

# Enabling lighter and faster simulations with repeated matrix blocks

**J. Plana-Riu**<sup>1</sup>, F.X. Trias<sup>1</sup>, G. Colomer<sup>1</sup>, À. Alsaltí-Baldellou<sup>2</sup>, X. Álvarez-Farré<sup>3</sup>, A. Oliva<sup>1</sup>

<sup>1</sup>Heat and Mass Transfer Technological Centre  
Technical University of Catalonia

<sup>2</sup>Department of Civil, Environmental and Architectural Engineering (ICEA)  
University of Padova

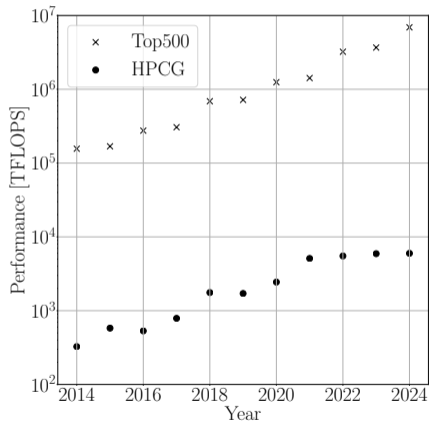
<sup>3</sup> High-Performance Computing and Visualization Team  
SURF

35th Parallel Computational Fluid Dynamics (ParCFD) International Conference  
Bonn, Germany  
Sept. 3rd, 2024



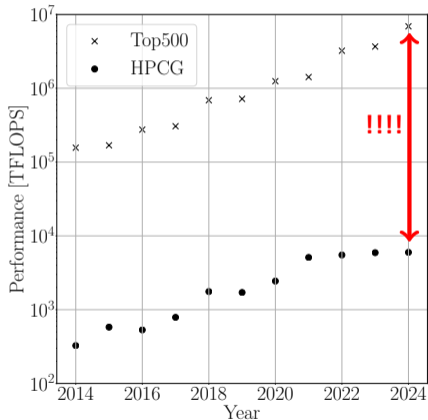
# Introduction

One big problem...



# Introduction

One big problem...

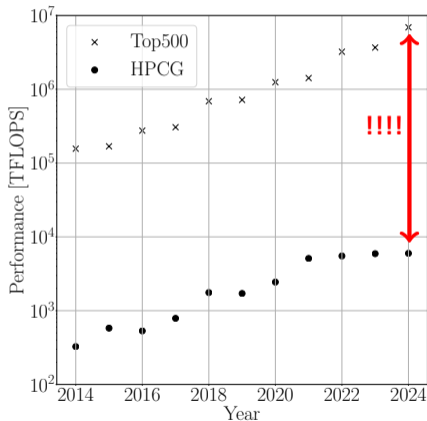


## Existing resources...

- Computational power from current top HPC systems is in the exaflop range...
- Sparse algebra, however...
  - has a low arithmetic intensity
  - is limited by memory bandwidth
- HPCG is the benchmark for us.

# Introduction

One big problem...



## Existing resources...

- Computational power from current top HPC systems is in the exaflop range...
- Sparse algebra, however...
  - has a low arithmetic intensity
  - is limited by memory bandwidth
- HPCG is the benchmark for us.

## Possible solutions...

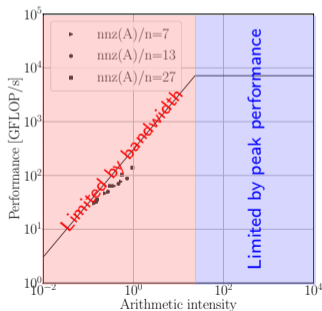
- Improving arithmetic intensity!

# Introduction

## Arithmetic intensity

### Roofline model<sup>1</sup>

- Memory-bound
- Compute-bound



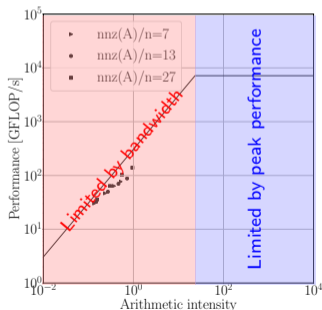
<sup>1</sup>S. Williams et al. "Roofline: an insightful visual performance for multicore architectures," *Commun. ACM* **52**, 2009

# Introduction

## Arithmetic intensity

### Roofline model<sup>1</sup>

- Memory-bound
- Compute-bound



### What is the arithmetic (or operational) intensity?

- Ratio between the number of operations and the amount of data that has to be handled (sent/received)

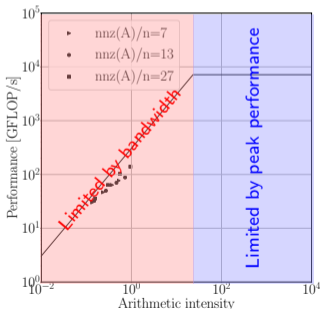
<sup>1</sup>S. Williams et al. "Roofline: an insightful visual performance for multicore architectures," *Commun. ACM* **52**, 2009

# Introduction

## Arithmetic intensity

### Roofline model<sup>1</sup>

- Memory-bound
- Compute-bound



### What is the arithmetic (or operational) intensity?

- Ratio between the number of operations and the amount of data that has to be handled (sent/received)

$$\begin{pmatrix} A & \\ & A \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}$$

2 SpMV

$$A(u_1 \ u_2)$$

1 SpMM with 2 RHS

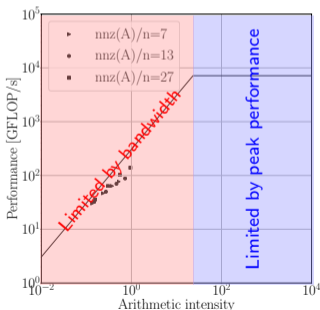
<sup>1</sup>S. Williams et al. "Roofline: an insightful visual performance for multicore architectures," *Commun. ACM* **52**, 2009

# Introduction

## Arithmetic intensity

### Roofline model<sup>1</sup>

- Memory-bound
- Compute-bound



### What is the arithmetic (or operational) intensity?

- Ratio between the number of operations and the amount of data that has to be handled (sent/received)

Oper.	# Mat.	# Vec. sent	# Vec. rcv	# Ops
2x SpMV	2	2	2	2
1x 2-SpMM	1	2	2	2
Oper.	Equivalent arithmetic intensity**			
2x SpMV	$2/(2+2+2)=1/3$			
1x 2-SpMM	$2/(1+2+2)=2/5$			

<sup>1</sup>S. Williams et al. "Roofline: an insightful visual performance for multicore architectures," *Commun. ACM* **52**, 2009



# Repeated matrix blocks

- CFD simulations are full of sparse matrix-vector products (SpMV):
  - $u^{n+1} = u^{*,n+1} - G\psi^{n+1}$
  - $u_i^* = u^n + \Delta t \sum_{j=1}^{i-1} a_{ij}(Du_j - C(u_j)u_j)$
  - $r_{k+1} = r_k - \alpha_k Ap_k$
  - ...

## Following the previous example...

- If some repeated matrix block structures are present SpMV can be translated to SpMM:
  - Symmetries
  - Repeated geometry patterns
  - Ensemble averaging parallel-in-time

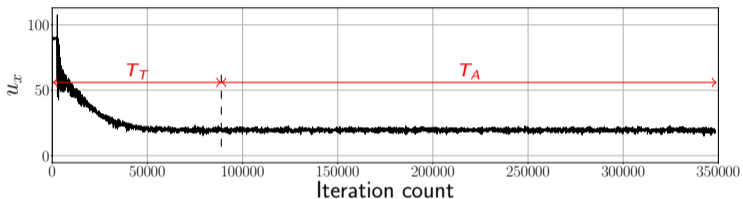
# Repeated matrix blocks

- CFD simulations are full of sparse matrix-vector products (SpMV):
  - $u^{n+1} = u^{*,n+1} - G\psi^{n+1}$
  - $u_i^* = u^n + \Delta t \sum_{j=1}^{i-1} a_{ij}(Du_j - C(u_j)u_j)$
  - $r_{k+1} = r_k - \alpha_k Ap_k$
  - ...

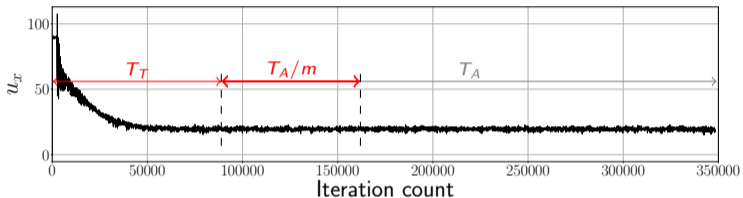
## Following the previous example...

- If some repeated matrix block structures are present SpMV can be translated to SpMM:
  - Symmetries
  - Repeated geometry patterns
  - **Ensemble averaging parallel-in-time**

# Ensemble averaging parallel-in-time



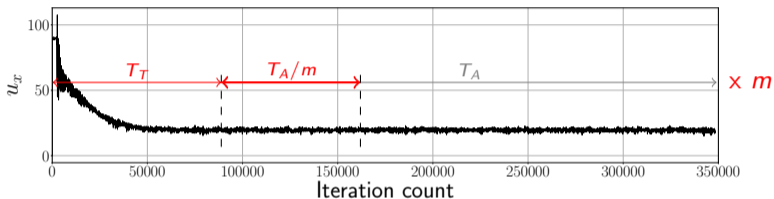
# Ensemble averaging parallel-in-time



## Ensemble average

$$U_x = \frac{1}{m} \sum_{i=1}^m \langle u_{x,i} \rangle = \frac{1}{m} \sum_{i=1}^m \frac{1}{T - T_T} \int_{T_T}^T u_{x,i} dt$$

# Ensemble averaging parallel-in-time



## Ensemble average

$$U_x = \frac{1}{m} \sum_{i=1}^m \langle u_{x,i} \rangle = \frac{1}{m} \sum_{i=1}^m \frac{1}{T - T_T} \int_{T_T}^T u_{x,i} dt$$

## $m$ simulations as a single one...

$$\mathbb{C} = I_m \otimes C$$

$$\mathbb{D} = I_m \otimes D$$

$$\mathbb{G} = I_m \otimes G$$

# Ensemble averaging parallel-in-time

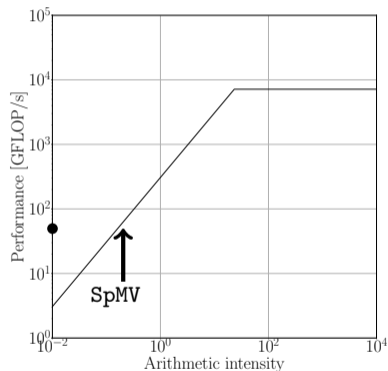
$$\frac{du}{dt} + C(u)u = -Gp + Du \rightarrow \frac{dU}{dt} + C(U)U = -GP + DU$$

- Block structures appear in  $C$ ,  $G$ ,  $D$
- SpMV's can be translated to SpMM!
  - Increases the arithmetic intensity

# Ensemble averaging parallel-in-time

$$\frac{du}{dt} + C(u)u = -Gp + Du \rightarrow \frac{dU}{dt} + C(U)U = -GP + \mathbb{D}U$$

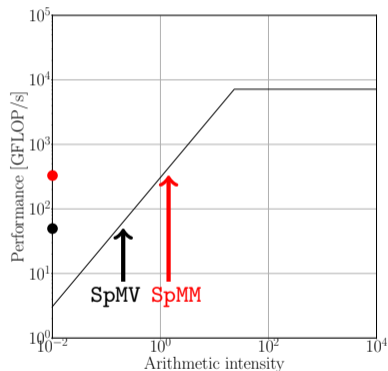
- Block structures appear in  $C, G, \mathbb{D}$
- SpMV's can be translated to SpMM!
  - Increases the arithmetic intensity



# Ensemble averaging parallel-in-time

$$\frac{du}{dt} + C(u)u = -Gp + Du \rightarrow \frac{dU}{dt} + C(U)U = -GP + DU$$

- Block structures appear in  $C, G, D$
- SpMV's can be translated to SpMM!
  - Increases the arithmetic intensity
  - Increases performance  $\rightarrow$  speed-up!





# Ensemble averaging parallel-in-time

## Generation of speed-up

- Speed-up in SpMV is guaranteed by increasing the AI
- How does this translate to the whole simulation and iteration?
  - $T_T$  will be simulated  $m$  times...
  - Speed-up in  $T_A$  has to be big enough!

---

<sup>2</sup>B.I. Krasnopolsky, "An approach for accelerating incompressible turbulent flow simulations based on simultaneous modelling of multiple ensembles," *Comput. Phys. Commun.* **229**, 2018

# Ensemble averaging parallel-in-time

## Generation of speed-up

- Speed-up in SpMV is guaranteed by increasing the AI
- How does this translate to the whole simulation and iteration?
  - $T_T$  will be simulated  $m$  times...
  - Speed-up in  $T_A$  has to be big enough!

## Times ratio $\beta$

$$\beta = \frac{T_A}{T_T}$$

## Estimation of simulation speed-up<sup>2</sup>

$$P_m = \frac{1 + \beta}{m + \beta} \frac{5m}{5m - 3\theta(m - 1)}$$

- Speed-up of the iteration,  $P_{m,ite}$
- Extension to the whole simulation

<sup>2</sup>B.I. Krasnopolsky, "An approach for accelerating incompressible turbulent flow simulations based on simultaneous modelling of multiple ensembles," *Comput. Phys. Commun.* **229**, 2018

# Ensemble averaging parallel-in-time

## Generation of speed-up

- Speed-up in SpMV is guaranteed by increasing the AI
- How does this translate to the whole simulation and iteration?
  - $T_T$  will be simulated  $m$  times...
  - Speed-up in  $T_A$  has to be big enough!

## Times ratio $\beta$

$$\beta = \frac{T_A}{T_T}$$

## Estimation of simulation speed-up<sup>2</sup>

$$P_m = \frac{1 + \beta}{m + \beta} \frac{5m}{5m - 3\theta(m - 1)}$$

- Speed-up of the iteration,  $P_{m,ite}$
- Extension to the whole simulation

<sup>2</sup>B.I. Krasnopolsky, "An approach for accelerating incompressible turbulent flow simulations based on simultaneous modelling of multiple ensembles," *Comput. Phys. Commun.* **229**, 2018

# Methodology

## Case definition

### Differentially heated cavity of aspect ratio 4

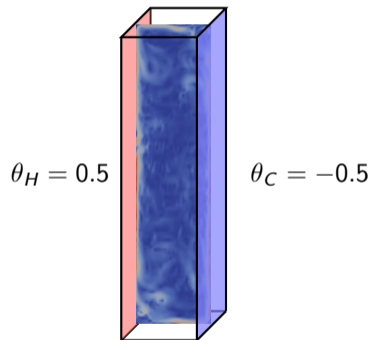
- Semi-discrete governing equations

$$M\mathbf{u}_s = 0_c,$$

$$\Omega \frac{d\mathbf{u}_c}{dt} + C(\mathbf{u}_s)\mathbf{u}_c - \frac{\text{Pr}}{\text{Ra}^{1/2}} D\mathbf{u}_c + \Omega G_c \mathbf{p}_c + \Omega \mathbf{f}_c = 0_c,$$

$$\Omega \frac{d\theta_c}{dt} + C(\mathbf{u}_s)\theta_c - \frac{1}{\text{Ra}^{1/2}} D\theta_c = 0_c$$

- 3rd-order Heun Runge-Kutta<sup>3</sup>, SAT<sup>4</sup>
- $\text{Ra}=10^{10}$ ,  $\text{Pr}=0.71$
- $\mathbf{f}_c = (0, \text{Pr}\theta, 0)$



<sup>3</sup>B. Sanderse and B. Koren, "Accuracy analysis of explicit Runge-Kutta methods applied to the incompressible Navier-Stokes equations," *J. Comput. Phys.* **231**, 2012

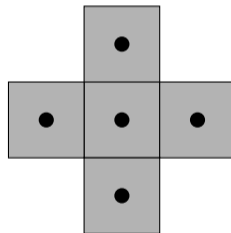
<sup>4</sup>J. Plana-Riu et al, "Cost-vs-accuracy analysis of self-adaptive time-integration methods," *10th THMT*, 2023

# Numerical tests in CPU

## Construction of test case

### Characteristics

- Run in 1 MN5 GPP-HighMem partition node:
  - 2x Intel Xeon Platinum 8480+ 56C 2GHz
  - 1024GB of RAM memory
  - 2x54 OpenMP threads within the socket
  - 2 MPI processes (1x socket)
- Mesh: 220x880x220 (42.6M cells)
  - 400k cells per CPU
- Run in three different discretizations for the Laplacian operator: 7p, 13p, and 27p
- Influence of Poisson solver iterations: (150, 350, 550)
- Run for 1, 2, 4, and 8 simultaneous flow states

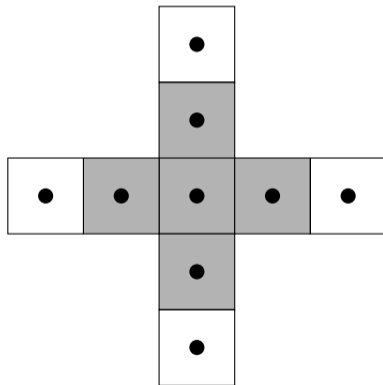


# Numerical tests in CPU

## Construction of test case

### Characteristics

- Run in 1 MN5 GPP-HighMem partition node:
  - 2x Intel Xeon Platinum 8480+ 56C 2GHz
  - 1024GB of RAM memory
  - 2x54 OpenMP threads within the socket
  - 2 MPI processes (1x socket)
- Mesh: 220x880x220 (42.6M cells)
  - 400k cells per CPU
- Run in three different discretizations for the Laplacian operator: 7p, 13p, and 27p
- Influence of Poisson solver iterations: (150, 350, 550)
- Run for 1, 2, 4, and 8 simultaneous flow states

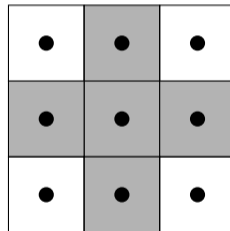


# Numerical tests in CPU

## Construction of test case

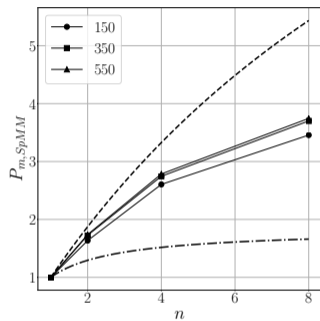
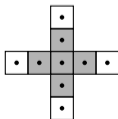
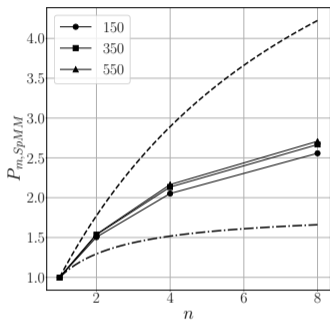
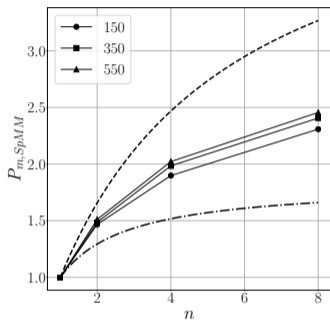
### Characteristics

- Run in 1 MN5 GPP-HighMem partition node:
  - 2x Intel Xeon Platinum 8480+ 56C 2GHz
  - 1024GB of RAM memory
  - 2x54 OpenMP threads within the socket
  - 2 MPI processes (1x socket)
- Mesh: 220x880x220 (42.6M cells)
  - 400k cells per CPU
- Run in three different discretizations for the Laplacian operator: 7p, 13p, and 27p
- Influence of Poisson solver iterations: (150, 350, 550)
- Run for 1, 2, 4, and 8 simultaneous flow states



# Numerical tests in CPU

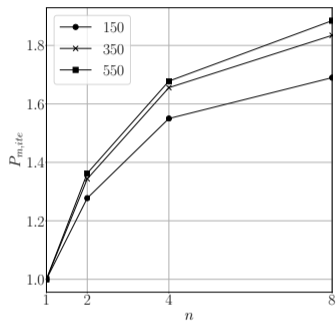
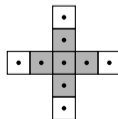
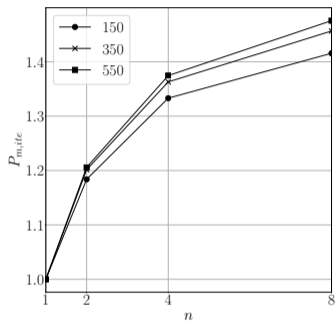
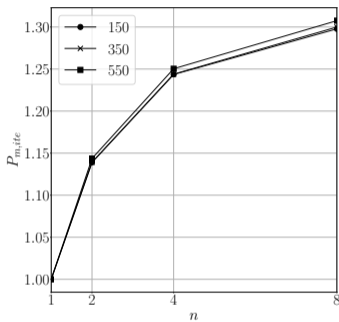
Results. Speed-up for SpMM





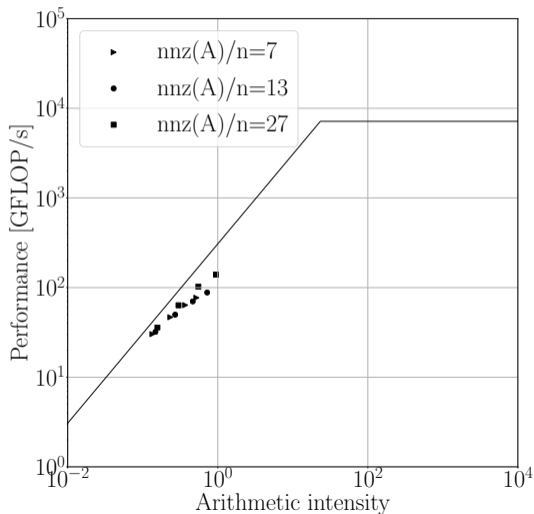
# Numerical tests in CPU

Results. Speed-up for iteration



# Numerical tests in CPU

Results. Roofline analysis



## System properties

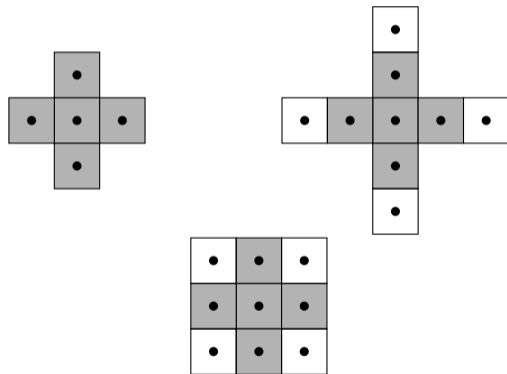
- 1x MareNostrum 5 GPP Node:
  - Peak performance:  $\approx 7155.86$  GFLOP/s
  - Memory bandwidth: 307.2 GB/s

# Numerical tests in GPU

## Construction of test case

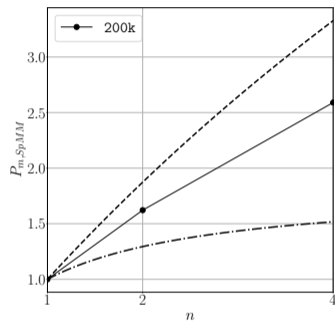
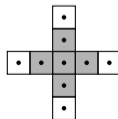
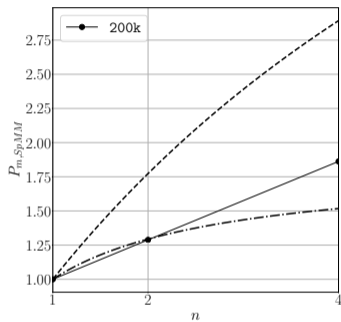
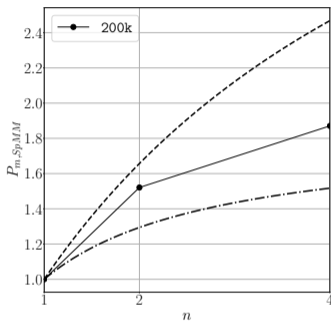
### Characteristics

- Run in 1 JFF GPU node:
  - 1x Nvidia A100 80GB PCIe
  - 2x Intel Xeon Gold 6442Y 48C
  - 1024GB of RAM memory
  - 96 OpenMP threads for preprocessing
  - 1 MPI process for computing + OpenCL
- Mesh: 176x704x176 (21M cells)
- Run in three different discretizations for the Laplacian operator: 7p, 13p, and 27p
- Run for 1, 2, 4 simultaneous flow states



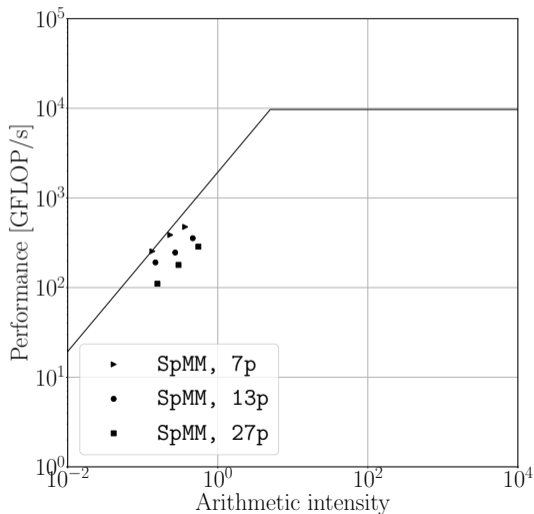
# Numerical tests in GPU

Results. Speed-up for SpMM



# Numerical tests in GPU

Results. Roofline analysis



## System properties

- 1x Nvidia A100 80GB PCIe:
  - Peak performance:  $\approx 9700$  GFLOP/s
  - Memory bandwidth: 1935 GB/s

# Conclusions

## Concluding remarks

- Method to exploit repeated block structures is presented, with PiT as an example
- Implementation in TFA+HPC<sup>2</sup> allows converting SpMV to SpMM without modifying call
- Ready to use in both CPU and GPU architectures
- Results with CPU and GPU as expected, within upper and lower bounds for all cases.
- As GPUs are loaded as much as possible, some latency issues reduces performance for higher nnz.