



University of Pisa

Department Of Computer Science

Parallel and Distributed Systems: Paradigms and Models

Project:
Distributed out-of-core MergeSort

Student:

Abderrahmane Salmi

704608

a.salmi1@studenti.unipi.it

Professor:

Massimo Torquati

ACADEMIC YEAR 2024/2025

Chapter 1

Single Node Solution (OpenMP and FastFlow)

1.1 Introduction

This project is focused on the design and implementation of a scalable, distributed, and out-of-core MergeSort for very large datasets that exceed the capacity of main memory. The sorting is performed on files that contain variable-sized records, each one contains a key (used for comparison) and a binary payload. The primary objective is to build a scalable MergeSort that works both on a single node (using shared-memory parallelism) and across multiple nodes (using distributed-memory models).

1.2 Shared-Memory Sorting with OpenMP

1.2.1 Overview

Given the memory limitation (e.g., $\leq 32\text{GB}$ RAM), the original input file is split into smaller chunks that can fit into memory, then these chunks are sorted independently in parallel using OpenMP, and finally merged into a single sorted output using a K-way merge.

1.2.2 Implementation Details

1. Chunking Phase:

The input file is first memory-mapped using `mmap()` for fast access. Then, the file is scanned sequentially to generate logical chunks (offsets), which are carefully computed so that record boundaries are not broken (all chunks contain complete records). The estimated chunk sizes are computed based on the total input size and available memory budget. Finally, the logical chunks are given to OpenMP threads where each thread creates a file for each chunk, which is done in parallel.

2. Sorting Phase:

Each chunk is read into memory, parsed into Record objects, and sorted using `std::sort` with a custom comparator based on the record key. This phase is parallelized using OpenMP (`#pragma omp parallel for`), where each thread processes one or more chunks.

After sorting, each sorted chunk is written to a temporary file to be used later in the merge.

3. Merging Phase:

Once all chunks are sorted, a k-way merge is performed to combine them into a single sorted output file. This phase is implemented using a min-heap (priority queue), where one record is read from each chunk and the smallest is selected and written to the output file. The process repeats until all chunks are handled and then the temporary chunk files are deleted (cleaned up).

PS: The entire sorting function is timed using "omp_get_wtime()" for accurate performance measurement. PS: To make sure the sorting is correct, a separate C++ utility code is used to verify that the output file is sorted correctly.

1.2.3 Challenges Encountered

Trying to make the chunking and merging phases parallel and efficient was one of the most difficult parts of the implementation. Different strategies were considered, and some were tried (tested) and discarded due to performance issues or implementation complexity.

Chunking: For example, for the chunking, two parallelization strategies were considered:

- **Strategy A: Chunking before sorting**

First, a sequential pass scans the file and computes all logical chunks. Then, multiple threads write each chunk to disk in parallel, then when all chunks are written, they are passed to the sorting function. This approach was effective and simple, and it was the one used in the project.

- **Strategy B: Chunking and sorting at the same time**

Same start with a sequential scan of the file and computeing all logical chunks, then each thread performs chunking and immediately launches an OpenMP task to sort the chunk, this way we don't wait for all chunks, each ready chunk is sorted immediatly. Although this idea is great in theory, in practice it introduced more overhead and was slower than the first one.

Merging:

The merging phase also posed a design challenge. The goal was to reduce total merge time by parallelizing it. I tried implementing a multi-stage parallel merge, where groups of chunk files were merged in parallel stages, and the intermediate files were merged again until only one final file remained.

While this would work well when the number of chunks is very large, it introduced too much I/O overhead in our project and test files. At the end, the multi-stage merging became slower than the simple sequential k-way merge.

1.3 Shared-Memory Sorting with FastFlow

1.3.1 Overview

Just like OpenMP, This solution implements a parallel external merge sort on large binary files that cannot fit entirely in memory, but it uses the FastFlow framework to do so.

1.3.2 Parallel Architecture: Farm Pattern

Our FastFlow implementation is developed using the farm pattern, which consists of:

- **Emitter:** Responsible for reading and splitting the input file into smaller chunks. Just like we did in OpenMP, we scan the file to determin the logical chunks and make sure we don't divide any records in the middle, then we pass those chunks to the workers.
- **Workers:** A pool of parallel workers that sort each chunk independently. Each worker receives a chunk, parses it into records, sorts the records in-memory using `std::sort`, and writes the sorted data to its temporary output file.
- **Collector:** Recieves the sorted chunks from the workers and performs a final k-way merge to produce a single output file. Temporary files are then deleted after the merge is complete.

PS: The total time for the sorting and merging phases is measured using the fastflow utility time functions for better performance accuracy.

Chapter 2

Multi-Node Hybrid Solution (MPI + FastFlow)

2.1 Overview

This section describes the implementation of the distributed sorting algorithm using MPI (Message Passing Interface) for inter-process parallelism and FastFlow for intra-process parallelism. The goal is always to efficiently sort large files that exceed the memory capacity of a single machine, using multiple nodes or cores.

The solution has two modes:

- **Single-node mode** (when MPI size = 1): the file is sorted locally on the only node available using the FastFlow sorter.
- **Distributed mode** (MPI size > 1): rank 0 becomes the coordinator, and all other ranks act as sorting workers.

2.2 Implementation Details

My solution was inspired from the farm pattern like we did in FastFlow, where we have 2 special nodes that are the emitter and the collector, and the rest are workers. Same logic applies here, rank 0 is the coordinator, which is the equivalent of both the emitter and collector (does both their jobs), and the rest of the ranks are workers.

Rank 0 - Coordinator:

- Logically splits the input file into chunks, making sure the record boundaries are preserved. Same chunking logic that was used in FastFlow and OpenMP is also used here.
- Writes each logical chunk to a physical file on disk.
- Distributes chunk files to workers in a round-robin way.
- Collects the paths to the sorted chunk files from each worker.
- Performs a final k-way merge to produce the fully sorted output.

- Deletes all temporary files after merging.

Other Ranks - Workers:

- Receives the list of chunk files assigned by the coordinator.
- For each file, uses the existing FastFlow code to sort it.
- Sends back to the coordinator the paths to the sorted output files.

PS: The MPI sorting solution is timed using `MPI_Wtime()` for more accurate performance measuring.

Chapter 3

Performance Evaluation

3.1 Overview

This section presents a detailed analysis of the performance of our parallel sorting implementations using FastFlow and OpenMP. The goal is to evaluate how the solutions scale when we change:

- Number of threads/workers (1, 2, 4, 8, 16, 32, 64)
- Input size (number of records: 10M, 50M)
- Payload size (64B, 512B)

3.2 Datasets Used

Filename	Records (N)	Payload (Bytes)	Size	Category
data_10M_p64.bin	10M	64	~725 MB	Small N, Small Payload
data_10M_p512.bin	10M	512	~4.9 GB	Small N, Large Payload
data_50M_p64.bin	50M	64	~3.6 GB	Large N, Small Payload
data_50M_p512.bin	50M	512	~25 GB	Large N, Large Payload

3.3 Methodology

Each dataset was sorted using both the FastFlow and OpenMP implementations. We measured the execution time for each run and computed the speedup and efficiency based on the single-thread result.

$$S(p) = \frac{T_{\text{seq}}}{T_{\text{par}}(p)} \quad (3.1)$$

$$E(p) = \frac{S(p)}{p} = \frac{T_{\text{seq}}}{T_{\text{par}}(p) \times p} \quad (3.2)$$

Where:

- $S(p)$: Speedup
- $E(p)$: Efficiency
- T_{seq} : Execution time of sequential version
- $T_{par}(p)$: Execution time with p threads/workers

3.4 Execution Time

Sequential Execution Time (s)	
Dataset	Time (s)
data_10M_p64	44.7748
data_10M_p512	78.5549
data_50M_p64	385.493
data_50M_p512	1171.51

Table 3.1: Sequential execution times (in seconds) for each dataset.

FastFlow Execution Time (s)							
Dataset	1W	2W	4W	8W	16W	32W	64W
data_10M_p64	12.76	6.99	5.24	4.53	4.91	5.13	5.37
data_10M_p512	69.76	27.99	18.61	15.73	14.73	13.27	28.68
data_50M_p64	126.61	94.45	40.88	29.61	24.84	27.87	28.29
data_50M_p512	994.54	404.19	388.39	212.45	256.12	249.03	360.53

Table 3.2: FastFlow execution times (in seconds) for varying number of workers.

OpenMP Execution Time (s)							
Dataset	1T	2T	4T	8T	16T	32T	64T
data_10M_p64	45.65	27.72	23.63	15.84	15.53	14.92	16.89
data_10M_p512	89.54	70.04	64.20	76.06	93.67	87.77	110.38
data_50M_p64	470.21	170.39	94.30	69.42	71.36	83.84	102.73
data_50M_p512	1304.60	466.79	157.40	104.17	168.32	466.12	462.60

Table 3.3: Execution times (in seconds) for OpenMP for varying number threads

3.5 Speedup and Efficiency Analysis

Table 3.4: Speedup Comparison

Threads	OMP	FF	Speedup
data_10M_p512			
1	0.88	1.13	0.88
2	1.12	2.81	1.12
4	1.22	4.22	1.22
8	1.03	4.99	1.03
16	0.84	5.33	0.84
32	0.89	5.92	0.89
64	0.71	2.74	0.71
data_10M_p64			
1	0.98	3.51	0.98
2	1.62	6.41	1.62
4	1.89	8.54	1.89
8	2.83	9.87	2.83
16	2.88	9.12	2.88
32	3.00	8.73	3.00
64	2.65	8.34	2.65
data_50M_p512			
1	0.90	1.18	0.90
2	2.51	2.90	2.51
4	7.44	3.02	7.44
8	11.25	5.51	11.25
16	6.96	4.57	6.96
32	2.51	4.70	2.51
64	2.53	3.25	2.53
data_50M_p64			
1	0.82	3.04	0.82
2	2.26	4.08	2.26
4	4.09	9.43	4.09
8	5.55	13.02	5.55
16	5.40	15.52	5.40
32	4.60	13.83	4.60
64	3.75	13.63	3.75

Table 3.5: Efficiency Comparison

Threads	OMP	FF	Efficiency
<i>data_10M_p512</i>			
1	0.88	1.13	0.88
2	0.56	1.40	0.56
4	0.31	1.06	0.31
8	0.13	0.62	0.13
16	0.05	0.33	0.05
32	0.03	0.18	0.03
64	0.01	0.04	0.01
<i>data_10M_p64</i>			
1	0.98	3.51	0.98
2	0.81	3.20	0.81
4	0.47	2.14	0.47
8	0.35	1.23	0.35
16	0.18	0.57	0.18
32	0.09	0.27	0.09
64	0.04	0.13	0.04
<i>data_50M_p512</i>			
1	0.90	1.18	0.90
2	1.25	1.45	1.25
4	1.86	0.75	1.86
8	1.41	0.69	1.41
16	0.43	0.29	0.43
32	0.08	0.15	0.08
64	0.04	0.05	0.04
<i>data_50M_p64</i>			
1	0.82	3.04	0.82
2	1.13	2.04	1.13
4	1.02	2.36	1.02
8	0.69	1.63	0.69
16	0.34	0.97	0.34
32	0.14	0.43	0.14
64	0.06	0.21	0.06

The following plots summarize the performance across all datasets for both fastflow and openmp:

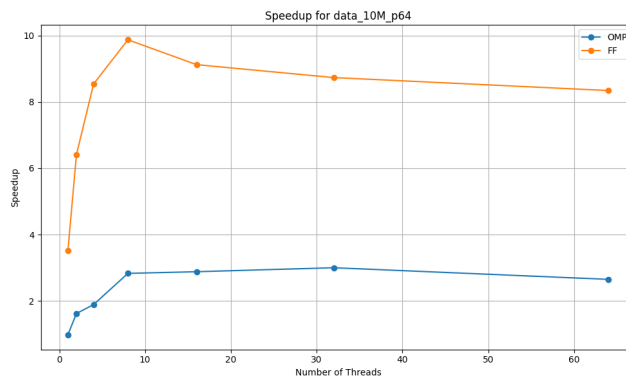


Figure 3.1: Speedup - 10M Records, 64B Payload

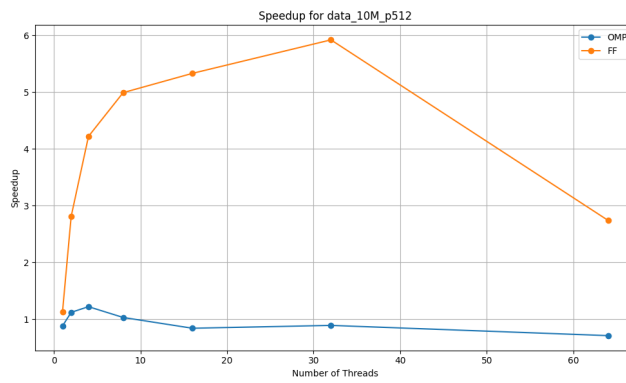


Figure 3.2: Speedup - 10M Records, 512B Payload

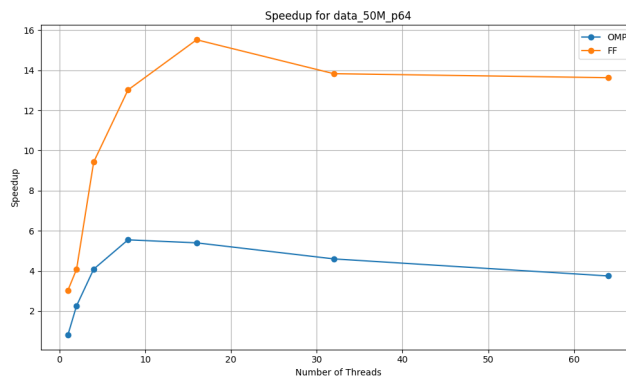


Figure 3.3: Speedup - 50M Records, 64B Payload

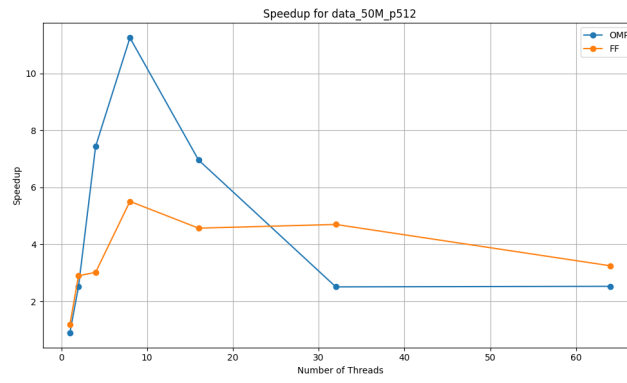


Figure 3.4: Speedup - 50M Records, 512B Payload

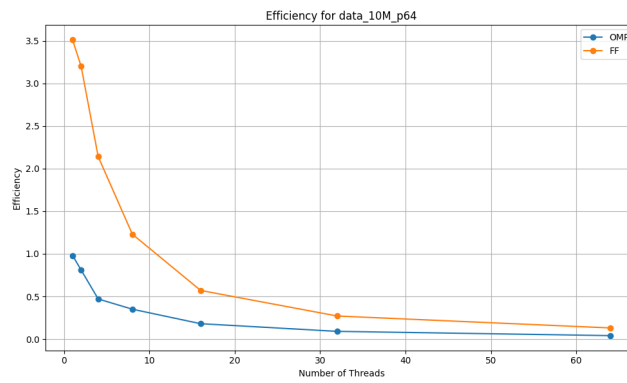


Figure 3.5: Efficiency - 10M Records, 64B Payload

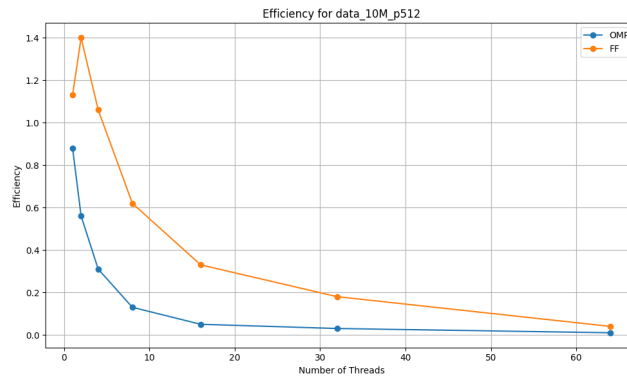


Figure 3.6: Efficiency - 10M Records, 512B Payload

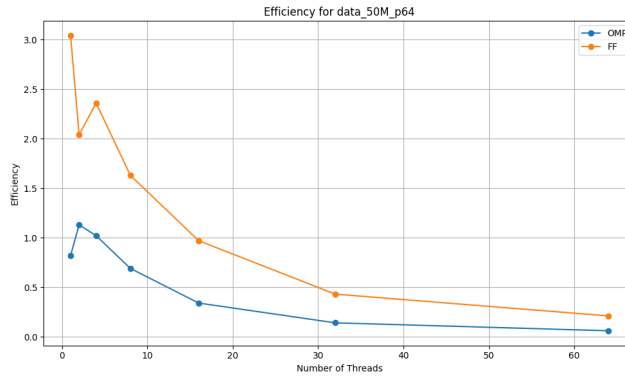


Figure 3.7: Efficiency - 50M Records, 64B Payload

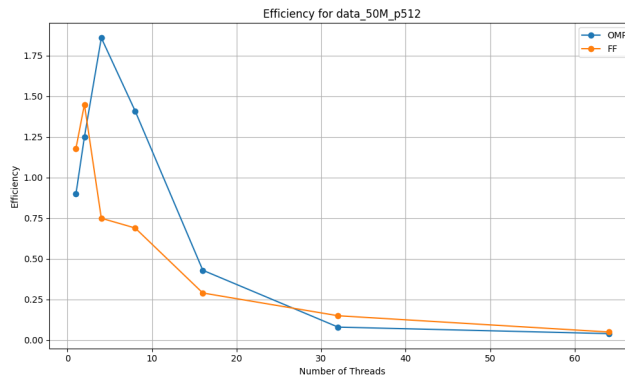


Figure 3.8: Efficiency - 50M Records, 512B Payload

3.6 Observations

Both FastFlow and OpenMP show significant speedup when increasing the number of workers/threads. This indicates that the parallelization strategy is effective for these sizes, and there is a good CPU-bound parallel scalability. However, sometimes as the number of workers/threads increases more, the speedup decreases, which means that the overhead of managing more threads/workers starts to outweigh the benefits of parallel execution. A possible reason for this decrease can be the increasing influence of I/O operations and memory bandwidth limitations.

Another important observation is that the payload size significantly affects performance. For 64-byte records, sorting is more compute-bound, and both systems scale well up. On the other hand for 512-byte records, the system becomes more memory and I/O bound, because read and write operations to disk dominate, especially during the chunking and merging phases.

3.7 MPI Speedup and Efficiency

Table 3.6: MPI Execution Time, Speedup, and Efficiency

Dataset	MPI Procs	Time (s)	Speedup	Efficiency
data_10M_p64	1	18.07	2.48	2.48
data_10M_p64	2	16.94	2.64	1.32
data_10M_p64	4	11.97	3.74	0.94
data_10M_p64	8	13.44	3.33	0.42
data_10M_p512	1	64.51	1.22	1.22
data_10M_p512	2	47.86	1.64	0.82
data_10M_p512	4	30.30	2.59	0.65
data_10M_p512	8	21.71	3.62	0.45
data_50M_p64	1	30.37	12.69	12.69
data_50M_p64	2	53.77	7.17	3.58
data_50M_p64	4	49.33	7.81	1.95
data_50M_p64	8	54.25	7.11	0.89
data_50M_p512	1	164.15	7.14	7.14
data_50M_p512	2	256.21	4.57	2.29
data_50M_p512	4	-	-	-
data_50M_p512	8	-	-	-

Note: Data for 4 and 8 processors for the *data_50M_p512* dataset were not available because it was hard to access the SPM cluster for that long since the file is too big.

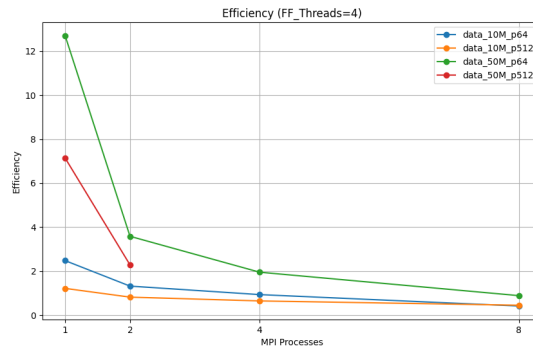


Figure 3.9: Efficiency - All Files

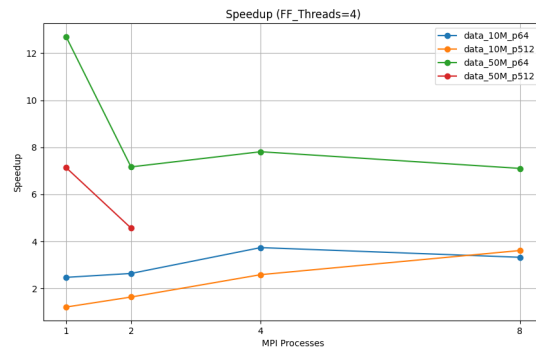


Figure 3.10: Speedup - All Files

Chapter 4

Cost Model and Analysis

4.1 Cost Model Approximation

The distributed external merge sort can be modeled using a combination of the **BSP** and **Farm** cost models that we saw in the course, because they represent well the hybrid architecture of MPI + FastFlow.

Let:

- p = number of MPI processes (1 coordinator + $p - 1$ workers)
- t = number of FastFlow threads per worker
- n = total number of records
- k = number of chunks
- M = total memory budget
- g = communication cost per word
- l = MPI synchronization latency

We can estimate the total cost as:

$$T_{\text{total}} = T_{\text{chunking}} + T_{\text{distribution}} + T_{\text{local-sort}} + T_{\text{merge}} + T_{\text{cleanup}}$$

1. Chunking (Coordinator only):

$$T_{\text{chunking}} = O(n) \quad (\text{sequential scan} + \text{disk I/O})$$

2. Distribution (MPI communication):

$$T_{\text{scatter}} = k \cdot (t_0 + d \cdot s)$$

where t_0 is startup latency, d is size of data, s is per-byte transmission cost ($\approx 1/B$ where B is the bandwidth)

This is inspired from the simple communication cost model seen in the course.

In our implementation, the coordinator sends just file paths, so T_{scatter} should be very small.

3. Local Sort (Workers):

Since each worker is sorting using FastFlow, and our FastFlow is using the farm pattern, so T_{sort} is just $T_c^{\text{farm}}(\text{worker_chunks}, k)$

$$T_{\text{sort}} = T_c^{\text{farm}}(\text{worker_chunks}, k) = (t + 1) \cdot T_{\text{comm}} + (\text{worker_chunks}/t) \cdot (T_s^w + T_{\text{comm}})$$

where t is the number of FastFlow workers (threads), worker_chunks is the number of chunks a rank worker has to process, and T_s^w is worker service time, which is $\max(T_{\text{seq}}, T_{\text{comm}})$

4. Merge (Coordinator):

Let k be the number of sorted chunks. The final k-way merge at the coordinator is:

$$T_{\text{merge}} = O(n \cdot \log k) + O(\text{I/O})$$

5. Cleanup:

Removing temporary files is done sequentially and its cost is negligible in most scenarios.

4.2 Bottleneck Phases

The most notable bottlenecks are:

- **Chunking:** Since its only performed by the coordinator (rank 0), it can be a performance limiter on very large files.
- **Final Merge:** Same as chunking, since it's only handled by rank 0, it can be a bottleneck. Although we considered a multi-stage parallel merge, but it introduced too much I/O and coordination overhead for this specific project.
- **File I/O:** I/O in general always remains a bottleneck especially when the memory budget is low, since it'll result in a larger number of chunks, which will result in too many I/O operations that will slow down the whole solution.

4.2.1 Computation-to-Communication Ratio

We define the computation-to-communication ratio (γ) as:

$$\gamma = \frac{\alpha}{\beta}$$

Where:

- α = total CPU work (chunking, sorting, merging)
- β = communication overhead (MPI messages, file path exchange, coordination)

In our case, γ is high, because:

- Communication usually involves only file paths, which are lightweight, no big data is exchanged.
- Each worker operates independently and locally, no communication between workers.
- Communication only occurs in setup and final gather phase.

Conclusion: The high γ value indicates good scalability, meaning the system is performing significantly more useful work than communication overhead.