

HPC Assignment Presentation

Simone Pio Maurutto
Sara Maddalena Piccinini
Francesca Scognamiglio

Politecnico di Torino

January 2026

- Systolic array structures for matrix multiplication with MPI

- Fundamentals of multi-dimensional data processing with CUDA

- Heat Diffusion Simulation in 2D Grid with OpenMP

General Overview

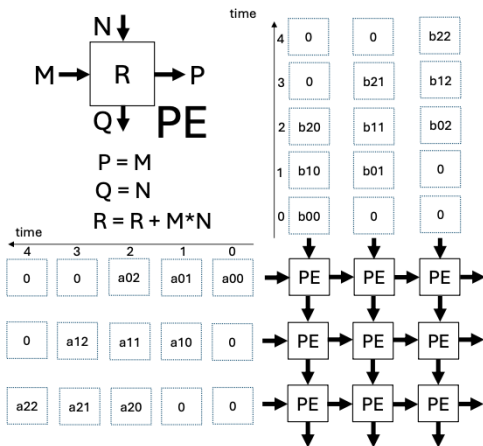
- 1 Introduction
- 2 Methodology
- 3 Metrics
- 4 Results and Analysis
- 5 Conclusions

- Design and implementation of a systolic array architecture for matrix multiplication
- Distributed realization using MPI processing elements
- Performance evaluation and scalability analysis
- Discussion of results and trade-offs

- ① Partition of input matrices among processing elements (PEs)
- ② Cannon's algorithm initialization for block alignment
- ③ Cannon's algorithm computation phase
- ④ Gathering of local results into the output matrix

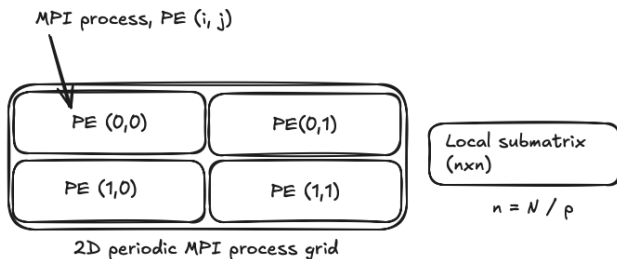
Methodology - Processing elements and Systolic Behavior

Processing Elements (PEs) are the fundamental units of the systolic array architecture. They perform local computations and propagate data across the array over time.



Methodology – Process Mapping

- Each Processing Element (PE) is mapped to one MPI process
- Block size: $n = N/p$ with N input matrices and p number of rocess per dimension
- Processes arranged in a 2D periodic grid
- Inter-process communication implemented via MPI



Methodology – Data Distribution and Initialization

- Input matrices are distributed among processing elements
- The PEs are initialized according to Cannon's algorithm requirements
- Data distribution and gathering implemented using MPI Scatter and Gather operations after the main computation

Data distribution among processing elements

```
MPI_Scatterv(A, counts, displs, blocktype2, A_block, block_size*block_size, MPI_DOUBLE, 0, MPI_COMM_WORLD);  
MPI_Scatterv(B, counts, displs, blocktype2, B_block, block_size*block_size, MPI_DOUBLE, 0, MPI_COMM_WORLD);
```

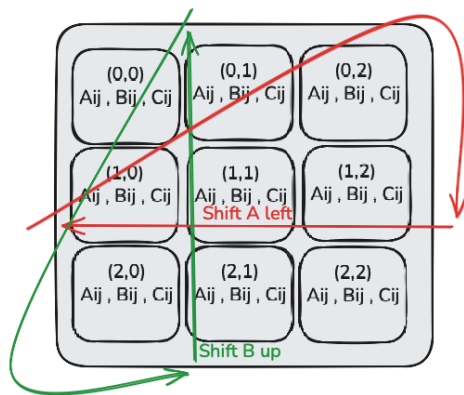
Result gathering after the computation

```
MPI_Gatherv(C_block, block_size*block_size, MPI_DOUBLE, C, counts, displs, blocktype2, 0, MPI_COMM_WORLD);
```


Methodology – Computation Phase

During each iteration, the processing elements:

- 1 Compute the local multiply
- 2 Exchange data blocks with neighboring processes (A left, B up on a 2D periodic grid)



2D periodic process grid (wrap-around)

- Experiments were executed on one or two nodes of the **Legion HPC cluster**
- The number of MPI processes was varied according to the matrix size and the 2D process grid configuration
- Execution time was measured on the root process using **MPI_Wtime()**

Results and Analysis

- Matrix-matrix multiplication
- Process grid and data movement
- Performance considerations

Conclusions

- Matrix-matrix multiplication
- Process grid and data movement
- Performance considerations

CUDA Filtering

- Image filtering with CUDA kernels
- Memory access patterns
- Speedup over CPU implementation

OpenMP Heat Diffusion

- Parallelization strategy
- Strong and weak scaling

Conclusions

- Summary of results
- Lessons learned
- Possible improvements