Available: https://developer.apple.com/videos/play/wwdc2025/308/.

[28] cppreference.com, "std::chrono::high_resolution_clock," December 2024. [Online]. Available: https://en.cppreference.com/w/cpp/chrono/high_resolution_clock.

[29] Orson R. L. Peters, "Pattern-defeating Quicksort," 2021.

[30] A. Inc., "C++ Language Support – Xcode," [Online]. Available: https://developer.apple.com/xcode/cpp/. [Accessed 24 October 2025].

[31] . Free Software Foundation, "GCC 14 Release Series – GNU Project," 23 May 2025. [Online]. Available: https://gcc.gnu.org/gcc-14/. [Accessed 24 October 2025].

[32] . Kitware, "CMake 4.0 Release Notes," [Online]. Available: https://cmake.org/cmake/help/latest/release/4.0.html. [Accessed 24 October 2025].