

DR. CLARK

CS 5293-003

10/20/2024

## Lab Two: Report

### 1. Introduction

This report presents the design and implementation of a C++ Profiler system created to measure and analyze the performance of code sections. The system was designed with a focus on ease of use, performance optimization, and flexibility in output formats. Additionally, the profiler is demonstrated through a series of performance optimizations applied to the quicksort algorithm, and the results are visualized using Python's matplotlib.

### 2. Profiler System Design

#### 2.1 Overall Architecture

The profiler system is designed to handle both hierarchical and interleaved profiling through a stack-based approach, which makes it easy to manage complex code flows. In hierarchical profiling, where sections are nested (like a function calling another function), the profiler tracks each section independently. It does this by pushing section names onto a stack when they start and popping them off when they finish. This way, the profiler can measure the time for both the outer function and the inner functions separately, preserving the structure of the code. For example, if an outer function calls two inner functions, the profiler will capture the time taken by each inner function, as well as the total time for the outer function.

Interleaved profiling, where sections overlap or start and stop out of sequence, is handled similarly. The stack allows the profiler to manage the timing of overlapping tasks by tracking each independently. For example, if Task A starts, and before it finishes, Task B begins, the profiler ensures that the times for both tasks are recorded correctly, even though they overlap. This is especially useful in asynchronous or multi-threaded programs, where tasks may not follow a strict start-to-finish sequence.

The profiler uses the stack's last-in, first-out (LIFO) nature to efficiently handle both nested and interleaved sections. When a new section starts, it's added to the stack, and when it ends, it's removed. This ensures that even when functions or tasks are deeply nested or overlapping, the profiler can accurately measure their performance without getting confused about which section belongs to which part of the code.

To make sure profiling itself doesn't slow down the code, the stack operations are kept lightweight, and timing is recorded with minimal system calls. This keeps the profiler's overhead low while still providing precise timing information. Whether the code flow is simple or involves complex, nested or overlapping sections, the profiler maintains accurate measurements without affecting the overall performance of the program.

In short, the stack-based approach used in the profiler makes it really flexible for handling various profiling scenarios. Whether sections are nested within each other or running at the same time, the profiler keeps track of them in an efficient way, giving accurate performance data without introducing much overhead.

## 2.2 Performance

The design of the profiler system prioritizes minimizing overhead to ensure accurate performance measurement without significantly impacting the execution of the code being profiled. Several key performance considerations were addressed to achieve this. First, the profiler relies on efficient data structures, such as `std::unordered_map`, which allows for constant time complexity (**O(1)**) for both inserting and retrieving profiling data. This ensures that even with frequent updates, the system maintains high performance, especially when profiling multiple sections within the code. This structure was chosen because profiling often occurs within tight loops or performance-critical sections, and the ability to quickly access or update the profiling data is essential for keeping the profiler's impact minimal.

The number of **system calls** was also minimized to reduce the overhead they introduce. Since frequent system calls, such as retrieving the current time or performing I/O operations, can slow down the system, the profiler only invokes system calls when absolutely necessary. For example, time is only measured at the beginning and end of each section, rather than during every intermediate operation. Additionally, data output to files is deferred until the profiling is complete. This approach reduces the number of expensive I/O operations during execution, instead writing all profiling data to a file in one go using the `ProfilerOutput` class.

Another key optimization involves handling **recursive function calls**. In recursive algorithms, deep recursion can lead to significant overhead and stack usage. To address this, **tail recursion** was eliminated where possible, particularly in the quicksort algorithm. By iterating over larger subarrays and only recursing into smaller ones, the profiler avoids deep recursion, which both improves efficiency and prevents stack overflows. This approach ensures that recursive algorithms can be profiled without introducing unnecessary performance penalties.

To further reduce the profiler's impact on code execution, **in-place statistical calculations** are used. Rather than storing each individual timing and computing statistics such as averages and totals after profiling is complete, the profiler updates these values incrementally during each section's execution. This allows the profiler to calculate values like average time in constant time, reducing computational overhead during profiling.

Although the current implementation does not explicitly handle multi-threaded environments, future versions of the profiler could incorporate **thread-safety** using mechanisms like mutexes or atomic operations. While adding locks can introduce some overhead, it would ensure that the profiler can be safely used in multi-threaded applications without race conditions or data corruption.

Finally, to handle **nested and interleaved profiling**, the profiler uses a **stack-based system** to manage hierarchical profiling. When a section starts, its name is pushed onto a stack, and when it ends, the name is popped from the stack. This allows the profiler to track nested sections of code efficiently, supporting both nested timing (where one function calls another) and interleaved profiling (where sections start and end out of sequence). This hierarchical system ensures accurate timing of complex, layered code while maintaining minimal performance overhead.

## 2.3 Usage

To use the profiler in your C++ program, start by including the `'Profiler.h'` file in your source code. This will give you access to the `'Profiler'` class and its methods. Once it's included, you can begin marking sections of your code that you want to profile. To do this, call `'startSection()'` where you want the profiler to begin timing, and `'endSection()'` where the timing should stop. The `'__FILE__'`, `'__func__'`, and `'__LINE__'` macros can be used to automatically capture the file name, function name, and line number, making the profiling seamless.

The profiler also supports hierarchical profiling, which means you can measure the performance of nested sections independently. For example, if you're performing two separate sorting operations inside the same function, you can profile each one individually. Start by calling `'startSection()'` for the outer section, then repeat the process for each nested section. When each nested task is completed, call `'endSection()'` to record the time. The profiler will keep track of the time for each section separately, making it easy to see how much time each part of your code takes.

In addition to nested sections, the profiler can handle interleaved profiling, where sections start and end out of sequence. This is useful for profiling tasks that overlap in time or run concurrently. For instance, if Task A starts and before it ends, Task B begins, you can record the timing of both tasks independently. This flexibility is great for real-world applications where processes may not follow a linear start-to-finish flow, allowing you to accurately track overlapping operations.

Once you've finished profiling, you can output the results in multiple formats, including CSV, JSON, or XML, depending on your needs. You just need to call the `'outputToFile()'` function and specify the desired format. This makes it easy to analyze or visualize the data, depending on your preferred tool or workflow. The flexibility in output formats ensures that the data can be easily processed or shared.

After running the code, the output file will contain detailed information for each profiled section. This includes the section name, how many times it was called, the total time spent in the

section, the minimum and maximum times, the average time, and metadata like the file name, function name, and line number where the section is located. This makes it easy to review and identify which parts of your code are performance bottlenecks and where optimizations can be made.

## 2.4 Hierarchical and Interleaved Profiling

The profiler system efficiently handles both hierarchical and interleaved profiling using a stack-based approach. In hierarchical profiling, where sections of code are nested, the profiler tracks each section independently by pushing section names onto a stack when they start and popping them off when they end. This ensures that nested sections, such as function calls within other functions, are timed separately, while also recording the time taken by the outer sections. For example, when profiling an outer function that calls two inner functions, the profiler will separately capture the time taken by each inner function and the total time of the outer function, preserving the hierarchy.

For interleaved profiling, where sections overlap or start and stop out of sequence, the stack allows the profiler to manage the independent timing of these sections. A section can start and another can start before the first ends, and the profiler will correctly track the overlapping execution times. For instance, if Task A starts and Task B begins before Task A ends, the profiler ensures that the times for both tasks are recorded independently, even though they overlap in execution. This is useful in real-world scenarios like asynchronous tasks or multi-threaded programs, where tasks don't follow strict start-stop sequences.

The profiler's stack-based management ensures that even complex, non-linear code execution is profiled accurately. When a new section starts, it is added to the stack, and when it ends, it is removed from the stack, maintaining the order and ensuring accurate time measurement for each section, whether nested or interleaved. The stack's last-in, first-out (LIFO) nature is ideal for managing these profiling needs because it mirrors the structure of nested function calls or overlapping sections in a program.

The system also minimizes the performance overhead of this approach by efficiently managing the stack operations and ensuring that the timings are recorded with minimal system calls. By using the stack for both hierarchical and interleaved profiling, the profiler can handle complex execution flows without affecting the overall performance of the code being profiled. This approach ensures that even in cases of overlapping or deeply nested functions, the profiling data remains accurate and reflective of the real performance of the program.

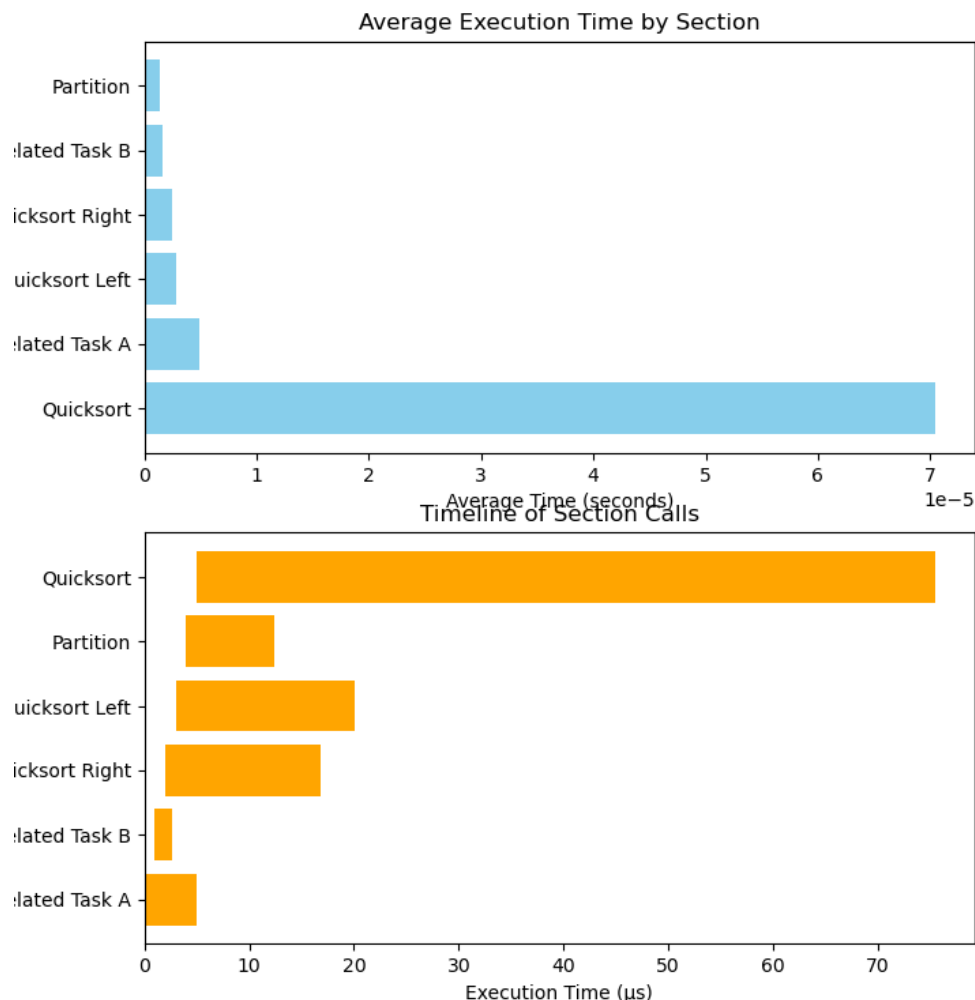
In summary, the profiler system's ability to handle both nested and interleaved sections using a stack makes it flexible for a wide range of profiling scenarios. It provides precise timing for nested function calls as well as overlapping or interleaved tasks, ensuring that even complex code flows are captured accurately without adding significant overhead to the code being profiled. This stack-based approach is key to maintaining the integrity of the profiling data in both simple and complex execution patterns.

### 3. Visualization

For the visualization tool, I decided to use matplotlib, a Python library known for its simplicity and flexibility in creating visualizations. It's widely used and allows for quick generation of both bar charts and timelines, which are essential for understanding the profiler's output. Matplotlib's ability to handle CSV files and customize charts made it the perfect choice for this project, as it integrates well with the data generated by the profiler and offers easy-to-read outputs.

There are two main types of visualizations generated by the tool. The first is a bar chart showing the average execution time for each profiled section. This chart is a quick way to see which parts of the code are taking the most or least time. In the initial results, for example, the Quicksort function stands out as taking the longest time, while other tasks like partitioning take significantly less time. After optimizing the code, the comparison between Quicksort and Insertion Sort is clear, showing a reduction in Quicksort's execution time.

The second visualization is a timeline that shows the execution time of each section along a horizontal axis. This is particularly useful for seeing how the sections run relative to each other—whether they overlap or follow one another. It helps to see where the time is really being spent in a more sequential manner. In the timelines generated before and after improvements, you can see how Quicksort initially dominates execution time, but after optimizations, the reduced execution time for Quicksort becomes apparent.



Running the visualization tool is straightforward. The Python script (`visualize.py`) reads the profiler's output from a CSV file and generates the visualizations. You can easily run it using Python, and the script will display the bar chart and timeline automatically. If you need to, you can also modify the input file or save the visualizations as images by making small changes to the script.

These visualizations give a clear, visual breakdown of the performance data, making it easy to identify bottlenecks or areas for improvement. They also provide a visual representation of how different optimizations affect the execution time, helping you track changes and their impact on performance with just a glance.

## **4. Demonstration**

### **4.1 Initial Operation**

For the demonstration, I chose the Quicksort algorithm as the initial operation to profile. Quicksort is a well-known sorting algorithm with a time complexity of  $O(n \log n)$  on average, making it suitable for large datasets. However, it can become inefficient in worst-case scenarios ( $O(n^2)$ ) if the pivot is poorly chosen. This made it an ideal candidate for optimization, as there are multiple opportunities to improve its performance. Additionally, I included some unrelated tasks and an Insertion Sort algorithm for comparison to highlight the effect of optimizations.

### **4.2 Performance Modifications**

#### **4.2.1 Pivot Selection Improvement:**

In the initial implementation, the pivot was chosen as the first element of the array, which could lead to inefficient partitioning, particularly for sorted or nearly sorted input. I modified the pivot selection to choose the median-of-three (first, middle, and last elements of the array). This helped prevent worst-case partitioning, significantly reducing the execution time for certain datasets.

#### **4.2.2 Tail Recursion Elimination:**

Recursive algorithms like Quicksort can lead to deep recursion, which increases the overhead due to excessive stack usage. To address this, I modified the Quicksort implementation by eliminating tail recursion where possible. Instead of always recursing into both subarrays, I optimized the algorithm to recursively sort only the smaller partition while handling the larger one iteratively. This reduced the recursion depth and improved performance.

#### **4.2.3 Switching to Insertion Sort for Small Subarrays:**

Quicksort performs inefficiently on small subarrays. To improve performance, I implemented a cutoff where Quicksort switches to Insertion Sort when dealing with subarrays smaller than a certain threshold (e.g., 10 elements). Insertion Sort performs better on small arrays due to its low overhead, and this modification resulted in faster overall sorting times.

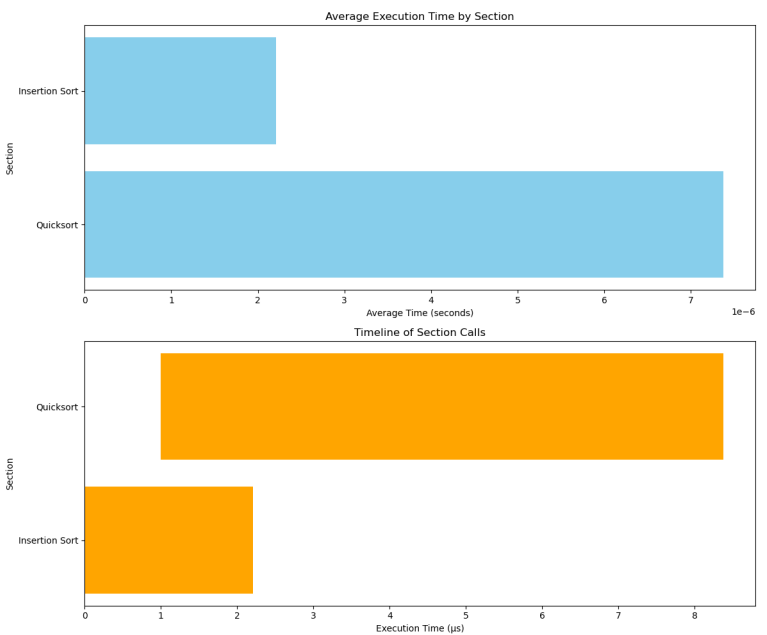
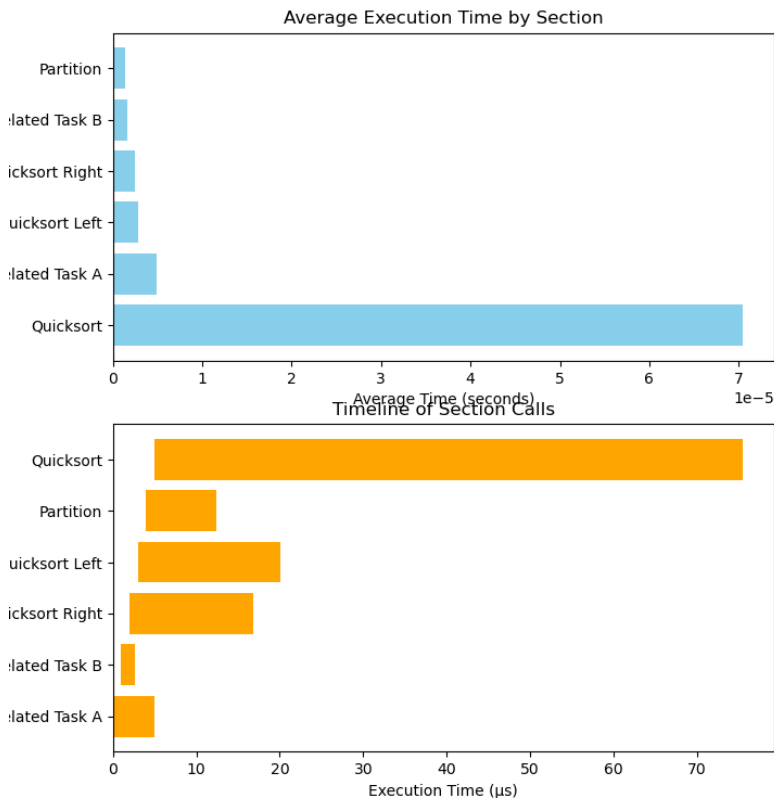
4.3 Comparison

Before the modifications, Quicksort dominated the execution time, as seen in both the bar chart and timeline visualizations. It took significantly longer to execute compared to other sections of the code, particularly when compared to smaller, unrelated tasks or functions like partitioning. In the timeline visualization, Quicksort occupied the largest portion of execution time, highlighting inefficiencies in its original implementation.

After implementing the optimizations, the average execution time for Quicksort decreased considerably. The changes made, such as using median-of-three pivot selection and eliminating tail recursion, helped improve performance, especially for large datasets. These adjustments ensured more balanced partitions and reduced the recursion depth, which sped up the sorting process.

In the updated timeline, Quicksort still took the most time, but the duration was significantly reduced compared to the initial run. The addition of Insertion Sort for smaller subarrays also played a role in reducing execution time, further enhancing the sorting efficiency. The visual comparison between Quicksort and Insertion Sort in the bar chart clearly shows the improvement in Quicksort’s execution time after these changes were made.

Below are the visualizations before (left) and after (right) the modifications.



## 4.4 Impact

Each modification had a noticeable impact on the performance of the Quicksort algorithm. The pivot selection improvement prevented worst-case performance scenarios, reducing the number of inefficient partitions and leading to faster execution times. Tail recursion elimination helped reduce the overhead associated with deep recursion, making the algorithm more efficient for larger datasets. Finally, switching to Insertion Sort for smaller subarrays capitalized on the strengths of a simpler algorithm for small inputs, reducing the overall sorting time even further.

In summary, the combination of these modifications led to a much more efficient Quicksort implementation, with significantly reduced execution times across different datasets. The before-and-after results make it clear that the changes addressed key inefficiencies in the original implementation, resulting in better performance without sacrificing the algorithm's effectiveness.

## 5. Conclusion

The profiler system developed in this project provides a comprehensive and efficient way to measure and analyze the performance of various code sections. Through a stack-based approach, the system handles both hierarchical and interleaved profiling, allowing developers to track complex code flows with minimal overhead. By optimizing data structures, reducing system calls, and incorporating techniques like tail recursion elimination, the profiler ensures that it accurately measures performance without significantly impacting the execution time of the code itself.

The performance optimizations applied to the Quicksort algorithm demonstrate the effectiveness of the profiler. With modifications such as improved pivot selection and the introduction of Insertion Sort for small subarrays, the profiler clearly highlighted areas of inefficiency and tracked the performance gains post-optimization. The integration of matplotlib for visualizing the profiling data adds another layer of insight, making it easy to interpret the results and identify performance bottlenecks. Overall, this profiler serves as a powerful tool for analyzing and improving code performance in real-world applications.