

---

# ADV-HPC

---

**Matteo Nunziante**  
Università degli studi di Trieste  
June 2024

---

# Contents

<b>I</b>	<b>Matrix Multiplication</b>	<b>3</b>
<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Methodology</b>	<b>3</b>
<b>3</b>	<b>Implementation</b>	<b>3</b>
3.1	Initialization . . . . .	3
3.2	Communication . . . . .	4
3.3	Computation . . . . .	5
<b>4</b>	<b>Scaling Results</b>	<b>5</b>
4.1	Matmul . . . . .	5
4.2	C-BLAS . . . . .	6
4.3	N = 10000 . . . . .	6
4.4	N = 80000 . . . . .	7
4.5	CU-BLAS . . . . .	9
<b>II</b>	<b>Jacobi</b>	<b>10</b>
<b>1</b>	<b>Introduction</b>	<b>10</b>
1.1	Domain Decomposition . . . . .	10
1.2	Initialization . . . . .	10
1.3	Evolution . . . . .	11
1.4	MPI Communications . . . . .	11
1.5	MPI I/O . . . . .	12
<b>2</b>	<b>Scaling Results</b>	<b>13</b>
2.1	CPU . . . . .	13
2.2	GPU . . . . .	14
2.3	OneSide Communcation . . . . .	15
<b>A</b>	<b>More Plots</b>	<b>16</b>

---

## Part I

# Matrix Multiplication

## 1 Introduction

This part is about the performance analysis of matrix multiplication in a distributed OpenMP + MPI + CUDA setting.

- Naive MATMUL
- C-BLAS
- CU-BLAS

The chosen function will perform the usual matrix multiplication  $C = A \cdot B$

## 2 Methodology

The goal is to implement multiple distributed matrix multiplication algorithms compatible with preexisting code. The main focus is on the performance of the algorithms in terms of execution time for the different sections as follows:

- Initialization
- Computation
- Communication

Analysis of performance is also based on the Speedup metric to measure the relative performance of the serial (1 process) and multiprocess implementations. For a given workload of  $k$  (for instance the dimension of the matrix) the speedup for  $size$  workers is defined as:

$$S_{size} = \frac{T_1}{T_{size}}$$

Note about timings: for each section execution time is the maximum of the times measured for each process.

## 3 Implementation

### 3.1 Initialization

Initialization is always performed on CPU using OMP threads, each process initializes a number of rows given by the formula:

$$N_{loc} = N/size + (rank < N \% size)$$

Where  $N$  is the number or rows of the matrix (and columns, assuming it is a square matrix),  $size$  is the number of mpi processes and  $rank$  is the id of the process.

This is repeated for matrix  $A$ ,  $B$  and  $C$ .

```

1 double* sequence_initialize(unsigned long long int N,unsigned long long int M
  ,int size, int rank, char* fname)
2 {
3     unsigned long long int N_loc=N, offset=0;
4
5     if (size > 1){
6         N_loc = N/size + 1*(rank < N%size);
7         offset = rank*N_loc + N%size*(rank >= N%size);
8     }
9
10    double *A = (double*)malloc(N_loc*M*sizeof(double));
11
12    #pragma omp parallel for
13    for (int i = 0; i < N_loc*M; i++)
14    {
15        A[i] = (offset*M+i);
16    }
17    return;
18 }

```

## 3.2 Communication

By construction, matrix matrix product iterates on rows of  $A$  and columns of  $B$ , this results in the need to gather the columns of  $B$  aggregating rows from every process. This is achieved with an All Gather.

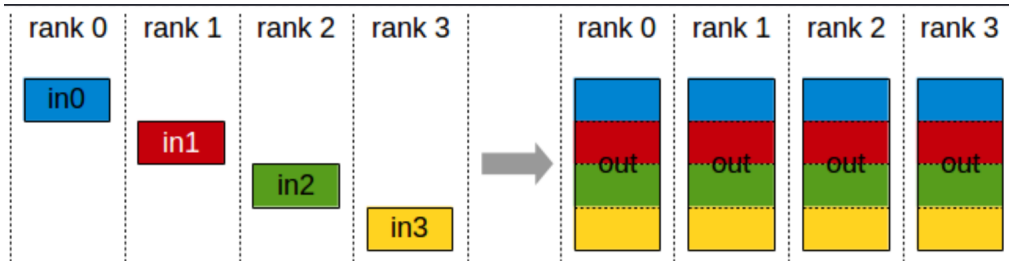


Figure 1: All Gather

Since it only works on contiguous blocks of memory, each process has to copy the selected columns in a new array, it can then proceed to call the `MPI_All_Gather_v` function (all gather v is necessary since the number of elements sent depends on the rows assigned to that process).

```

1 void cpu_copy_gather(double* restrict B,double* restrict B_loc, double*
  restrict B_mult, unsigned long long int N_loc, unsigned long long int N,
  unsigned long long int ncol, int *nrows, int *displs)
2 {
3     #pragma omp parallel for collapse(2)
4     for(unsigned long long int k = 0; k < N_loc; k++){
5         for(unsigned long long int j = 0; j < ncol; j++){
6             B_loc[k*ncol+j] = B[ N*k + j ];
7         }
8     }
9 }

```

---

```

10 MPI_Allgatherv(B_loc, N_loc*ncol, MPI_DOUBLE, B_mult, nrows, displs,
11 MPI_DOUBLE, MPI_COMM_WORLD);
12 }

```

Matrix allocation is performed only once at the beginning of the program, since tests showed it is a quite consuming operation.

### 3.3 Computation

The actual product is performed block-wise among rows of matrix  $A$  stored in the process memory and columns of matrix  $B$  gathered as above. This is repeated  $size$  times as columns of  $B$  are gathered in chunks of number of columns  $ncol = N/size + (i < N\%size)$  where  $i$  is an index ranging from 0 to  $size - 1$

---

#### Algorithm 2 Matrix Multiplication

---

```

1: for  $i = 1, 2, \dots, size$  do
2:   B_mult = Gather(B,i)
3:   Copy B_mult to Device
4:   mult(A,B_mult,C,i)
5: end for

```

---

Figure 2: Matrix Multiplication

## 4 Scaling Results

### 4.1 Matmul

Scaling is performed mapping processes to nodes (2x56 cores Intel Sapphire Rapids at 2.00 GHz) using all threads available (112 - no hyperthreading).

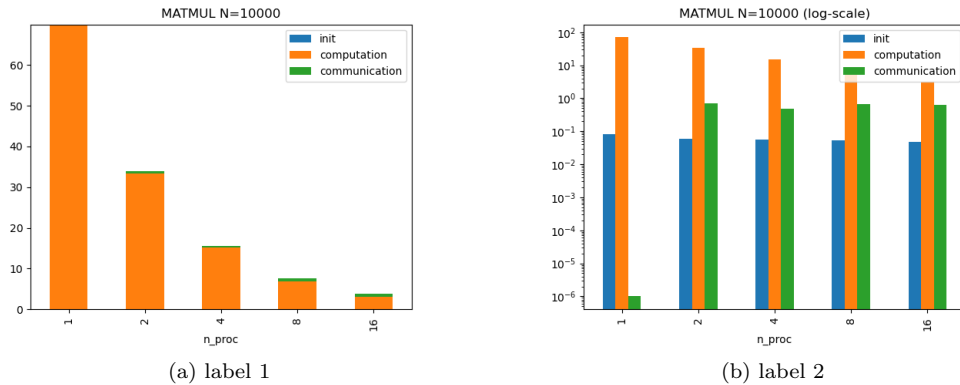


Figure 3: Matmul Scaling

- **Initialization** doesn't scale well with the number of processes, this is due to the fact that the size of the matrix doesn't justify the overhead of threads creation.

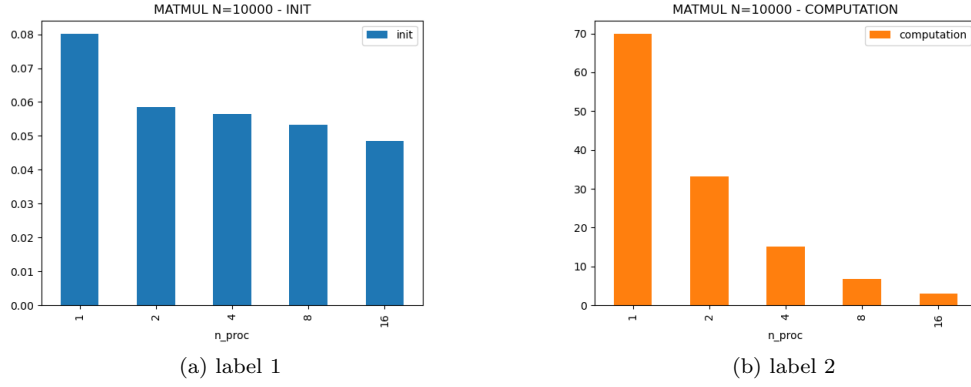


Figure 4: Matmul init and comp

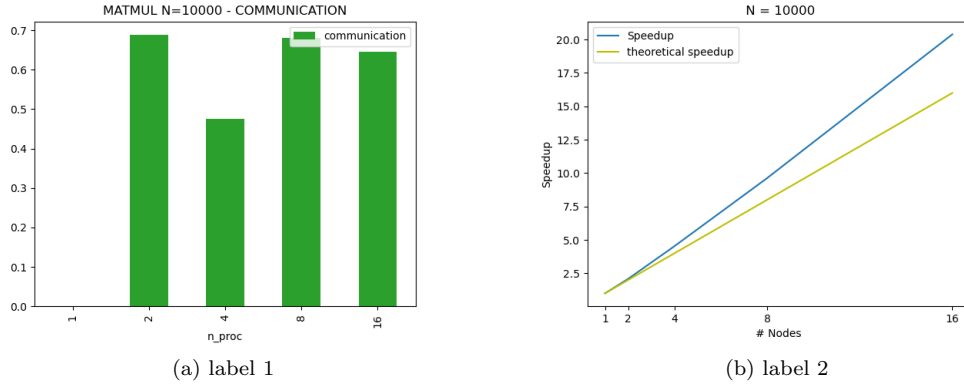


Figure 5: Matmul comm and speedup

## 4.2 C-BLAS

Scaling is performed mapping processes to nodes (2x56 cores Intel Sapphire Rapids at 2.00 GHz) using all threads available (112 - no hyperthreading).

## 4.3 N = 10000

- **Communcation** is quite unstable, it is probably due to Leonardo's network congestion. Working on 16 nodes results in the lowest communication time, this is again an unexpected phenomena.

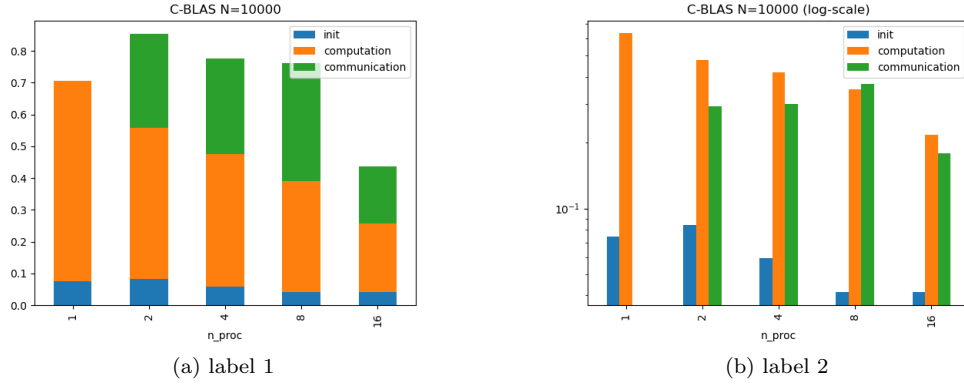


Figure 6: C-BLAS Scaling

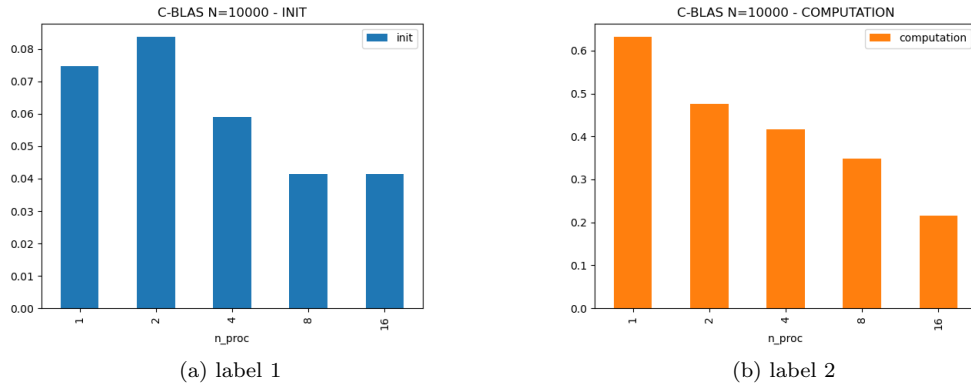


Figure 7: C-BLAS init and comp

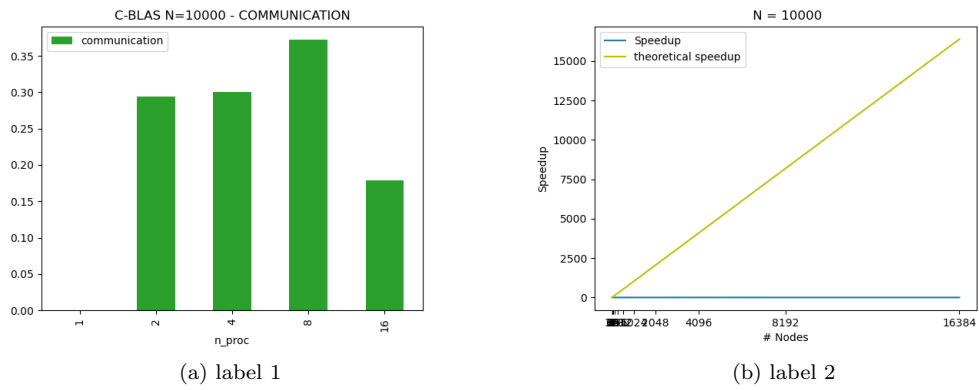


Figure 8: C-BLAS comm and speedup

#### 4.4 $N = 80000$

- **Initialization** at 2 cpu nodes on 80000 size results in unexpected low timings. this is not a coincidence, the initialization has been performed multiple times and the same

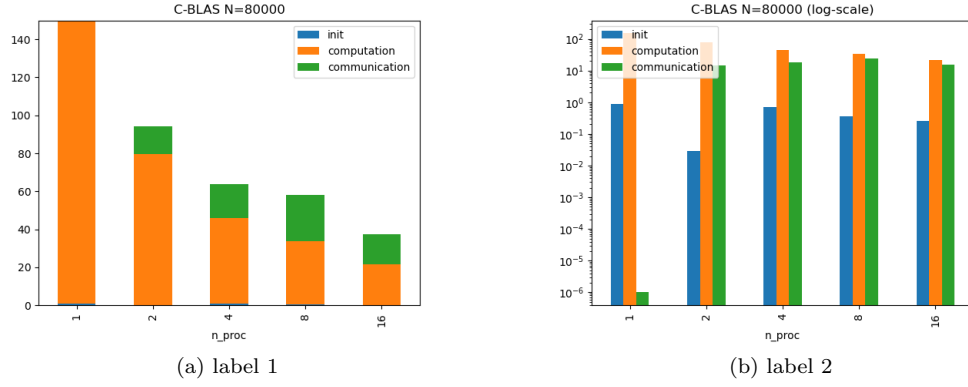


Figure 9: C-BLAS Scaling

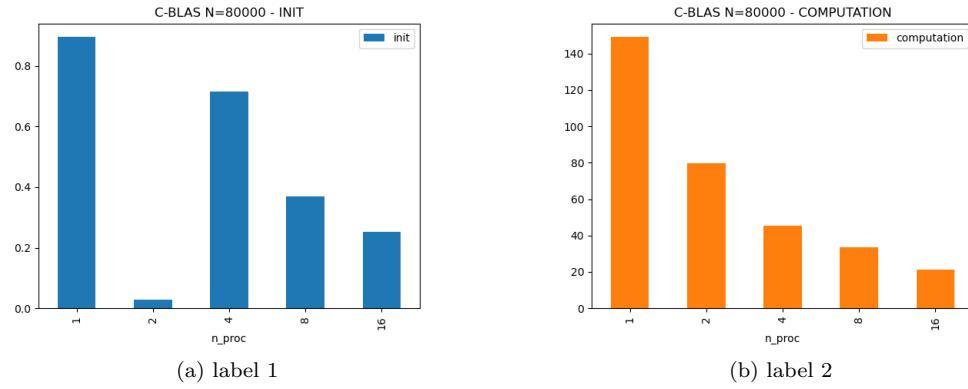


Figure 10: C-BLAS init and comp

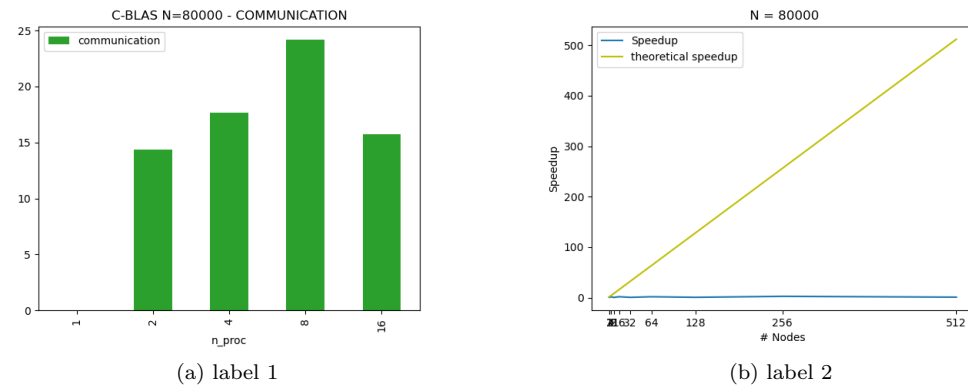


Figure 11: C-BLAS comm and speedup

phenomena has been observed. The size of the matrix (equivalent to a 55568 size on a single node) probably resonates with the memory architecture of the nodes even though analysis of cache misses (using perf) shows no significant difference between the two cases (40000 vs 80000 size).



## 4.5 CU-BLAS

Scaling is performed mapping processes to GPUs (4 x NVIDIA Ampere A100 GPUs, 64GB per node) assigning 8 cores per process (32 cores Ice Lake at 2.60 GHz per node)

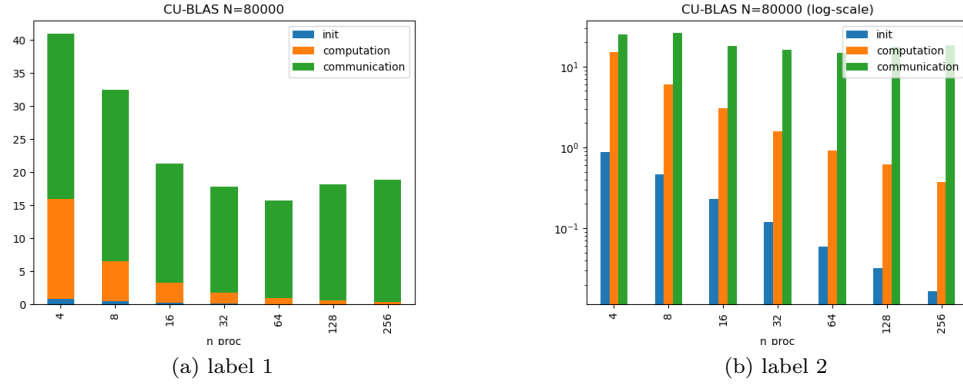


Figure 12: CU-BLAS Scaling

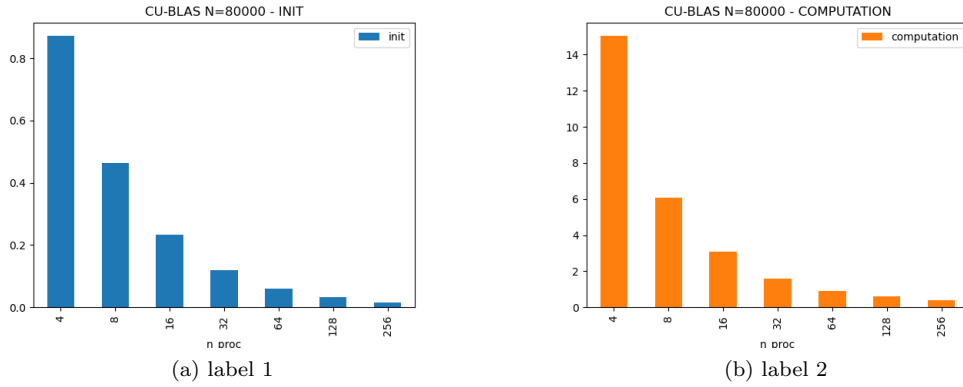


Figure 13: CU-BLAS init and computation

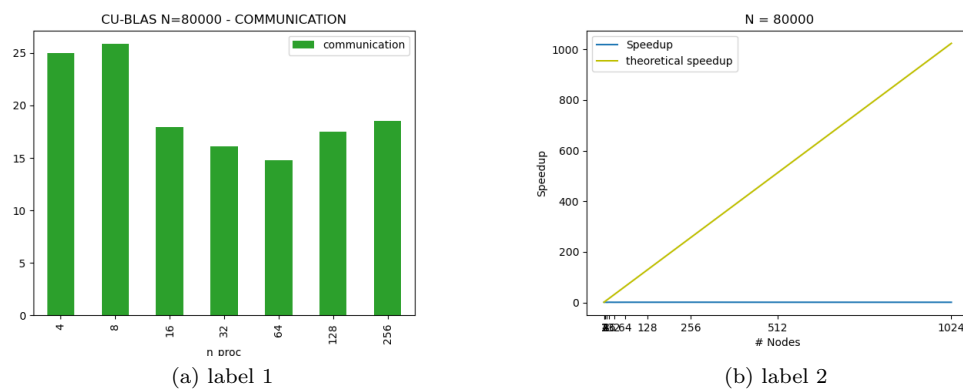


Figure 14: CU-BLAS comm and speedup

## Part II

# Jacobi

## 1 Introduction

### 1.1 Domain Decomposition

Domain is partitioned as in Figure 15. Halo zones are needed since each cell needs to be updated using information from up, down, left and right neighbour.

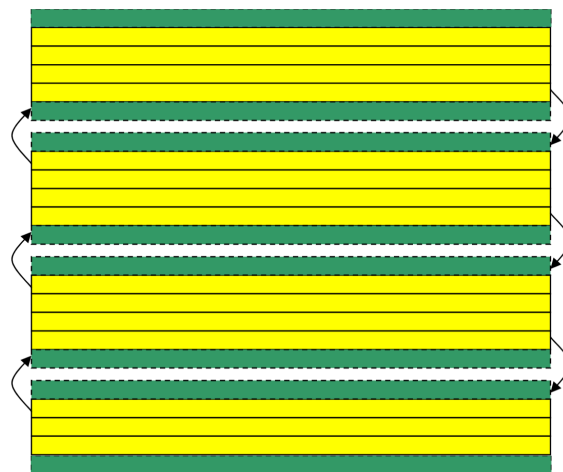


Figure 15: Domain Decomposition

### 1.2 Initialization

Initialization of the Matrix is performed using OMP / OACC when GPUs are available [16].

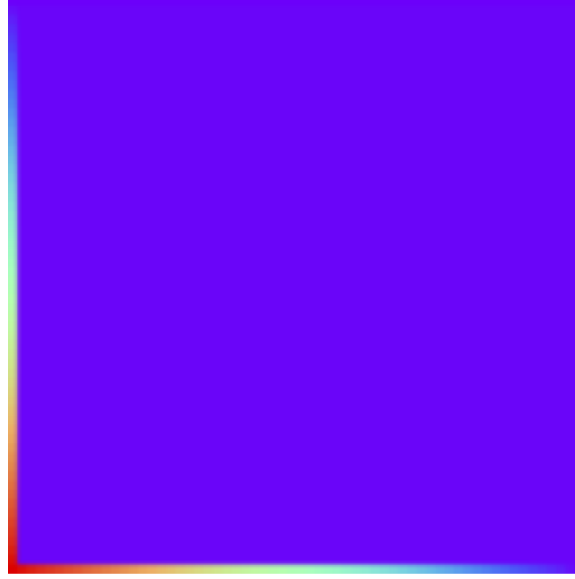


Figure 16: Initial Configuration - Boundary Conditions - 60x60 matrix

### 1.3 Evolution

Evolution of the matrix is performed again using OMP / OACC when GPUs are available.

```

1 void evolve( double* restrict matrix, double* restrict matrix_new, const int
    N_loc, const size_t dimension )
2 {
3     #pragma acc parallel loop collapse(2) present( matrix, matrix_new )
4     for(size_t i = 1 ; i < N_loc-1; ++i ){
5         for(size_t j = 1; j < dimension+1; ++j ){
6             matrix_new[ ( i * ( dimension + 2 ) ) + j ] = ( 0.25 ) *
7                 (matrix[ ( ( i - 1 ) * ( dimension + 2 ) ) + j ] +
8                  matrix[ ( i * ( dimension + 2 ) ) + ( j + 1 ) ] +
9                  matrix[ ( ( i + 1 ) * ( dimension + 2 ) ) + j ] +
10                 matrix[ ( i * ( dimension + 2 ) ) + ( j - 1 ) ] );
11         }
12     }
13 }

```

### 1.4 MPI Communications

Communication among processes allows to update the "shared" rows after the update of the "private"<sup>1</sup> ones. Each process is set to send first and last internal (meaning row 1 and *my\_row*) rows to previous and next process defined by the formulas

$$\text{next} = \text{rank} == 0 ? \text{MPI\_PROC\_NULL} : \text{rank} - 1$$

$$\text{prev} = \text{rank} == \text{size} - 1 ? \text{MPI\_PROC\_NULL} : \text{rank} + 1$$

When working with GPUs communication is **Cuda-Aware**, avoiding the extra copy

<sup>1</sup>shared and private do not refer to any memory policy but to the rows of the matrix that show up in multiple or just one process

---

to Host.

```
1 void update_boundaries(int rank, int size, double* local_matrix, int my_row,
2 int N)
3 {
4     MPI_Request send_request;
5     MPI_Status recv_status;
6
7     int tag1=0,tag2=1;
8
9     int prev = rank == 0? MPI_PROC_NULL : rank-1;
10    int next = rank == size-1? MPI_PROC_NULL : rank+1;
11
12    // sending last row to rank +1
13    #pragma acc host_data use_device( local_matrix )
14    {
15        MPI_Isend(&local_matrix[(my_row-2)*N],N, MPI_DOUBLE, next, tag1,
16        MPI_COMM_WORLD, &send_request);
17        // sending first row to rank -1
18        MPI_Isend(&local_matrix[N],N, MPI_DOUBLE, prev, tag2, MPI_COMM_WORLD,
19        &send_request);
20
21        // reveiving first row from rank -1
22        MPI_Recv(local_matrix,N ,MPI_DOUBLE,prev,tag1,MPI_COMM_WORLD,&
23        recv_status);
24        // receiving last row from rank +1
25        MPI_Recv(&local_matrix[(my_row-1)*N],N ,MPI_DOUBLE,next,tag2,
26        MPI_COMM_WORLD,&recv_status);
27    }
28 }
```

The communication involves two **Non-Blocking Send** and two **Blocking Receive** for process, This choice avoids any possible deadlock and assures the process doesn't start a new iteration before getting updated shared rows.

## 1.5 MPI I/O

Parallel write to binary file is also provided. As before, if working on GPUs Cuda-Aware write is performed.

```
1 void mpi_io(int rank, int size, double* local_matrix, int N_loc, int LDA,
2 char* fname_snap, MPI_File file)
3 {
4     char filename[256];
5     sprintf(filename, "%s.bin", fname_snap);
6
7     MPI_Offset offset = ((N_loc-2) * rank + (LDA-2)%size*(rank>=(LDA-2)%size)
8     ) * LDA * sizeof(double);
9
10    MPI_File_open(MPI_COMM_WORLD, filename, MPI_MODE_CREATE | MPI_MODE_WRONLY
11    , MPI_INFO_NULL, &file);
12
13    // Set the file pointer to the correct position
14    MPI_File_seek(file, offset, MPI_SEEK_SET);
15
16    // Set the file view
```

---

```
14 MPI_File_set_view(file, offset, MPI_DOUBLE, MPI_DOUBLE, "native",
15 MPI_INFO_NULL);
16
17 // Write the local matrix block to the file
18 #pragma acc host_data use_device( local_matrix )
19 MPI_File_write(file, &local_matrix[LDA], (N_loc-2)*LDA, MPI_DOUBLE,
20 MPI_STATUS_IGNORE);
21
22 // Close the file
23 MPI_File_close(&file);
24 }
```

The final image [17] can be showed using *Gnuplot*.

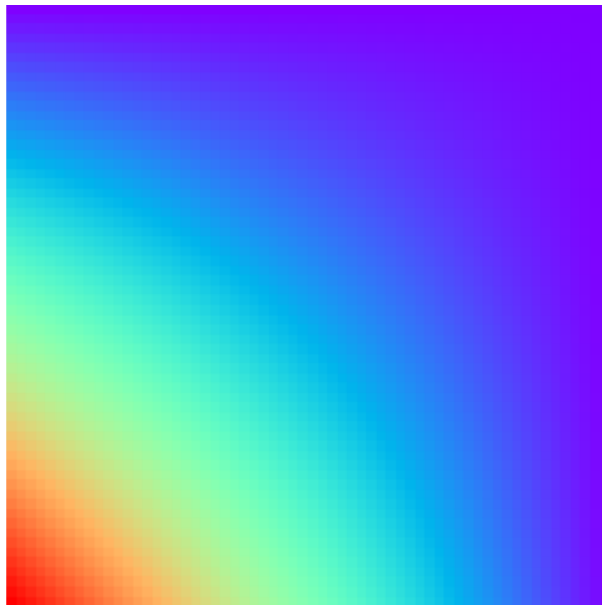


Figure 17: Final Configuration - 2000 iterations - 60x60 matrix

## 2 Scaling Results

### 2.1 CPU

Scaling is performed mapping processes to nodes (2x56 cores Intel Sapphire Rapids at 2.00 GHz) using all threads available (112 - no hyperthreading).

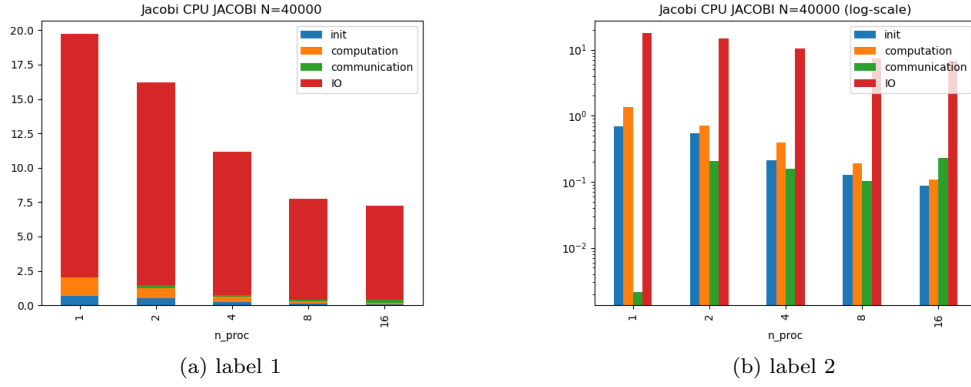


Figure 18: CPU Scaling

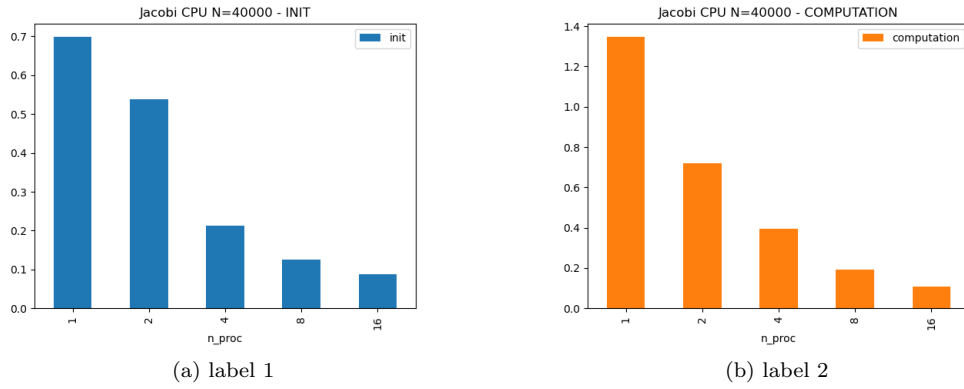


Figure 19: CPU Scaling

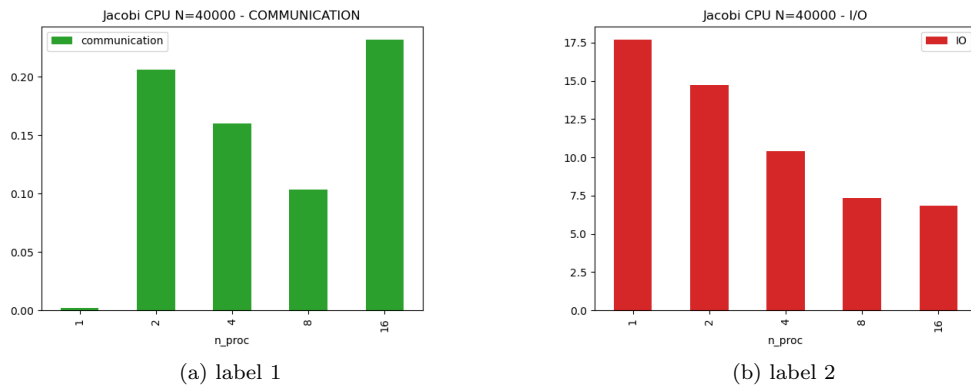


Figure 20: CPU Scaling

## 2.2 GPU

Scaling is performed mapping processes to GPUs (4 x NVIDIA Ampere A100 GPUs, 64GB per node) assigning 8 cores per process (32 cores Ice Lake at 2.60 GHz per node)

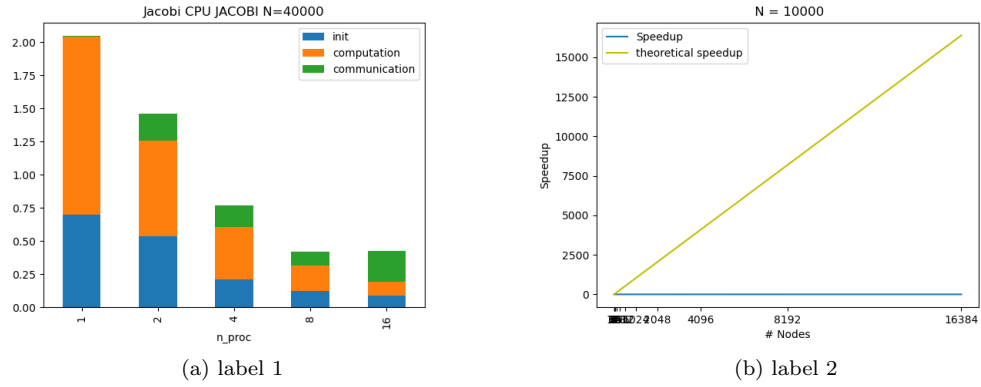


Figure 21: CPU Scaling

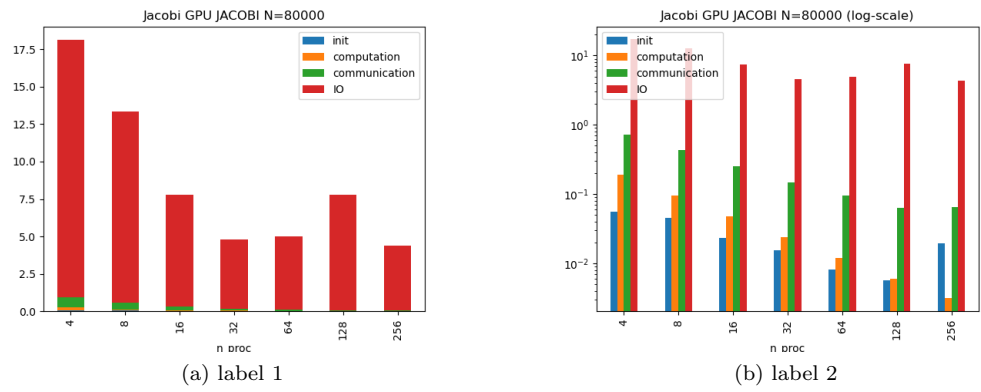


Figure 22: GPU Scaling

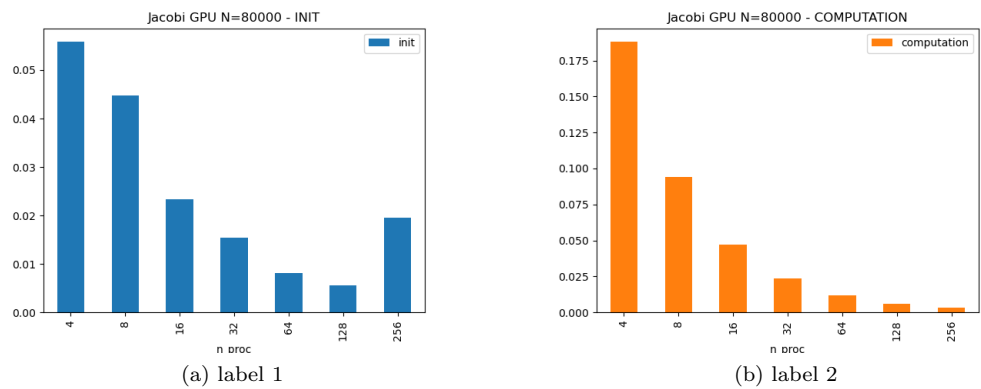


Figure 23: GPU Scaling

## 2.3 OneSide Communcation

```
1 MPI_Win win_u_n, win_d_n;
```

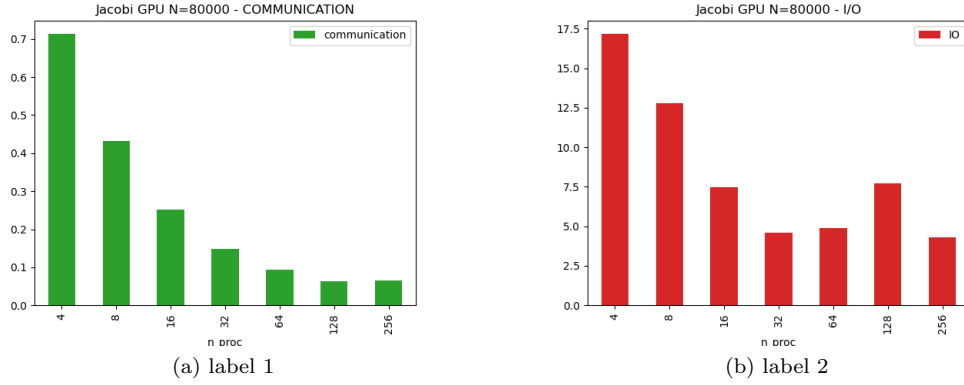


Figure 24: GPU Scaling

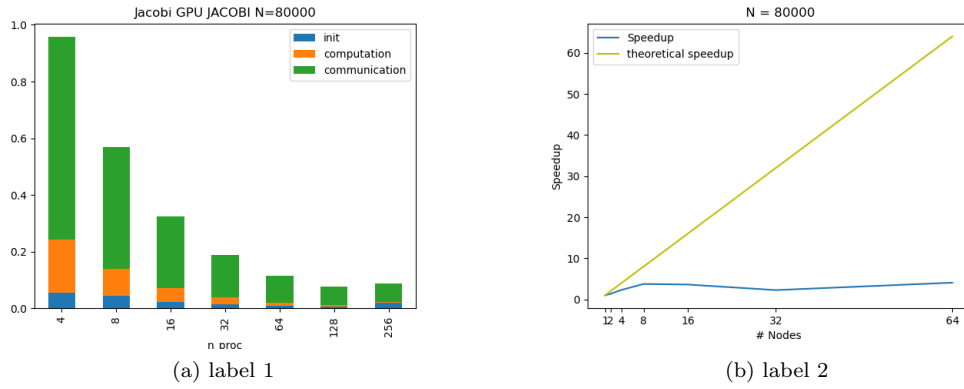


Figure 25: CPU Scaling

```

2 MPI_Win_create(&matrix_new[dimension+2], (MPI_Aint)(dimension+2) * sizeof(
    double), sizeof(double), MPI_INFO_NULL, MPI_COMM_WORLD, &win_u_n);
3 MPI_Win_create(&matrix_new[(N_loc-2)*(dimension+2)], (MPI_Aint)(dimension+2)
    * sizeof(double), sizeof(double), MPI_INFO_NULL, MPI_COMM_WORLD, &win_d_n)
    ;
4
5 MPI_Win win_u, win_d;
6 MPI_Win_create(&matrix[dimension+2], (MPI_Aint)(dimension+2) * sizeof(double)
    , sizeof(double), MPI_INFO_NULL, MPI_COMM_WORLD, &win_u);
7 MPI_Win_create(&matrix[(N_loc-2)*(dimension+2)], (MPI_Aint)(dimension+2) *
    sizeof(double), sizeof(double), MPI_INFO_NULL, MPI_COMM_WORLD, &win_d);

```

- Communication is less invasive than but we have to account also for window creation

## Appendix A More Plots



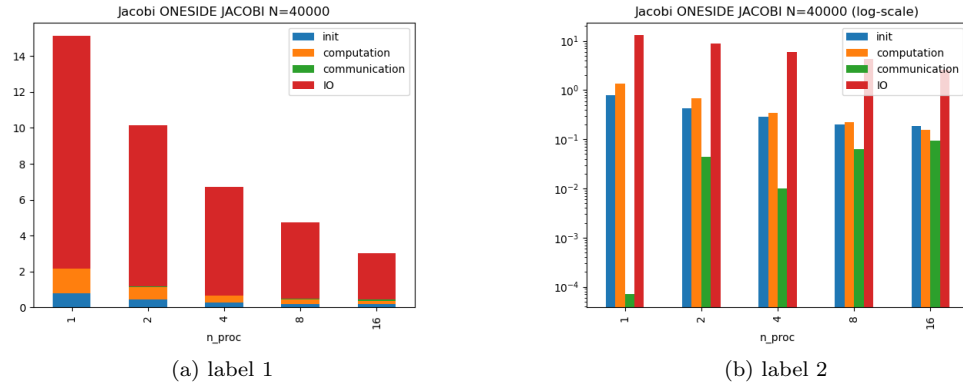


Figure 26: ONESIDE Scaling

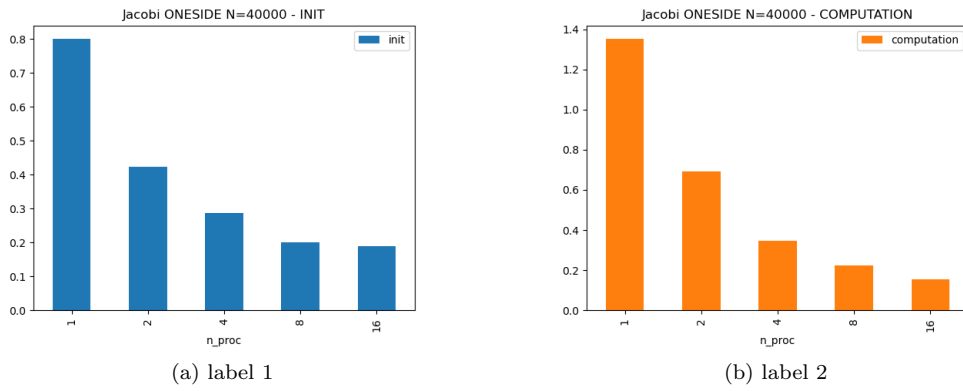


Figure 27: ONESIDE Scaling

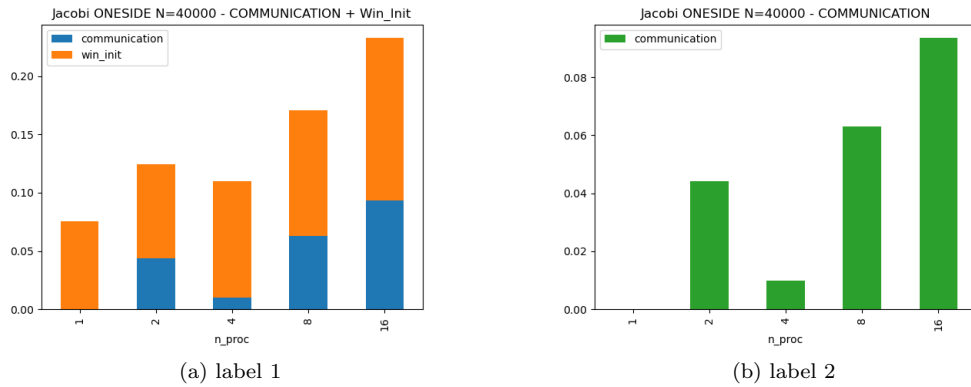


Figure 28: ONESIDE Scaling

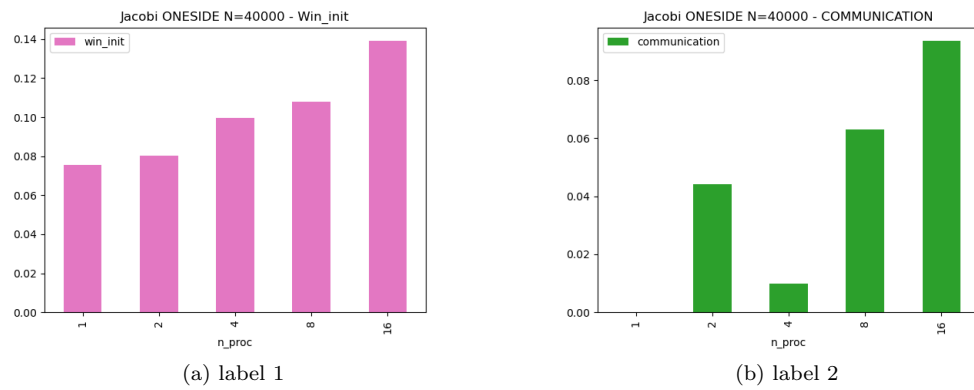


Figure 29: ONESIDE Scaling

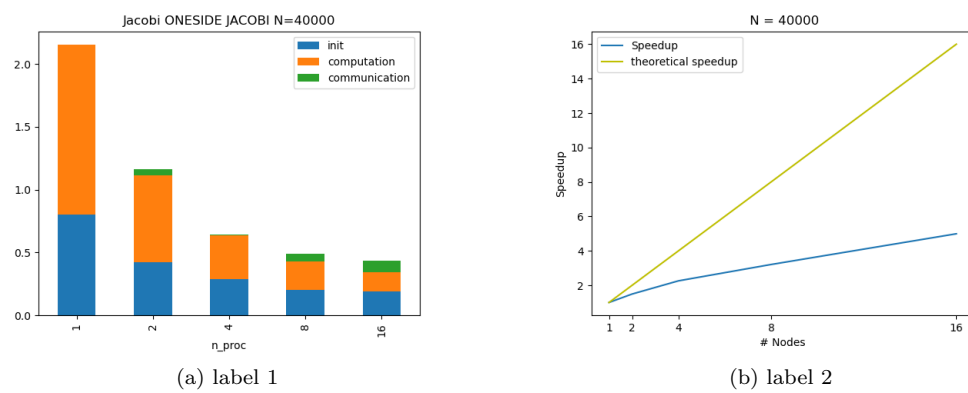


Figure 30: CPU Scaling