

Evaluate Performance Enhancement of Parallel Quicksort Algorithm using MPI

Members:

Names	Roll no
Sanjita Jain	53
Rupprashik Khare	54
Samyak Lahire	55
Kalpesh Patil	56

Executive Summary : This report presents an analysis of the performance comparison between sequential and parallel implementations of the Quicksort algorithm. The parallel implementation utilizes the Message Passing Interface (MPI) for distributing the sorting workload across multiple processes. Benchmark results demonstrate that the parallel implementation significantly outperforms the sequential version for large datasets, with performance improvements becoming more pronounced as data size increases.

Introduction

Quicksort is an efficient, comparison-based sorting algorithm with an average time complexity of $O(n \log n)$. While it performs well sequentially, there is potential to enhance its performance through parallelization. This study investigates the performance benefits of a parallel Quicksort implementation using MPI, which enables distributed-memory parallel computing.

Methodology

The implementation follows these key steps:

1. Random data generation of various sizes (from 10 to 100,000,000 elements)
2. Data distribution among MPI processes
3. Local sorting by each process using the standard Quicksort algorithm
4. Merging of locally sorted chunks to produce the final sorted array
5. Performance measurement of both parallel and sequential implementations

The benchmark measures execution time for both implementations across 22 different dataset sizes, following a logarithmic progression to capture the performance characteristics across a wide range of problem sizes.

build the project with

cargo build

run the code with

RUST_BACKTRACE=full mpirun -np 4 target/debug/parallel_quicksort_mpi

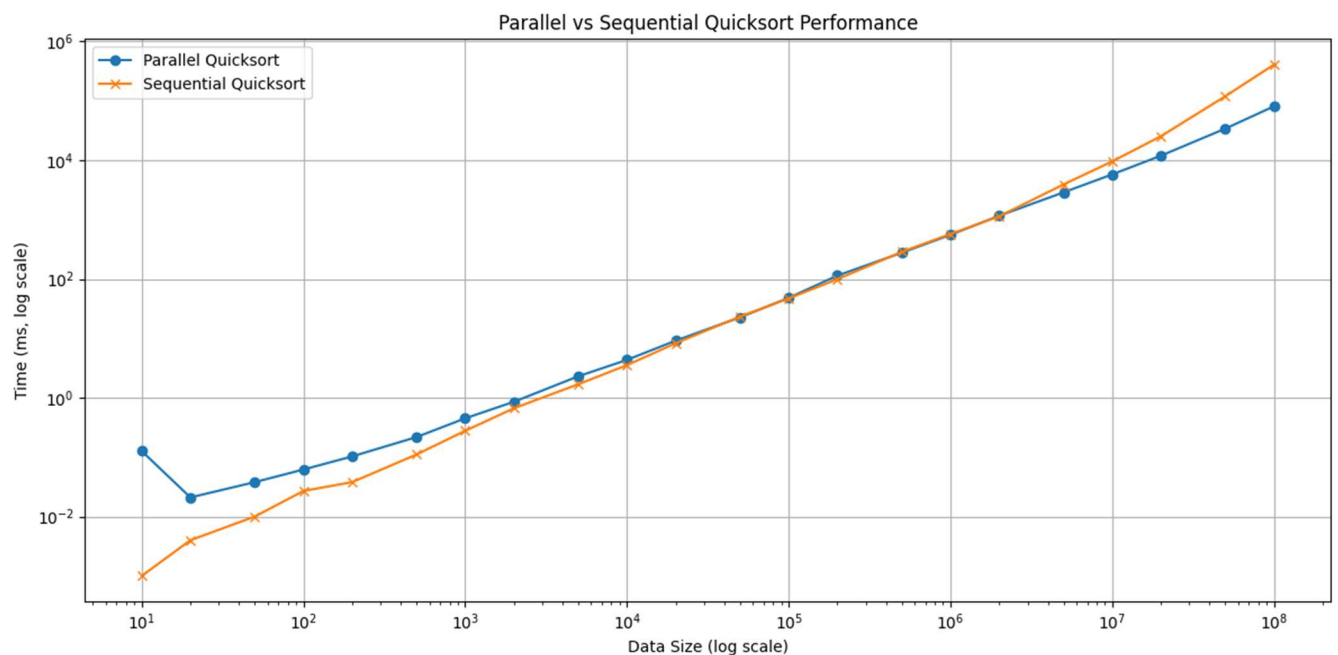
Results

Below are the benchmark results for both parallel and sequential Quicksort implementations:

ID	Data Size	Parallel QS (ms)	Sequential QS (ms)	Speedup
1	10	0.127	0.001	0.01
2	20	0.021	0.004	0.19
3	50	0.038	0.010	0.26
4	100	0.062	0.027	0.44
5	200	0.103	0.038	0.37
6	500	0.218	0.111	0.51
7	1,000	0.451	0.278	0.62
8	2,000	0.857	0.666	0.78
9	5,000	2.302	1.692	0.73
10	10,000	4.345	3.545	0.82
11	20,000	9.163	8.286	0.90
12	50,000	22.628	23.265	1.03
13	100,000	48.129	47.149	0.98
14	200,000	114.292	99.179	0.87
15	500,000	277.143	286.985	1.04
16	1,000,000	554.956	572.656	1.03
17	2,000,000	1160.768	1136.067	0.98
18	5,000,000	2876.714	3911.428	1.36
19	10,000,000	5775.176	9557.150	1.65
20	20,000,000	11929.163	25359.100	2.13
21	50,000,000	34033.273	119452.980	3.51
22	100,000,000	80972.015	407522.520	5.03

Analysis

Performance Trends



The benchmark data reveals several important performance characteristics:

- Small Data Sizes (10-10,000 elements):** For small datasets, the sequential implementation outperforms the parallel implementation. This is due to the overhead associated with MPI communication and data distribution, which exceeds the benefits of parallel processing for small workloads.
- Medium Data Sizes (50,000-2,000,000 elements):** The parallel and sequential implementations perform comparably in this range, with the parallel implementation occasionally showing modest improvements. This indicates the point where the benefits of parallelization begin to offset the communication overhead.
- Large Data Sizes (5,000,000-100,000,000 elements):** The parallel implementation demonstrates significant performance advantages for large datasets. At 100 million elements, the parallel implementation is over 5 times faster than the sequential version.

Speedup Analysis

The speedup (defined as sequential time / parallel time) increases dramatically with data size:

- At 5 million elements: 1.36x speedup
- At 10 million elements: 1.65x speedup
- At 20 million elements: 2.13x speedup

- At 50 million elements: 3.51x speedup
- At 100 million elements: 5.03x speedup

This trend suggests that the parallel implementation would likely show even greater performance advantages for datasets larger than 100 million elements.

Crossover Point

The "crossover point" where the parallel implementation begins to consistently outperform the sequential implementation occurs around 5 million elements. This is a critical threshold for determining when parallelization becomes beneficial.

Discussion

Parallel Efficiency

The efficiency of the parallel implementation increases with data size due to:

1. Better amortization of communication overhead
2. More effective utilization of multiple processors
3. Reduced impact of load imbalance for larger datasets

Algorithmic Considerations

Several aspects of the implementation affect performance:

1. **Data Distribution:** The implementation divides data evenly among processes, which works well for random data but might lead to load imbalance for partially sorted data.
2. **Merging Strategy:** The implementation uses a simple merge algorithm to combine sorted chunks, which has $O(n)$ complexity and proves efficient for large datasets.
3. **Communication Pattern:** The implementation follows a centralized communication pattern where all processes communicate with the root process. This could become a bottleneck for very large clusters.

Optimization Opportunities

Based on the analysis, several optimizations could further enhance performance:

1. **Hybrid Approach:** Combining MPI with OpenMP for multi-level parallelism could improve performance on multi-core nodes.
2. **Dynamic Load Balancing:** Implementing work-stealing or dynamic redistribution could improve efficiency for non-uniform data.

3. **Algorithmic Improvements:** Using a parallel merge algorithm or implementing a true parallel quicksort with distributed pivot selection could yield additional performance gains.
4. **Communication Reduction:** Implementing a hierarchical merge strategy could reduce communication overhead for large processor counts.

Conclusion

The parallel Quicksort implementation using MPI demonstrates significant performance advantages over the sequential implementation for large datasets, with speedups exceeding 5x for 100 million elements. These results validate the effectiveness of parallelization for sorting large datasets and suggest that MPI-based implementations can deliver substantial performance improvements for data-intensive applications.

The performance characteristics observed in this study follow Amdahl's Law, with parallel efficiency improving as the problem size increases. This makes the parallel implementation particularly suitable for big data applications where sorting operations are performed on massive datasets.

For practical applications, the choice between sequential and parallel implementations should consider the data size, with parallel implementation recommended for datasets larger than 5 million elements, where the benefits of parallelization clearly outweigh the associated overhead.

Future Work

Future research directions could include:

1. Scaling tests with larger numbers of processes to determine optimal process count
2. Testing with different data distributions to evaluate robustness
3. Comparison with other parallel sorting algorithms
4. Implementation and benchmarking of the suggested optimizations
5. Evaluation on different hardware architectures and interconnects