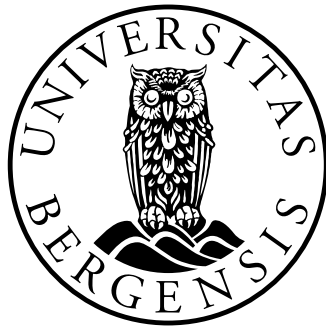


# Performance of Distributed and Shared Memory Parallel Sparse Matrix Vector Multiplication

Kristian Sordal



Thesis for Master of Science Degree at the University  
of Bergen, Norway

2025

©Copyright Kristian Sordal

The material in this publication is protected by copyright law.

Year: 2025

Title: Performance of Distributed and Shared Memory  
Parallel Sparse Matrix Vector Multiplication

Author: Kristian Sordal

# Acknowledgements

I want to thank mummi gætta  
and especially shimmy shimmy yeah shimmy yeah



# Abstract

SPARSE MATRIX VECTOR MULTIPLICATION is an important kernel used in scientific computing. It is a problem that lends itself well to parallelization. The problem is bounded by, and scales with the memory bandwidth of the system. Therefore in order to efficiently perform *SpMV* on large distributed memory systems, it is important to reduce the communication between nodes, in order to extract as much as possible out of the memory bandwidth of the system.

This thesis aims to investigate the results of *SpMV* when ran using different communication strategies.



# Contents

<b>Acknowledgements</b>	<b>iii</b>
<b>Abstract</b>	<b>v</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Sparse Matrix Vector Multiplication . . . . .	1
1.1.1 Computational Density . . . . .	1
<b>2 Background</b>	<b>3</b>
2.1 CSR Storage Format . . . . .	3
2.2 Load Balancing . . . . .	4
2.3 Distributed Memory CSR . . . . .	4
<b>3 Theory</b>	<b>5</b>
3.1 Definitions . . . . .	5
3.2 Amdahl's Law . . . . .	5
3.3 Partitioning Graph . . . . .	6
3.4 Separator . . . . .	6
<b>4 Communication Strategies</b>	<b>7</b>
4.1 1a - Exchange entire vector . . . . .	7
4.2 1b - Exchange only separators . . . . .	8
4.3 1c - Exchange only required separators . . . . .	9
4.4 1d - Exchange only required separator values . . . . .	10
<b>5 Previous Work</b>	<b>11</b>
<b>6 Results</b>	<b>13</b>





# Chapter 1

## Introduction

### 1.1 Sparse Matrix Vector Multiplication

Sparse Matrix Vector Multiplication (SpMV) is a common operation encountered in many areas of scientific computation. The matrices used in such operations are simultaneously large, but also sparse. *Sparse* in this context means that the average number of nonzeros per row is on the order of  $\mathcal{O}(1)$ , and thus it is worthwhile to treat nonzeros differently from non-valued entries. The performance of Sparse Matrix Vector Multiplication can be improved by utilizing parallel computing.

SpMV is known to be quite hard to optimize, both in sequential and parallel implementations. This is in part due to the fact that the operation has a low computational density.

#### 1.1.1 Computational Density

The *computational density* of an operation is defined by the relation between the number of floating point operations (FLOPS) and the number of memory accesses. We can then represent the computational density as in 1.1.

$$\text{Computational density} = \frac{\text{FLOPS}}{\text{Memory accesses}} \quad (1.1)$$

It becomes clear that operations with low computational density wont scale with the increased computing power



# Chapter 2

## Background

### 2.1 CSR Storage Format

CSR (Compressed Sparse Row) is the most widely used storage format for sparse matrices. As its name suggests, it compresses the amount of memory used to store a matrix without loss of information. It does so by utilizing three vectors  $A_p, A_j, A_x$ . Figure 2.1 shows an example of a matrix stored in CSR format, adapted from [1].

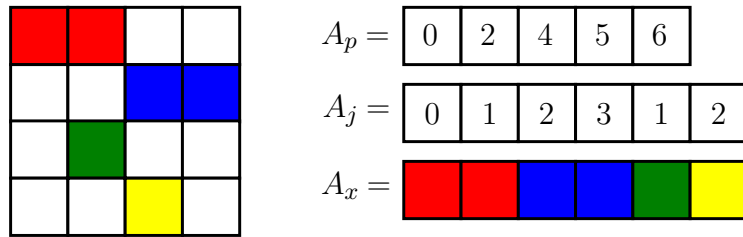


Figure 2.1: Example matrix represented in CSR Format.

The first vector,  $A_p$  stores the indices of the first nonzero in the vectors  $A_p$  and  $A_x$ . For a given entry  $A_p[i]$ ,  $A_p[i]$  is the index of the first nonzero in the  $i^{\text{th}}$  row.  $A_j[j]$  and  $A_x[j]$  denotes the column index and value of the  $j^{\text{th}}$  nonzero, respectively.

A sequential implementation of SpMV on a matrix stored in the CSR format can be implemented in the following manner:

---

**Algorithm 1:** Sequential CSR-based SpMV

---

**Input** :  $A_p, A_j, A_x, x$

**Output** :  $y$

```

for  $i \leftarrow 0$  to  $n$  do
   $y[i] \leftarrow 0$ 
  for  $j \leftarrow A_p[i]$  to  $A_p[i+1]$  do
     $y[i] \leftarrow y[i] + A_x[j] \cdot x[A_j[j]]$ 

```

---

This algorithm can be parallelized using shared memory parallelization using OpenMPs `#pragma omp parallel for` directive in the following manner:

---

**Algorithm 2:** Sequential CSR-based SpMV

---

**Input** :  $A_p, A_j, A_x, x$

**Output** :  $y$

```

#pragma omp parallel for
for  $i \leftarrow 0$  to  $n$  do
   $y[i] \leftarrow 0$ 
  for  $j \leftarrow A_p[i]$  to  $A_p[i+1]$  do
     $y[i] \leftarrow y[i] + A_x[j] \cdot x[A_j[j]]$ 

```

---

## 2.2 Load Balancing

## 2.3 Distributed Memory CSR

# Chapter 3

## Theory

### 3.1 Definitions

**Definition 3.1.1** (Separator). In the context of SpMV, a separator is a node in the graph that has an edge that strides between two partitions.

### 3.2 Amdahl's Law

Amdahl's Law provides a theoretical framework for understanding the limits of performance improvement when additional computational resources are applied to a given problem. It quantifies the potential speedup achieved by optimizing a specific portion of a system, emphasizing that the overall gain is constrained by the proportion of time the optimized component contributes to execution.

**Definition 3.2.1** (Amdahl's Law). The maximum achievable speedup of a computation is limited by the fraction of execution time that remains sequential, even when an arbitrarily large number of parallel resources is employed.

In the context of parallel computing, this principle highlights that while increasing the number of processing units can accelerate the parallelizable portion of a workload, the sequential fraction imposes a fundamental performance ceiling. Formally, if  $S$  denotes the total speedup,  $t_p$  represents the fraction of execution time that can be parallelized, and  $s_p$  is the speedup achieved for that parallelizable portion, Amdahl's Law is

expressed as:

$$S = \frac{1}{(1 - t_p) + \frac{t_p}{s_p}} \quad (3.1)$$

This equation reveals that as  $s_p \rightarrow \infty$ , the theoretical maximum speedup approaches  $\frac{1}{1 - t_p}$ , illustrating that the non-parallelizable portion becomes the dominant limiting factor in scalability.

### 3.3 Partitioning Graph

When performing Sparse Matrix Vector Multiplication in parallel, it is crucial to partition the matrix in such a way that the computation difference between all ranks is as small as possible. This is done by using a graph partitioner. In this project, the METIS graph partitioner was used. In particular, METIS' provides a function `METIS_partgraphkway`, given a parameter *nprocs* representing the number of processes the program is ran on, attempts to partition the graph in *nprocs* equal sized parts. It gives no guarantees that the partition will be optimal, because this problem is NP-Hard. It does however give an approximation which is good enough for all intents and purposes

### 3.4 Separator

When a graph is partitioned into different parts, there will inevitably be some edges which strides across different partitions. The endpoints of these edges are called separators, and will become important when it comes to reducing the communication load of the SpMV computation.

# Chapter 4

## Communication Strategies

This thesis evaluates and compares several communication strategies for Sparse Matrix-Vector Multiplication (SpMV) in a parallel, distributed memory setting. During each iteration of SpMV, every process computes a partial result of the output vector  $y$ .

In subsequent iterations, these computed values may be required by other processes to proceed with their own calculations. To ensure correctness, it is therefore necessary to communicate between processes so that each has access to the values it depends on. This section outlines a progression of increasingly efficient strategies for managing this communication.

### 4.1 1a - Exchange entire vector

The most straightforward approach is to have each rank send all of its computed values of  $y$  to every other rank. This ensures that all processes possess a complete and updated copy of the output vector before the next iteration. This strategy can be implemented using MPI's collective communication operation `MPI_Allgatherv`, which accommodates variable message sizes from each rank. Figure 4.1 illustrates the state of the  $y$  vector before and after communication using this strategy.

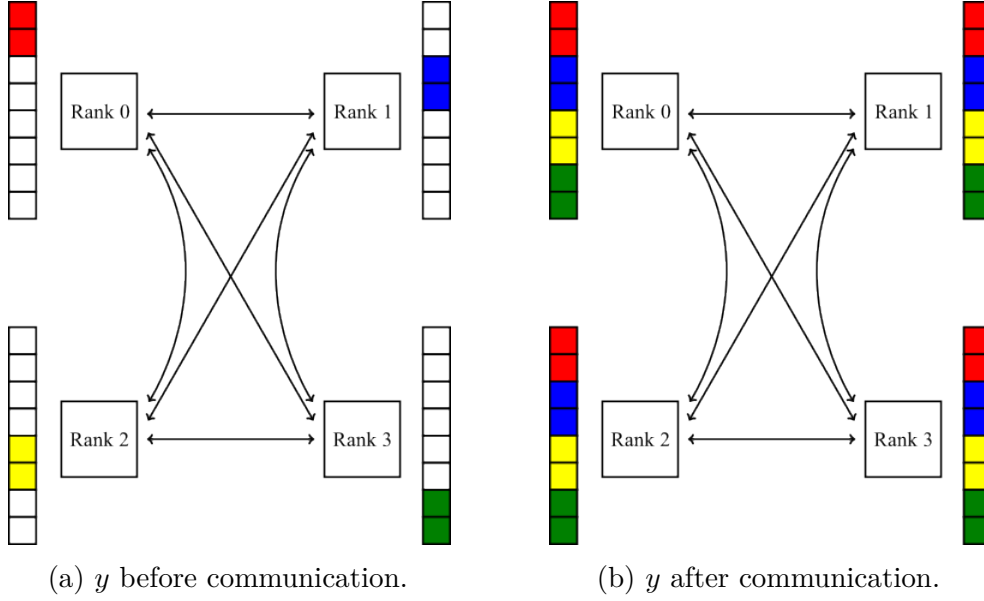


Figure 4.1: Visual representation of the  $y$  vector in communication strategy 1a.

---

**Algorithm 3:** 1a - Exchange entire vector

---

```

for each iteration do
    spmv(g,x,y)
    MPI_Allgatherv(local_y, sendcount, MPI_DOUBLE, y,
        recvcunts, displs, MPI_DOUBLE, MPI_COMM_WORLD)
    swap pointers of  $x$  and  $y$ 

```

---

## 4.2 1b - Exchange only separators

An improvement to the previous strategy can be achieved by recognizing that only separator values—those required by multiple processes—must be communicated. Non-separator values are used exclusively by the process that computed them and therefore do not need to be communicated.

To facilitate this strategy, separator values are reordered such that they appear at the beginning of each process's local segment of  $y$ . Once this structure is established, communication is performed using `MPI_Allgatherv`, transmitting only the subset of  $y$  that contains separator values. The number of separators on each process must be known beforehand, which can be computed by counting the number of elements that have neighbours belonging to a different partition.



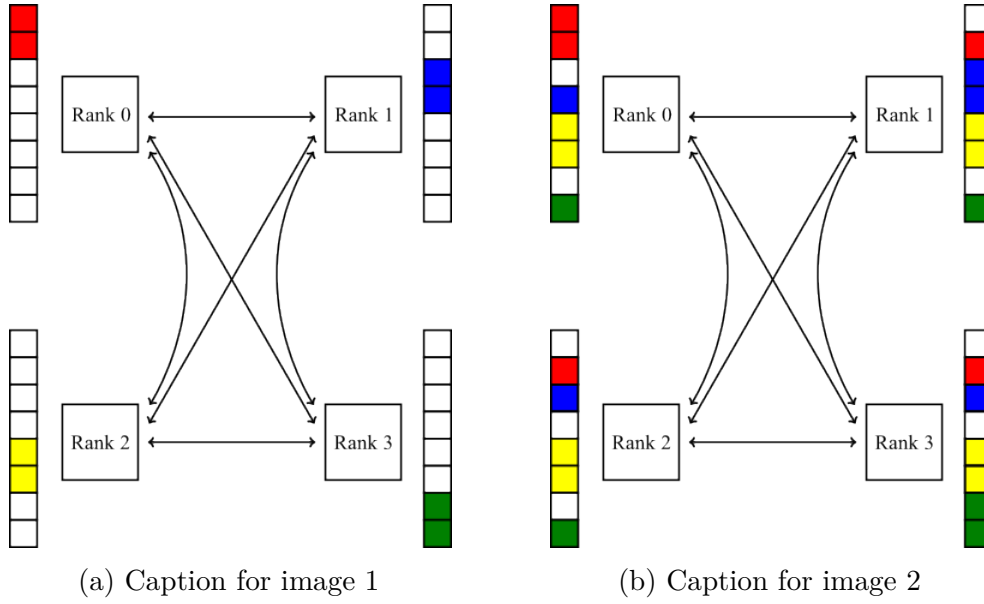


Figure 4.2: Main figure caption describing both subfigures

### 4.3 1c - Exchange only required separators

Further reduction to the communication volume can be achieved by observing that not all separator values are required by every process. As the number of partitions increases, the set of dependencies between partitions tends towards sparsity. As a consequence of this, certain sets of separators may only need to be communicated to a given subset of processes. Using this strategy, each process only communicates its set of separator values to the processes that require them.

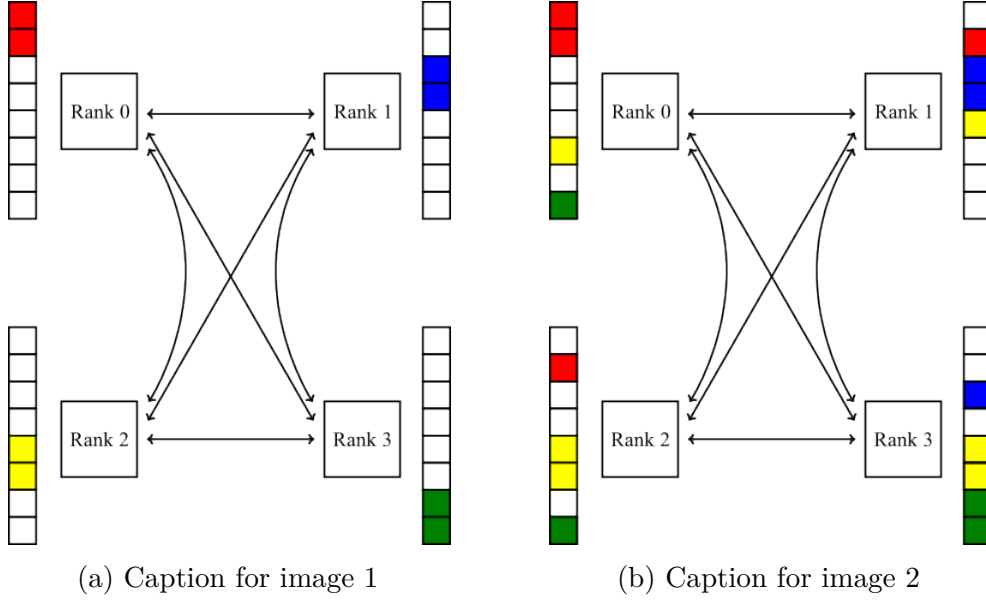


Figure 4.3: Main figure caption describing both subfigures

#### 4.4 1d - Exchange only required separator values

The final strategy aims to minimize communication overhead by transmitting only the exact subset of separator values that are both computed by and required for inter-process computation. If a specific separator value computed by one process is needed by exactly one other process, then only that single recipient receives the value.

This approach eliminates all unnecessary data transfers but introduces additional complexity in managing communication schedules. Dependencies must be mapped at a fine-grained level, and communication patterns must be explicitly tailored to the structure of the matrix and its partitioning.

## Chapter 5

### Previous Work



# Chapter 6

## Results

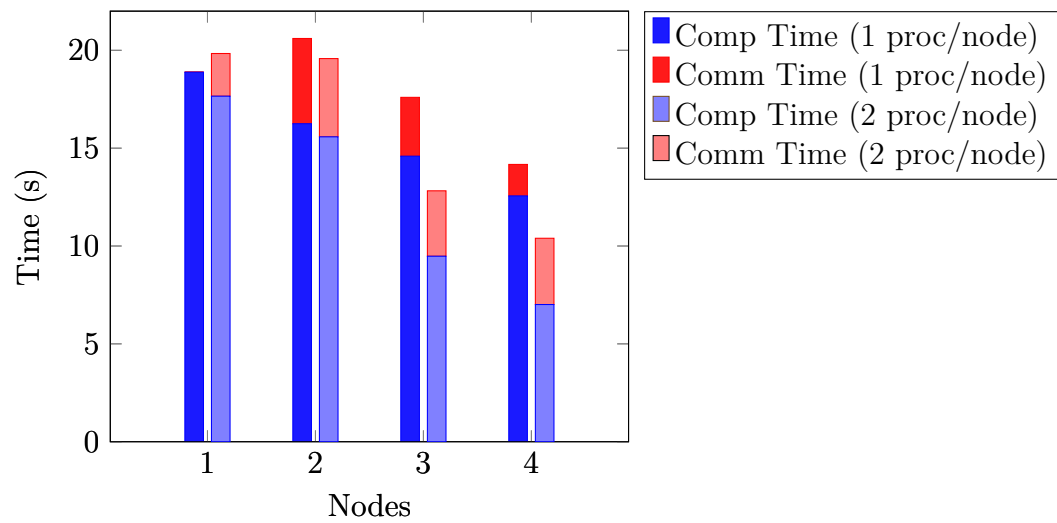


Figure 6.1: Stacked bar chart of SpMV and Halo Time.

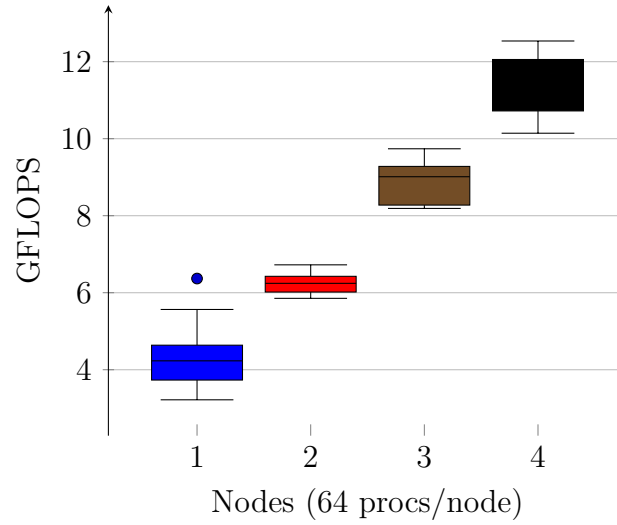


Figure 6.2: Aggregated Results of SpMV on defq chip

Table 6.1: Hardware used in our experiments. STREAM Triad was run with the `-march=native` and `-O3` compilation flags.

	Xeon-A	Xeon-B	Naples	Rome	Milan	TX2	Hi1620
CPUs	Intel Xeon Gold 6130	Intel Xeon Platinum 8168	AMD Epyc 7601	AMD Epyc 7302P	AMD Epyc 7763	Cavium TX2 CN9980	HiSilicon Kunpeng 920-6426
Instr. set	x86-64	x86-64	x86-64	x86-64	x86-64	ARMv8.1	ARMv8.2
Microarch.	Skylake	Skylake	Zen	Zen 2	Zen 3	Vulcan	TaiShan v110
Sockets	2	2	2	1	2	2	2
Cores	$2 \times 16$	$2 \times 24$	$2 \times 32$	$1 \times 16$	$2 \times 64$	$2 \times 32$	$2 \times 64$
Freq. [GHz]	1.9–3.6	2.5–3.7	2.7–3.2	1.5–3.3	2.5–3.5	2.0–2.5	2.6
L1I/core [KiB]	32	32	64	32	32	32	64
L1D/core [KiB]	32	32	32	32	32	32	64
L2/core [KiB]	1024	1024	512	512	512	256	512
L3/socket [MiB]	22	33	64	16	256	32	64
Mem. channels	$2 \times 6$	$2 \times 6$	$2 \times 8$	$1 \times 8$	$2 \times 8$	$2 \times 8$	$2 \times 8$
Bandwidth [GB/s]	256	256	342	204.8	409.6	342	342
Triad [GB/s]	147.1	137.4	169.7	90.9	256.5	236.4	260.4

# Bibliography

- [1] Gautam Gupta, Sivasankaran Rajamanickam, and Erik G. Boman. GAMGI: Communication-reducing algebraic multigrid for gpus. In *Proceedings of the 28th ACM SIGPLAN Annual Symposium on Principles and Practice of Parallel Programming (PPoPP)*, pages 61–75. ACM, 2024.