

# HIGH PERFORMANCE UNSTRUCTURED SPMM COMPUTATION USING TENSOR CORES

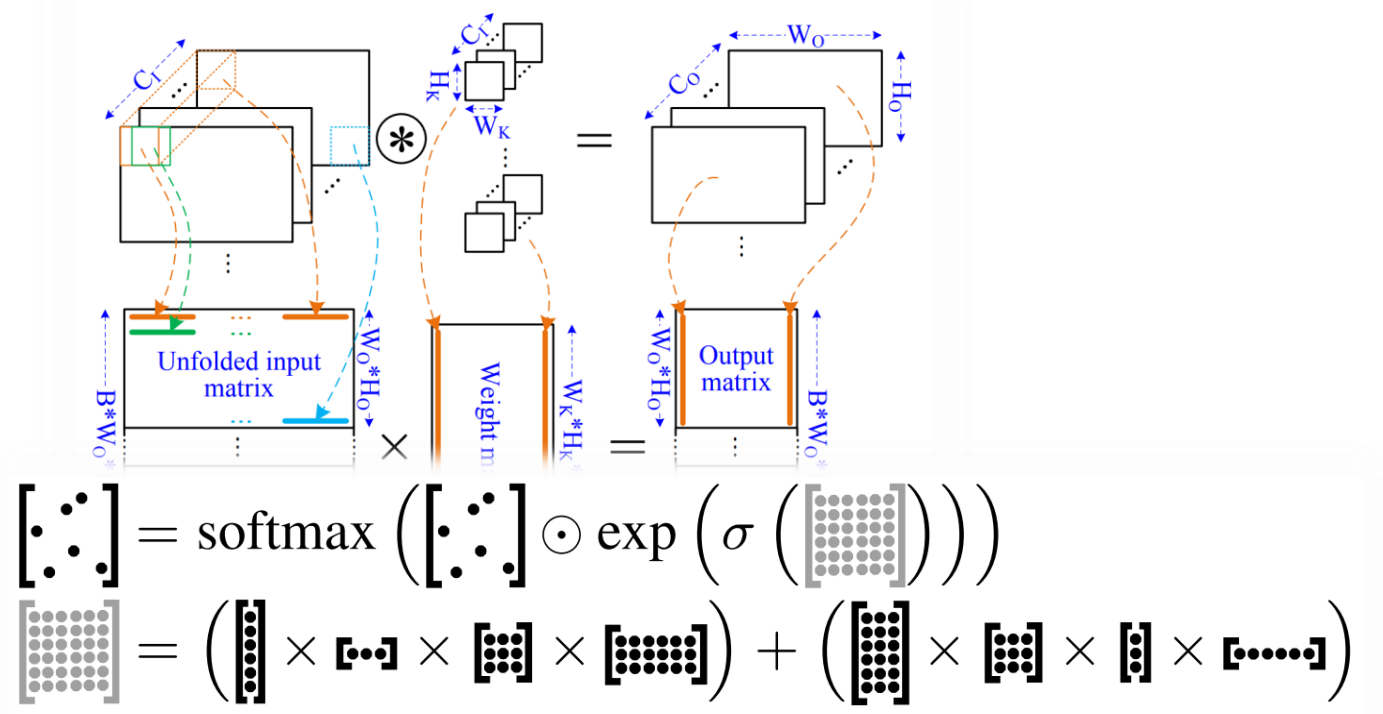
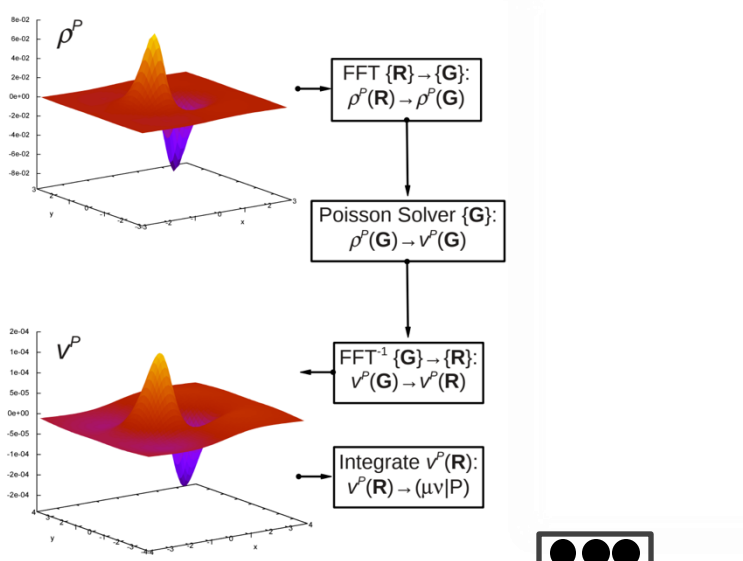
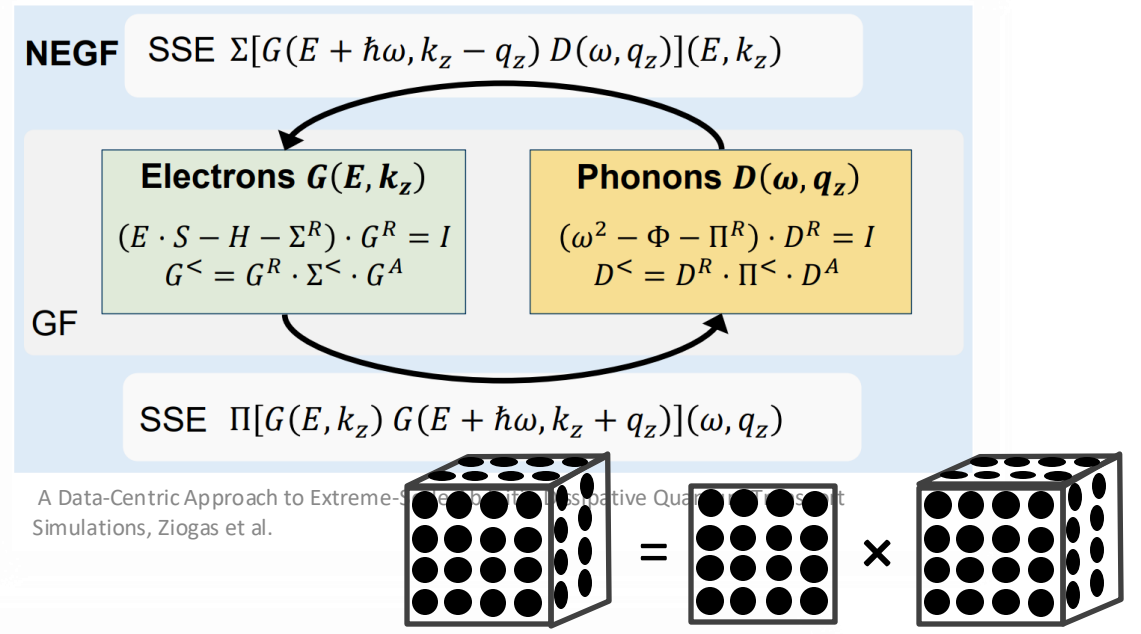
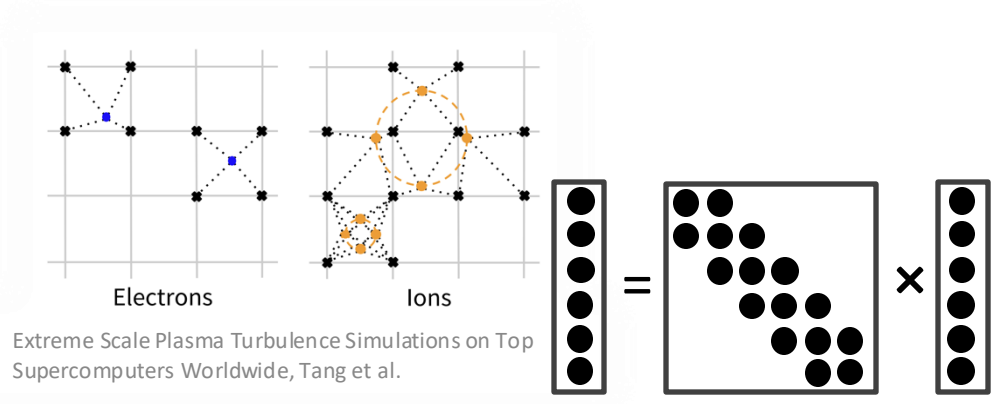
Patrik Okanovic<sup>1</sup>, Grzegorz Kwasniewski<sup>1</sup>, Paolo Sylos Labini<sup>2</sup>, Maciej Besta<sup>1</sup>,  
Flavio Vella<sup>3</sup>, Torsten Hoefler<sup>1</sup>

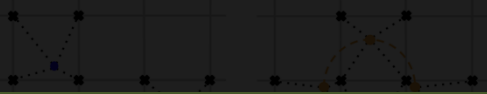
<sup>1</sup> ETH Zurich

<sup>2</sup> Free University of Bozen-Bolzano

<sup>3</sup> University of Trento



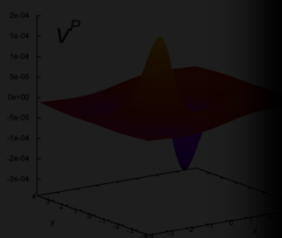




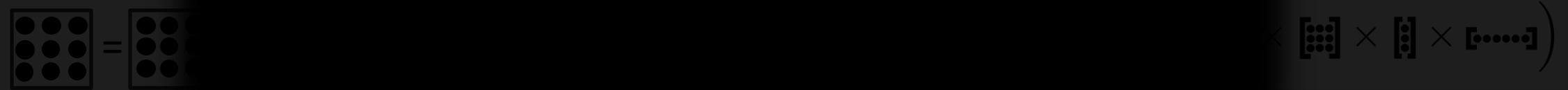
NEGF SSE  $\Sigma[G(E + \hbar\omega, k_z - q_z) D(\omega, q_z)](E, k_z)$

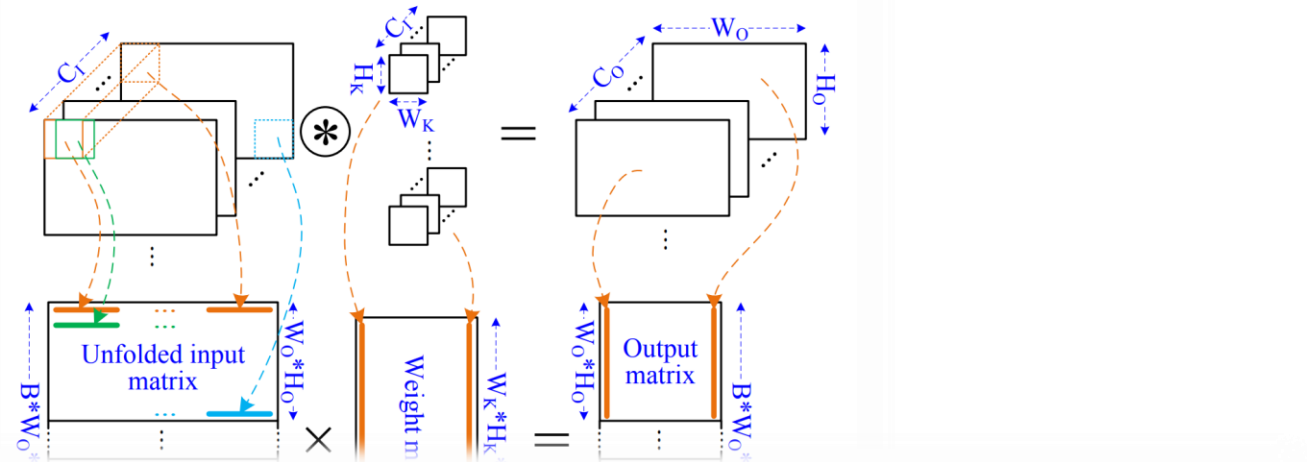
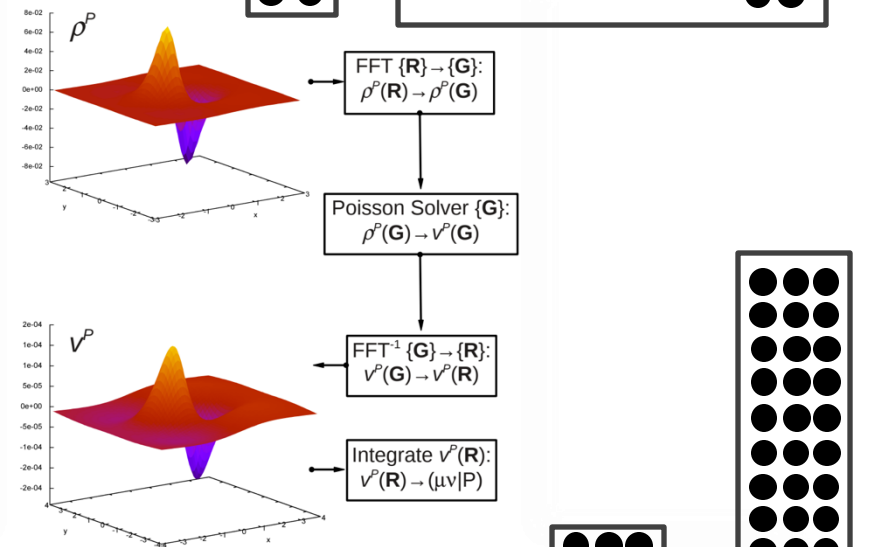
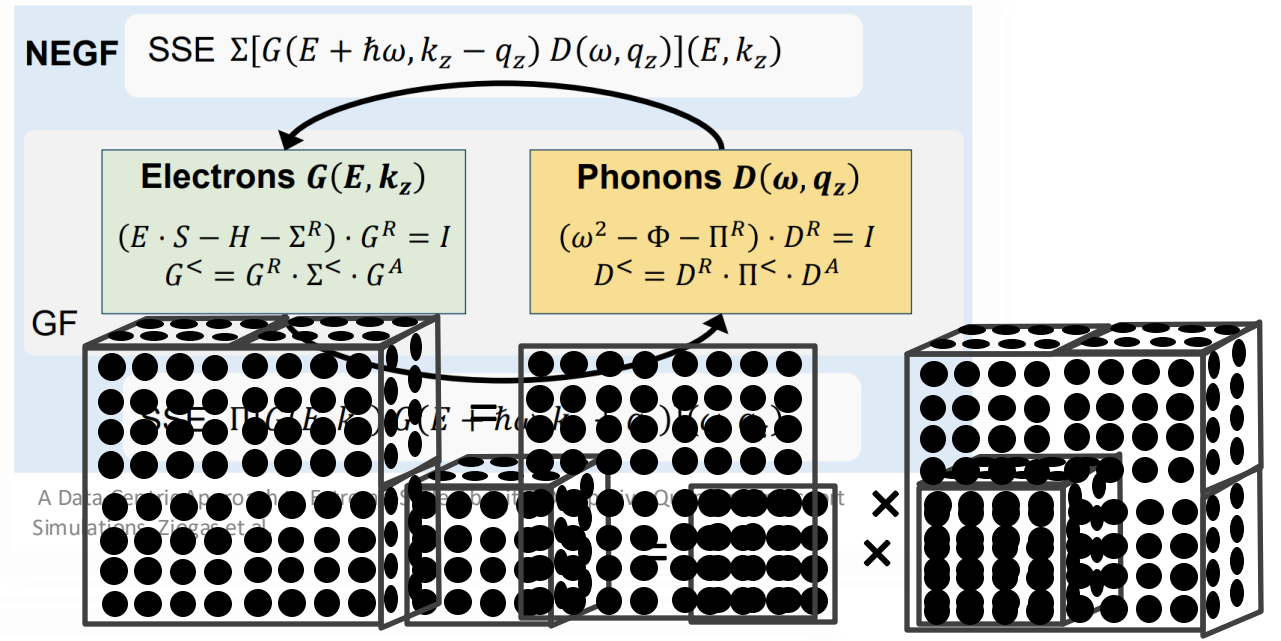
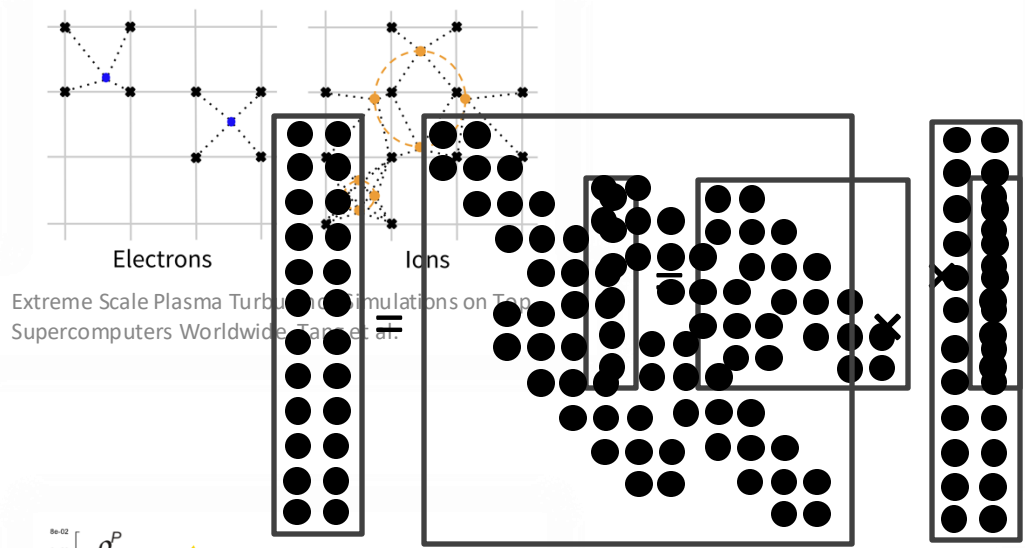
*Ideal for hardware!*

*...All done?...*



Enabling Simulation at the Fifth Rung of DFT: From Problem to Solution, Del Ben et al.





Enabling Simulation at the Fifth Rung of DFT: Large Scale RPA Calculations with Sub-Linear Time to Solution, Del Ben et al.

$$[\cdot \cdot \cdot] = \text{softmax} \left( [\cdot \cdot \cdot] \odot \exp \left( \sigma \left( \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} \right) \right) \right)$$

$$\begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} = \left( \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \times \begin{bmatrix} \cdot & \cdot & \cdot \end{bmatrix} \times \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} \times \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} \right) + \left( \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} \times \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} \times \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \times \begin{bmatrix} \cdot & \cdot & \cdot \end{bmatrix} \right)$$

$$\text{NEGF} \quad \text{SSE} \quad \Sigma[G(E + \hbar\omega, k_z - q_z) D(\omega, q_z)](E, k_z)$$

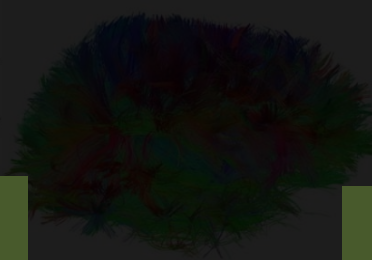
$$\Sigma^R \cdot G^R = I$$

$$\Sigma^< \cdot G^A$$

Phonons  $D(\omega, q_z)$

$$(\omega^2 - \Phi - \Pi^R) \cdot D^R = I$$

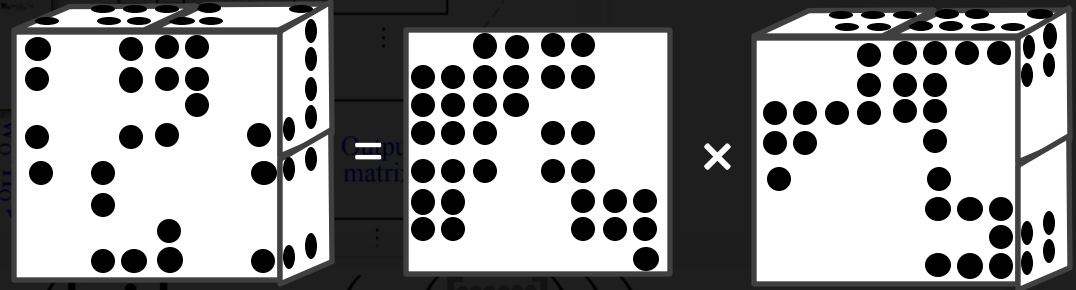
$$D^< = D^R \cdot \Pi^< \cdot D^A$$



Brain/neurons



Social

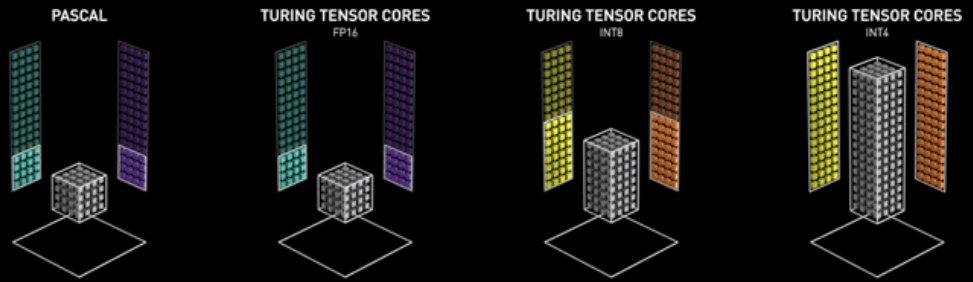


$$= \text{softmax} \left( \left[ \cdot \cdot \cdot \right] \odot \exp \left( \sigma \left( \left[ \cdot \cdot \cdot \right] \right) \right) \right)$$

$$= \left( \left[ \cdot \cdot \cdot \right] \times \left[ \cdot \cdot \cdot \right] \times \left[ \cdot \cdot \cdot \right] \times \left[ \cdot \cdot \cdot \right] \right) + \left( \left[ \cdot \cdot \cdot \right] \times \left[ \cdot \cdot \cdot \right] \times \left[ \cdot \cdot \cdot \right] \times \left[ \cdot \cdot \cdot \right] \right)$$

Hardware requirement:  
Dense matrices

Problems to solve:  
Large, sparse data  
structures



## ...but... we do have solutions... right?

### Libraries exist!

- *cuSPARSE*
- *hipSPARSE*

### Hardware exists!

- *Sparse Tensor Cores*

### Dedicated solutions exist!

- *DASP* [2]
- *Magicube* [3]
- *VENOM* [4]
- *cuSPARSE/Lt* [5]
- ...

### cuSPARSE [1]:

Up to 73,350  
(73 thousand!)  
times slower than the  
achievable peak!

### Sparsity requirements:

2:4  
*EVERY* four elements  
*EXACTLY* two nonzeros

### Lacking performance:

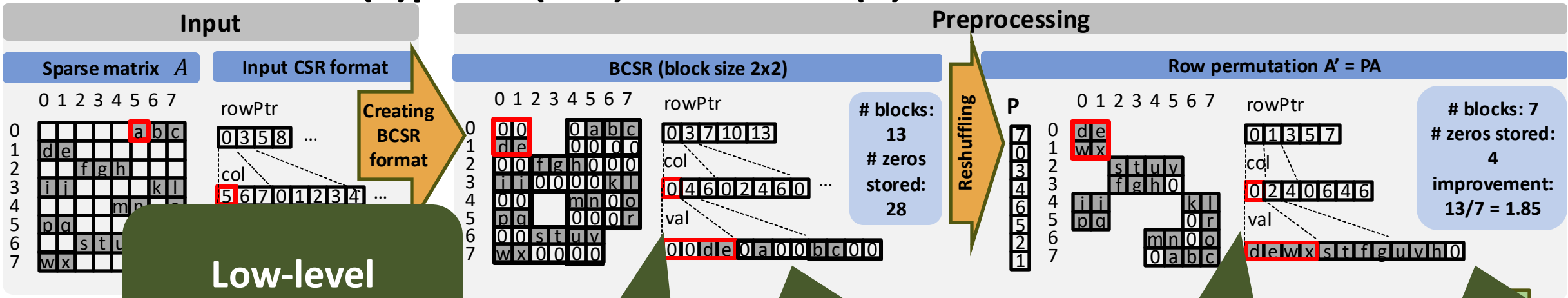
*At least 200x slower than the  
peak*

### Narrow applicability:

Only machine learning,  
“dense” sparse matrices (up  
to 80-90% sparsity)

[1] cuSPARSE v12.0 vs CUBLAS c12.0 on dense matrices, NVIDIA A100

# SMaT: (S)parse (Ma)trix Matrix (T)ensor Core-accelerate



Low-level optimizations?

Arbitrary unstructured sparsity

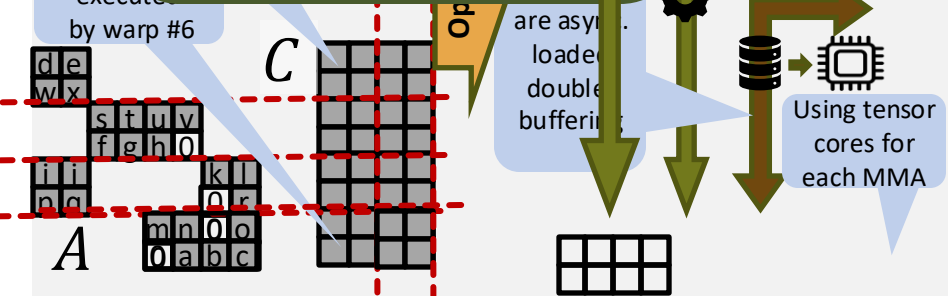
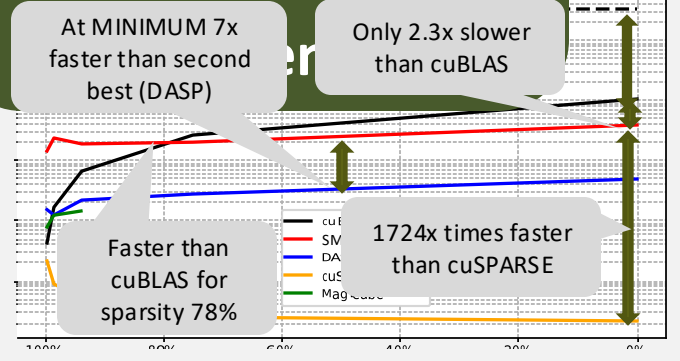
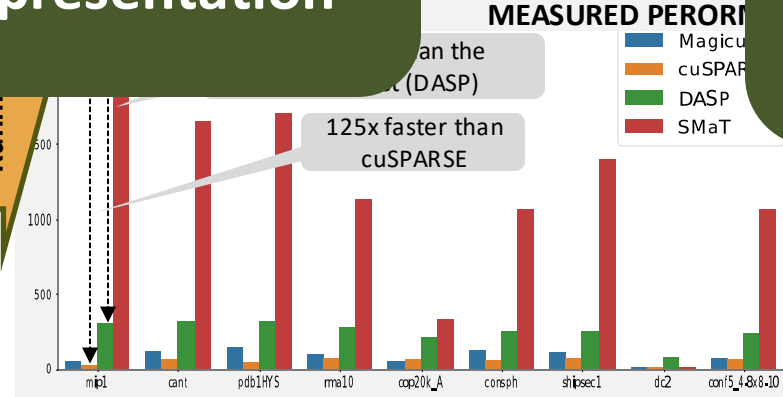
Hardware-friendly format?

0 representation

SmaT is FAST. But why?

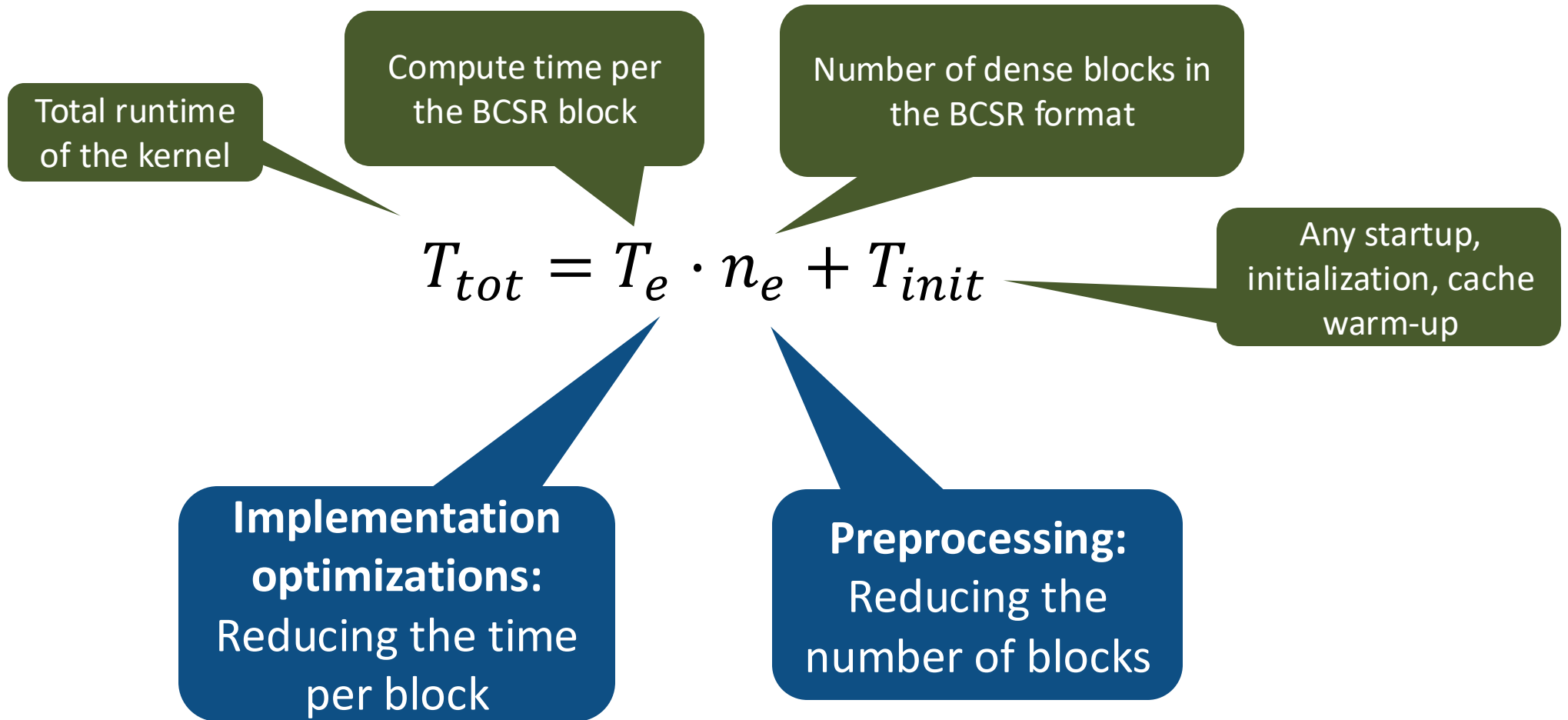
Efficient preprocessing?

Maximize compute





# Performance Model

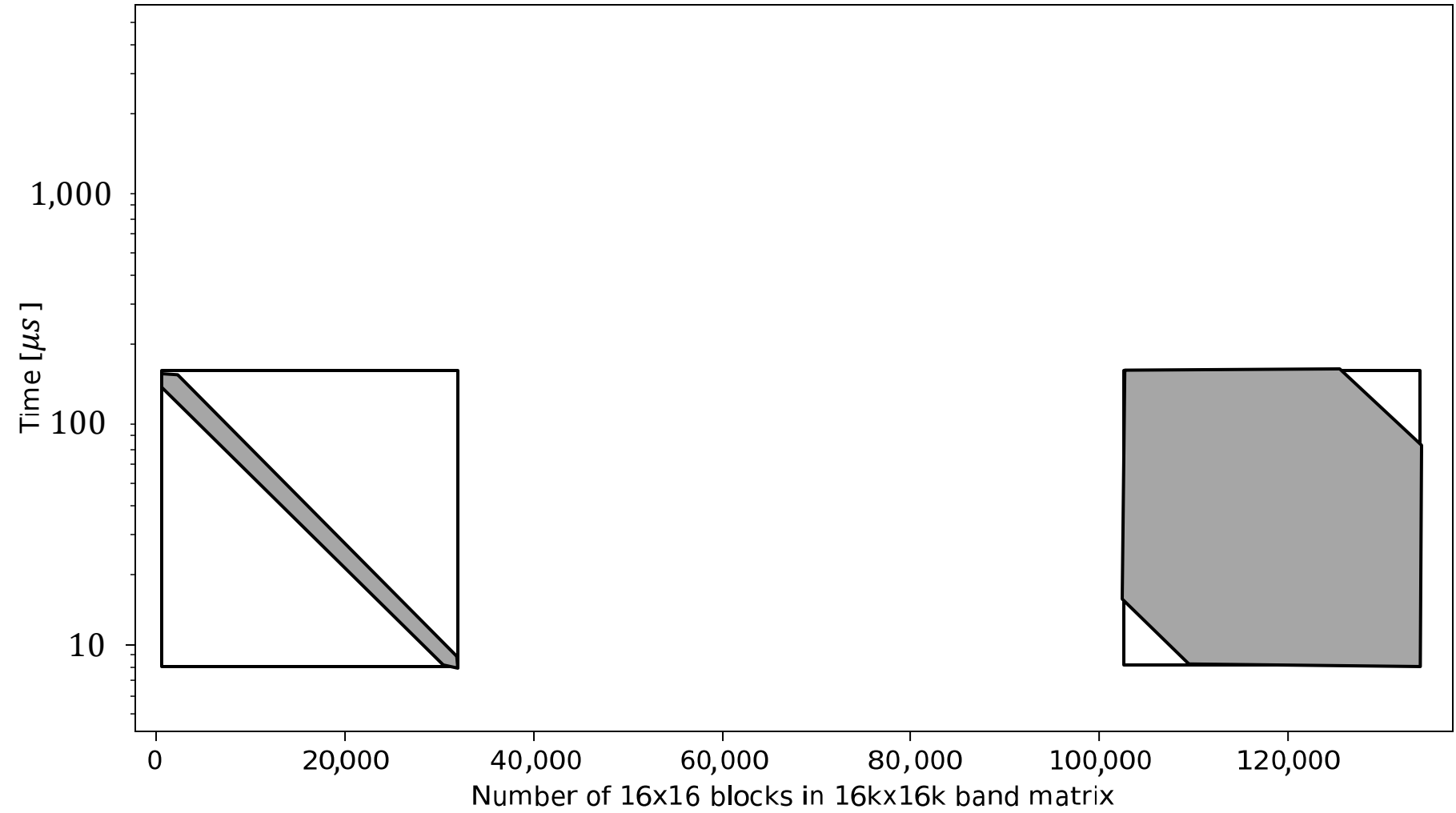
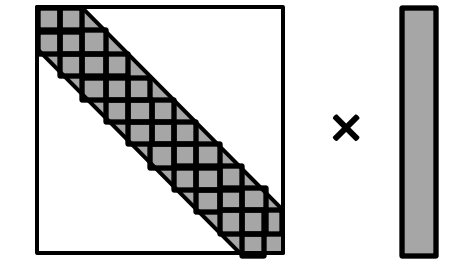




# Performance Model

$$T_{tot} = T_e \cdot n_e + T_{init}$$

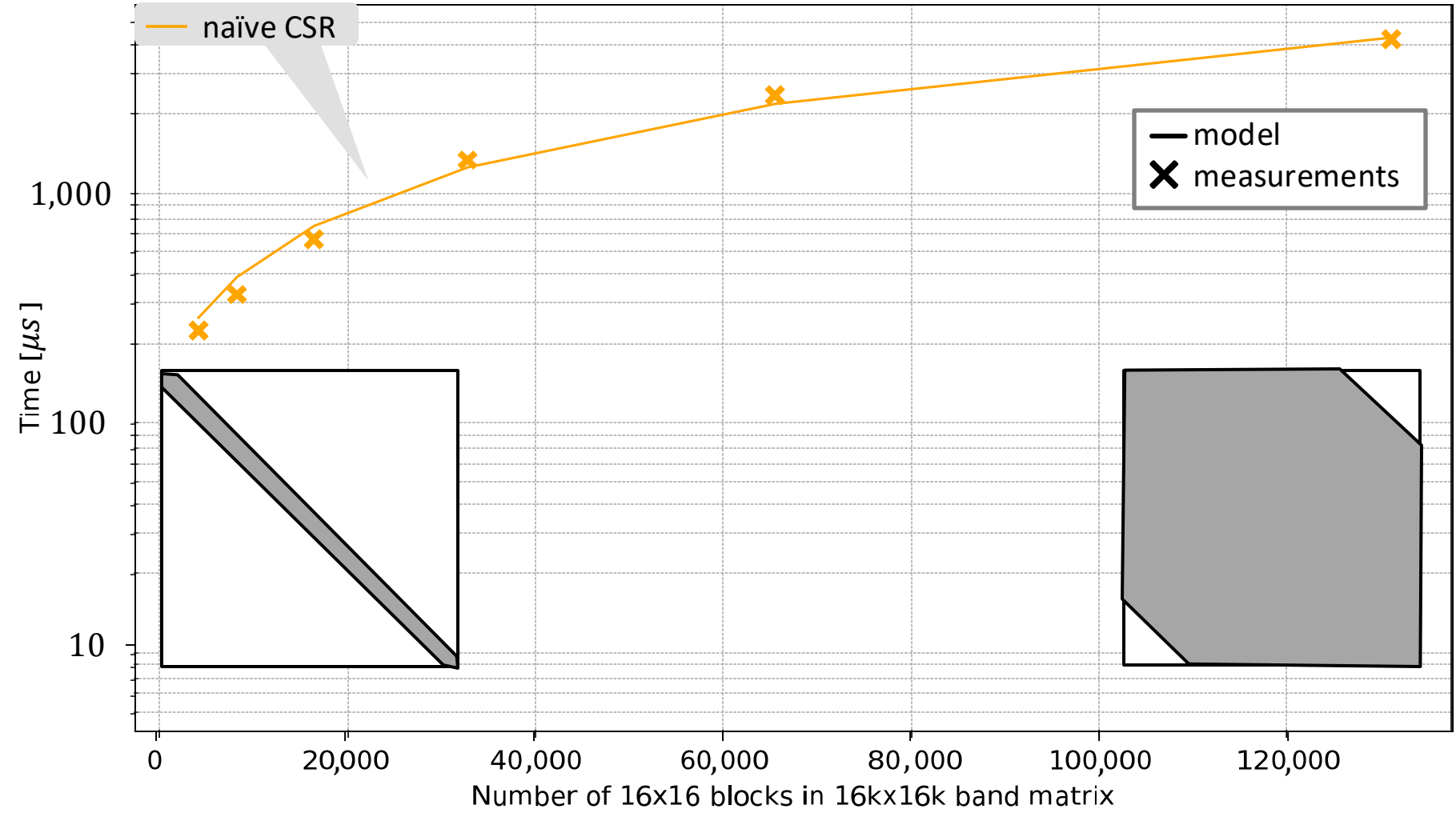
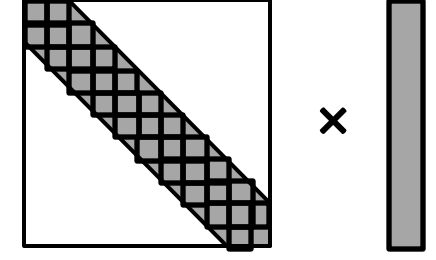
Number of columns in dense matrix  $N = 8$



# Performance Model

$$T_{tot} = T_e \cdot n_e + T_{init}$$

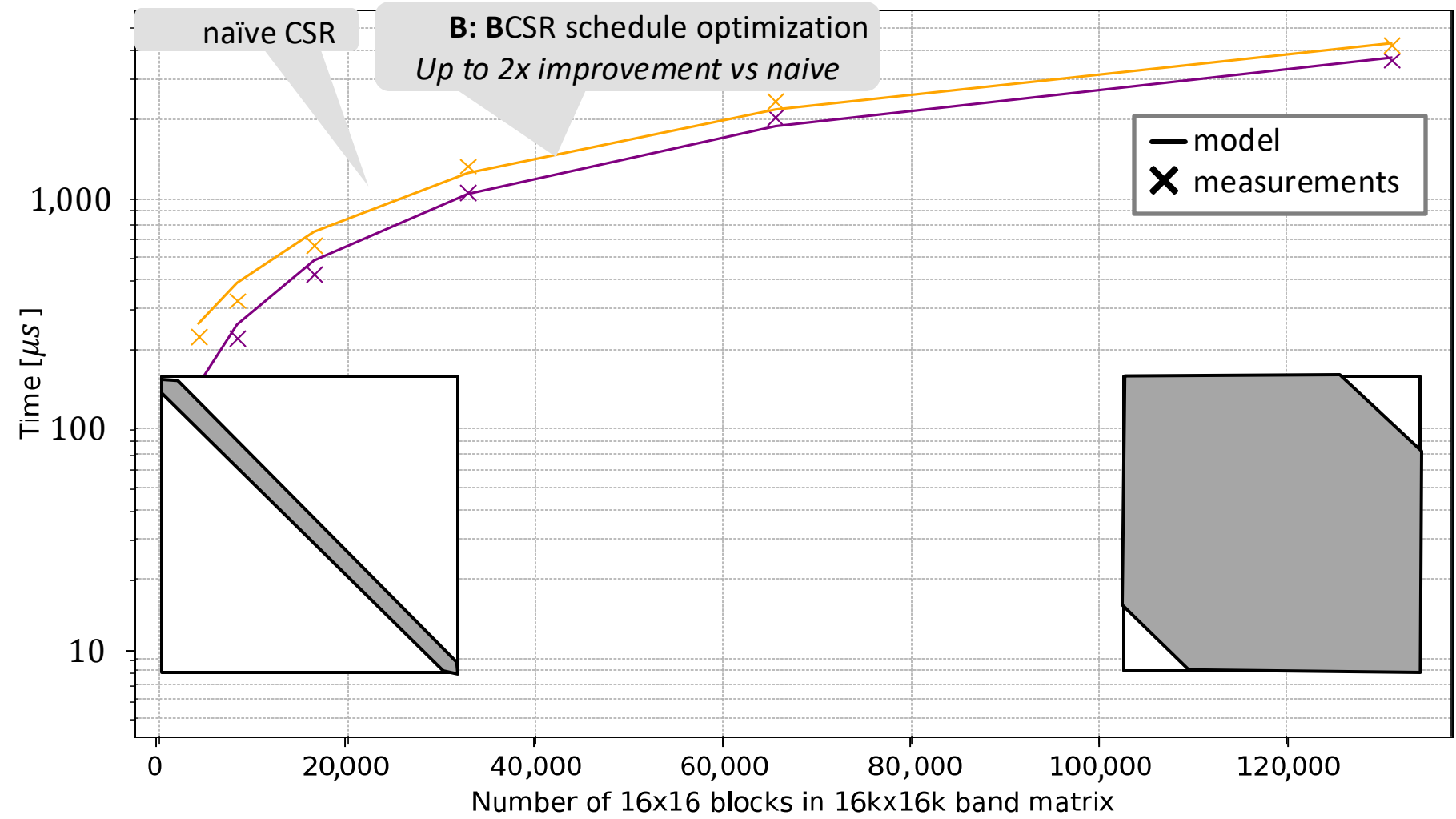
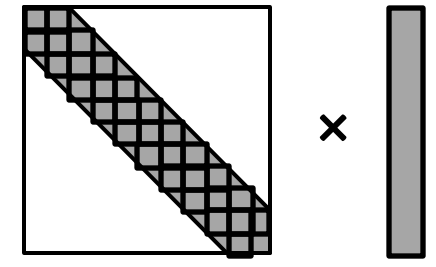
Number of columns in dense matrix  $N = 8$



# Performance Model

$$T_{tot} = T_e \cdot n_e + T_{init}$$

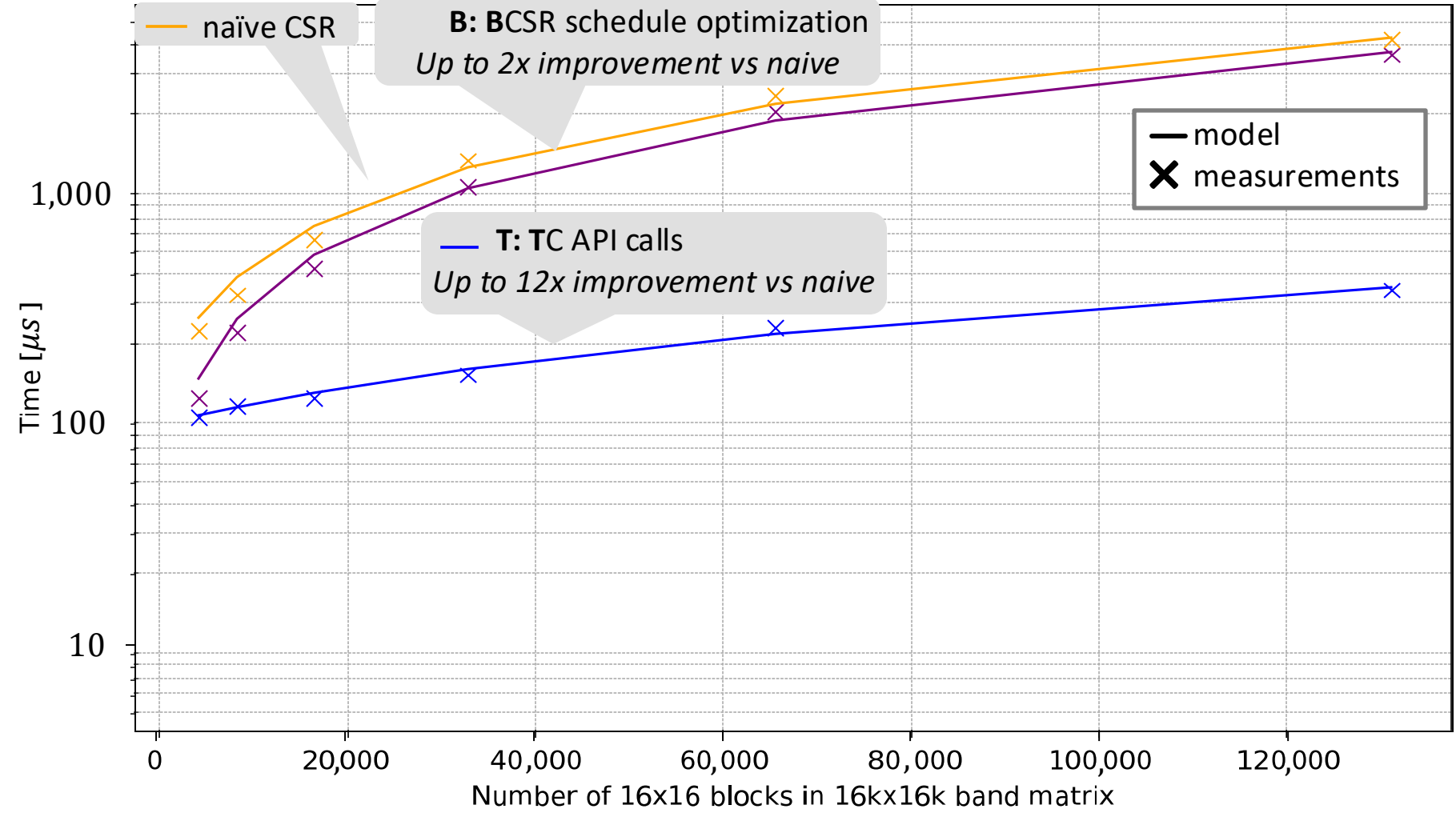
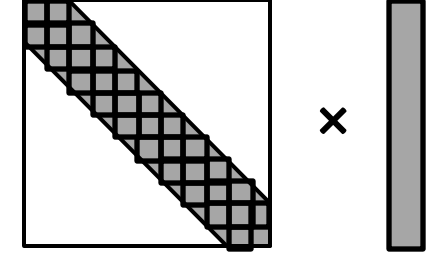
Number of columns in dense matrix  $N = 8$



# Performance Model

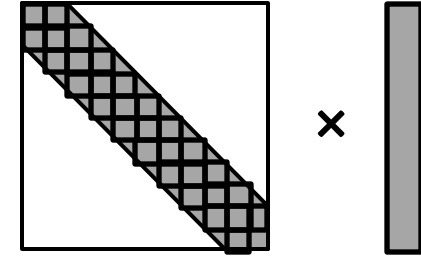
$$T_{tot} = T_e \cdot n_e + T_{init}$$

Number of columns in dense matrix  $N = 8$

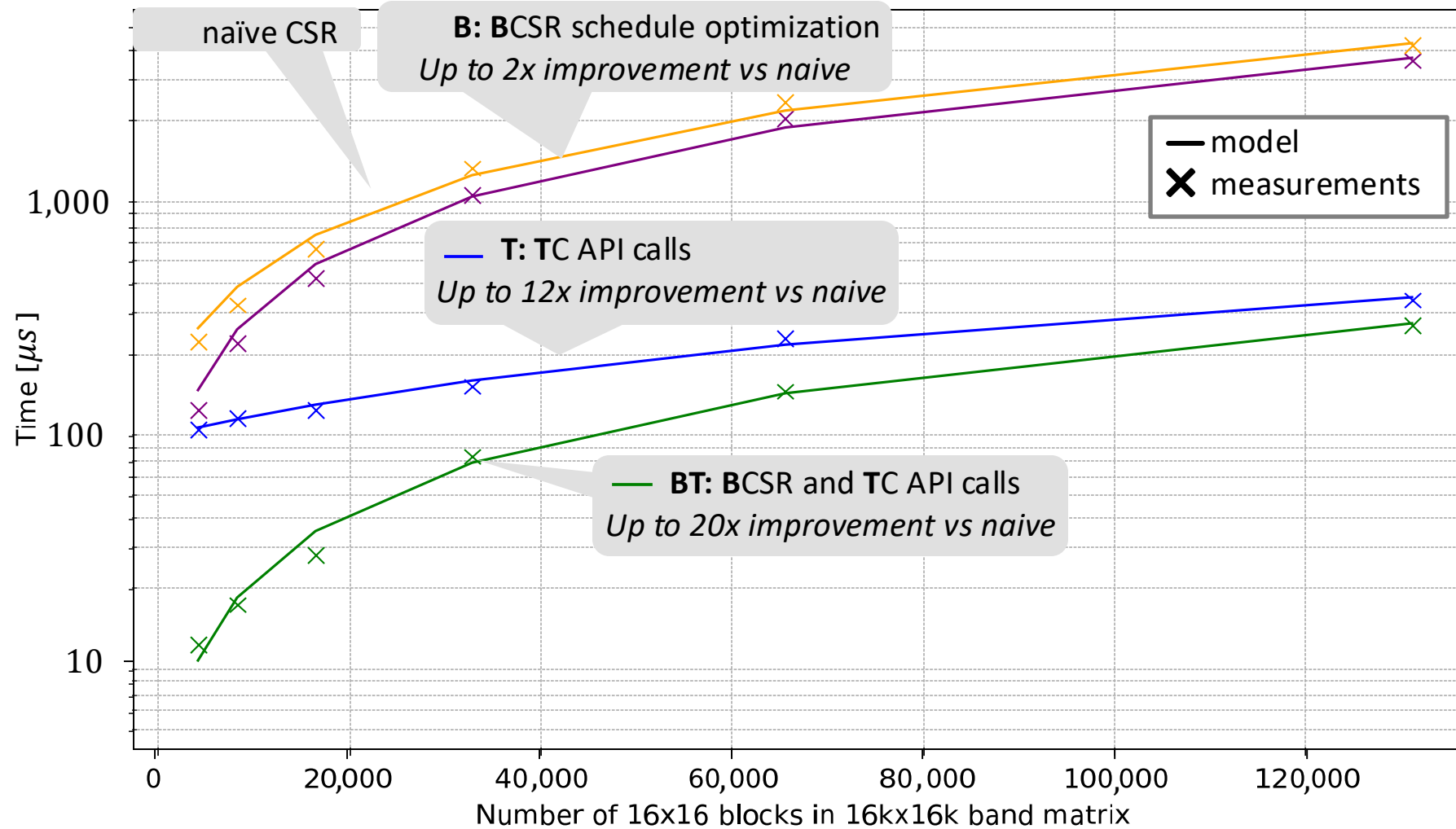


# Performance Model

$$T_{tot} = T_e \cdot n_e + T_{init}$$

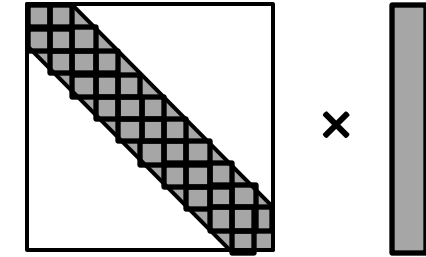


Number of columns in dense matrix  $N = 8$

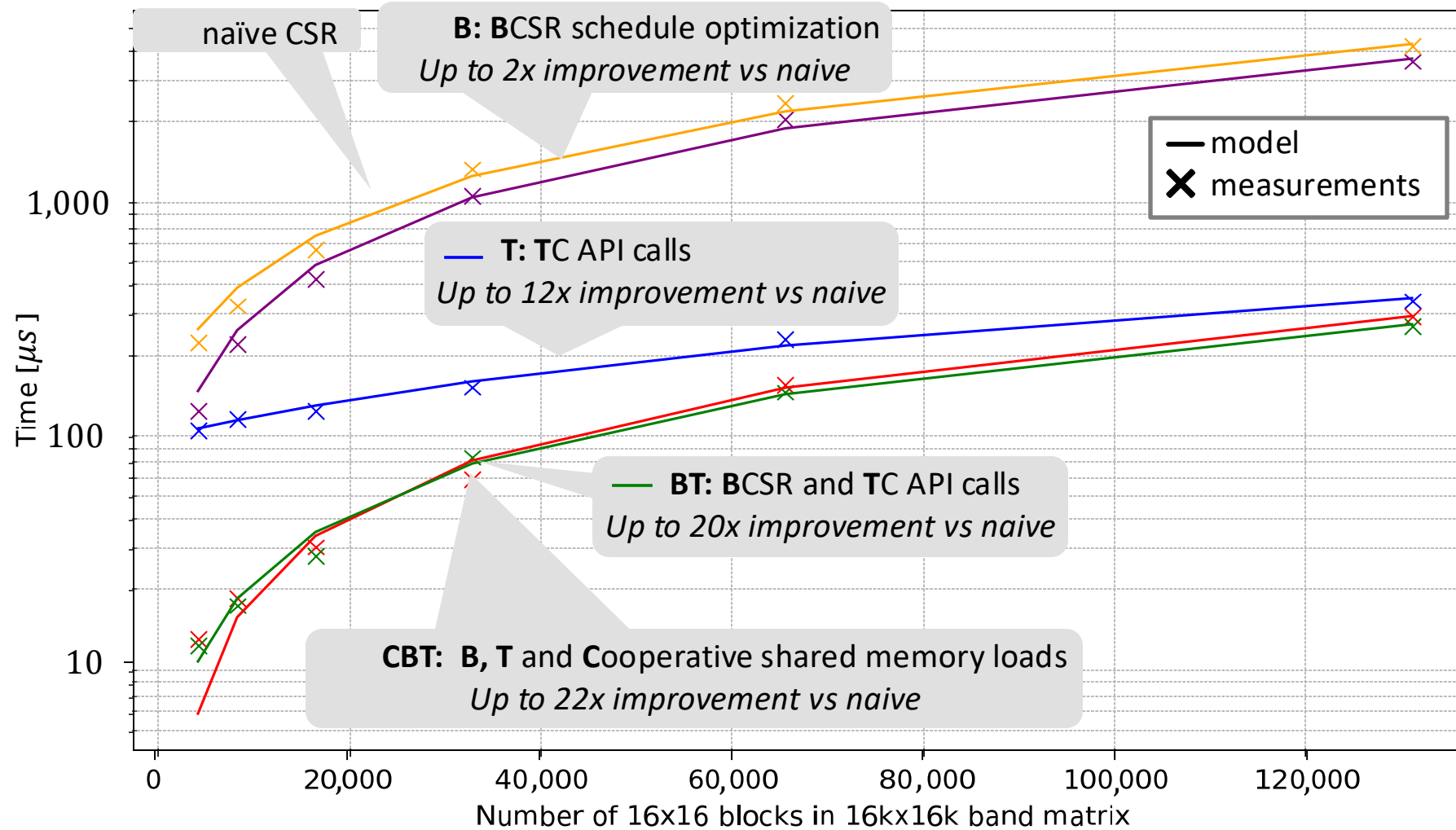


# Performance Model

$$T_{tot} = T_e \cdot n_e + T_{init}$$



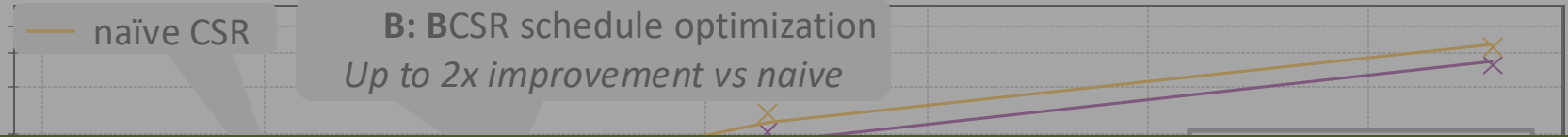
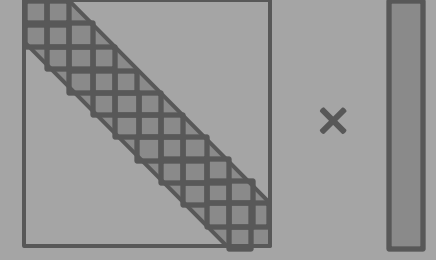
Number of columns in dense matrix  $N = 8$



# Performance Model

$$T_{tot} = T_e \cdot n_e + T_{init}$$

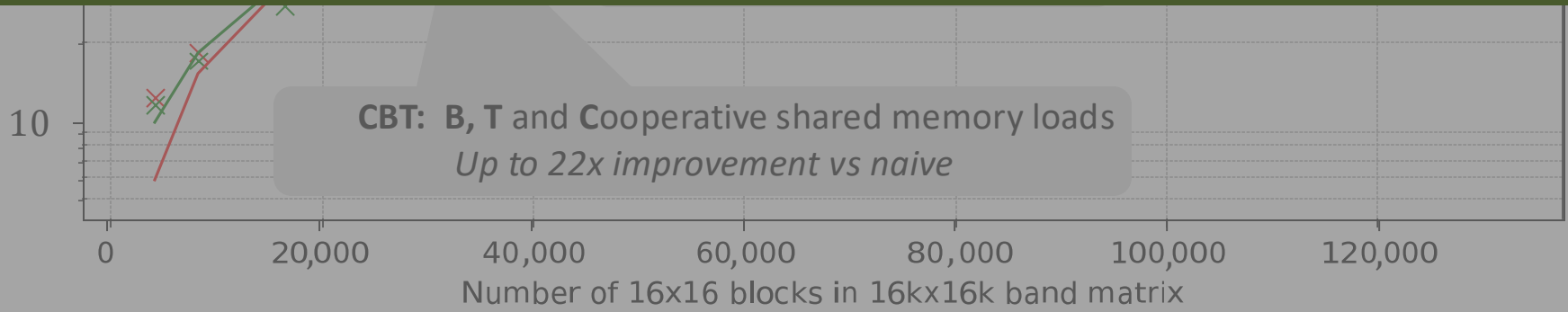
Number of columns in dense matrix  $N = 8$



**Low-level kernel optimizations are more important than high-level preprocessing!**

**Optimal matrix permutation:** up to 2.5x reduction in number of blocks (*across Suitesparse matrices*)

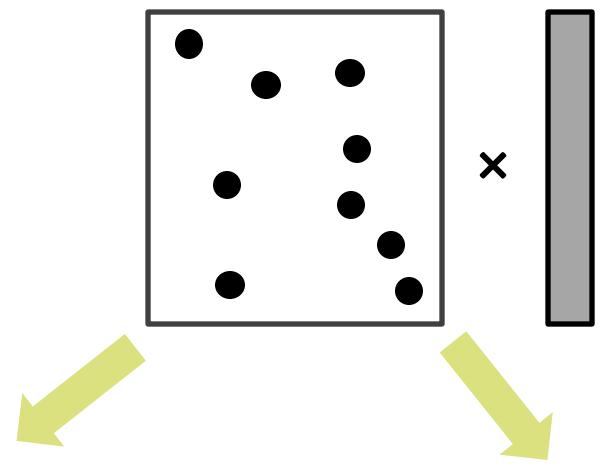
**Performance improvement: naïve CSR vs CBT up to 22x speedup**





# Evaluation

**State-of-the-art benchmark:**  
 Large, widely-used repository of sparse matrices from real-world applications



**Motivation:**  
 1. Remove nondeterminism and isolate performance-critical aspects:
 

- Asynchronous, pipelined loads
- Compute efficiency (tensor cores)
- Stress test on “dense” sparse matrices

 2. Unstructured sparsity kernels on structured sparsity matrices: *e.g., HPCG*

## SuiteSparse Collection

Domain	Name	Size	nnz	Sparsity
optimization	mip1	66K × 66K	10.4M	99.76%
quantum chem.	conf5_4-8x8	49K × 49K	1.9M	99.92%
2D/3D mesh	cant	62K × 62K	4M	99.89%
weighted graph	pdb1HYS	36K × 36K	4.3M	99.67%
fluid dynamics	rma10	46.8K × 46.8K	2.3M	99.89%
2D/3D mesh	cop20k_A	121K × 121K	2.6M	99.98%
2D/3D mesh	consph	83K × 83K	6M	99.91%
structural	shipsec1	140K × 140K	7.8M	99.96%
circuit simulation	dc2	116K × 116K	766K	99.99%

## Synthetic Band Matrices

# Preprocessing reordering

BCSR (b)

	0	1	2	3	4	5	6	7
0	0	0			0	a	b	c
1	d	e			0	0	0	0
2	0	0	f	g	h	0	0	0
3	i	i	0	0	0	0	k	l
4	0	0			m	n	0	o
5	p	q			0	0	0	r
6	0	0	s	t	u	v		
7	w	x	0	0	0	0		

Many blocks are almost

**DISCLAIMER:**

The preprocessing time is **not** included in the performance measurements

- Creating BCSR format:
  - Less than 1% of the runtime
  - Fair comparison: related works also measure only compute time

$A' = PA$

# blocks: 7  
# zeros stored: 4  
improvement:  $13/7 = 1.85$

4 6

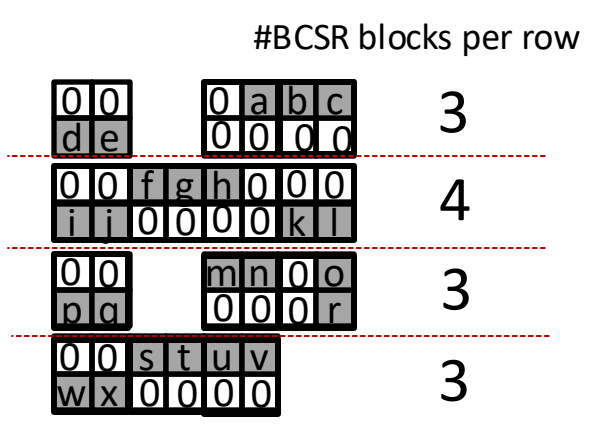
t f g u v h 0

pooling-based bit masking  
clustering based on  
rd distance

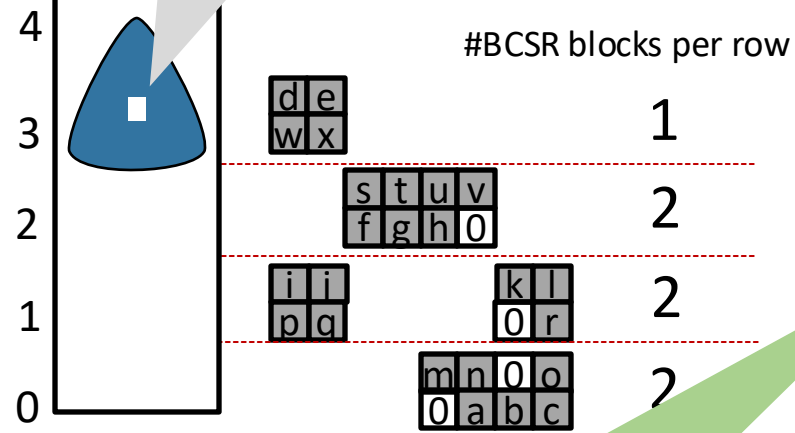
# Preprocessing reordering

Arithmetic mean = 3.25

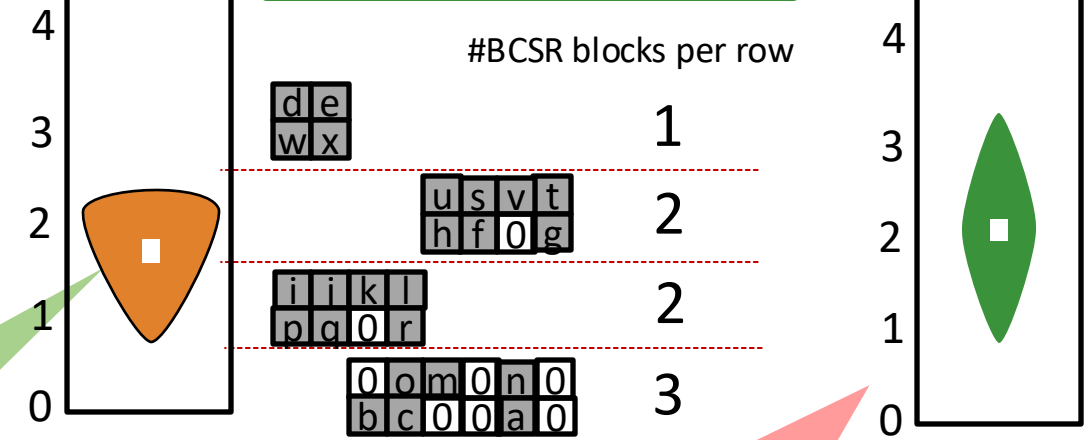
## Original input



## Row permutation



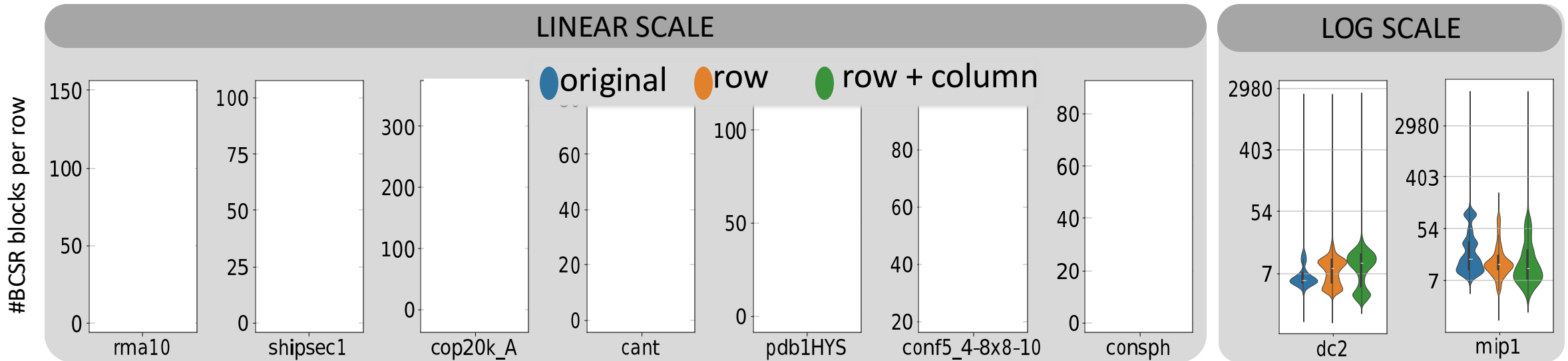
## Row + column permutation





**Smaller average #blocks/row:**  
 👍 less work  
 👍 smaller memory footprint



**Greater variance:**  
 👎 worse load balancing



# Preprocessing reordering



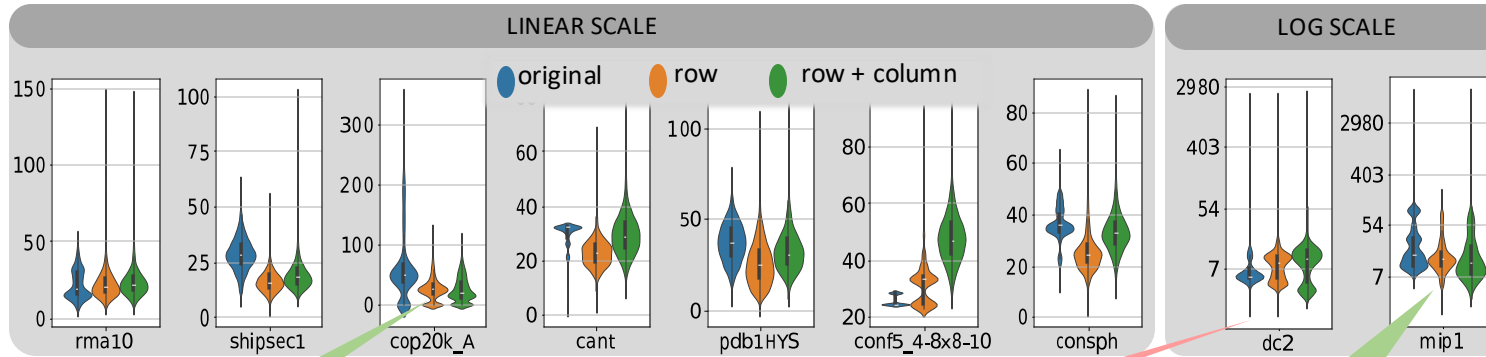
Distribution of the #BCSR blocks per row

 2.5x fewer blocks  
 3x smaller standard deviation

 99.994% sparsity  
 Coefficient of variation: 10.9

 1.8x fewer blocks  
 8.4x smaller standard deviation

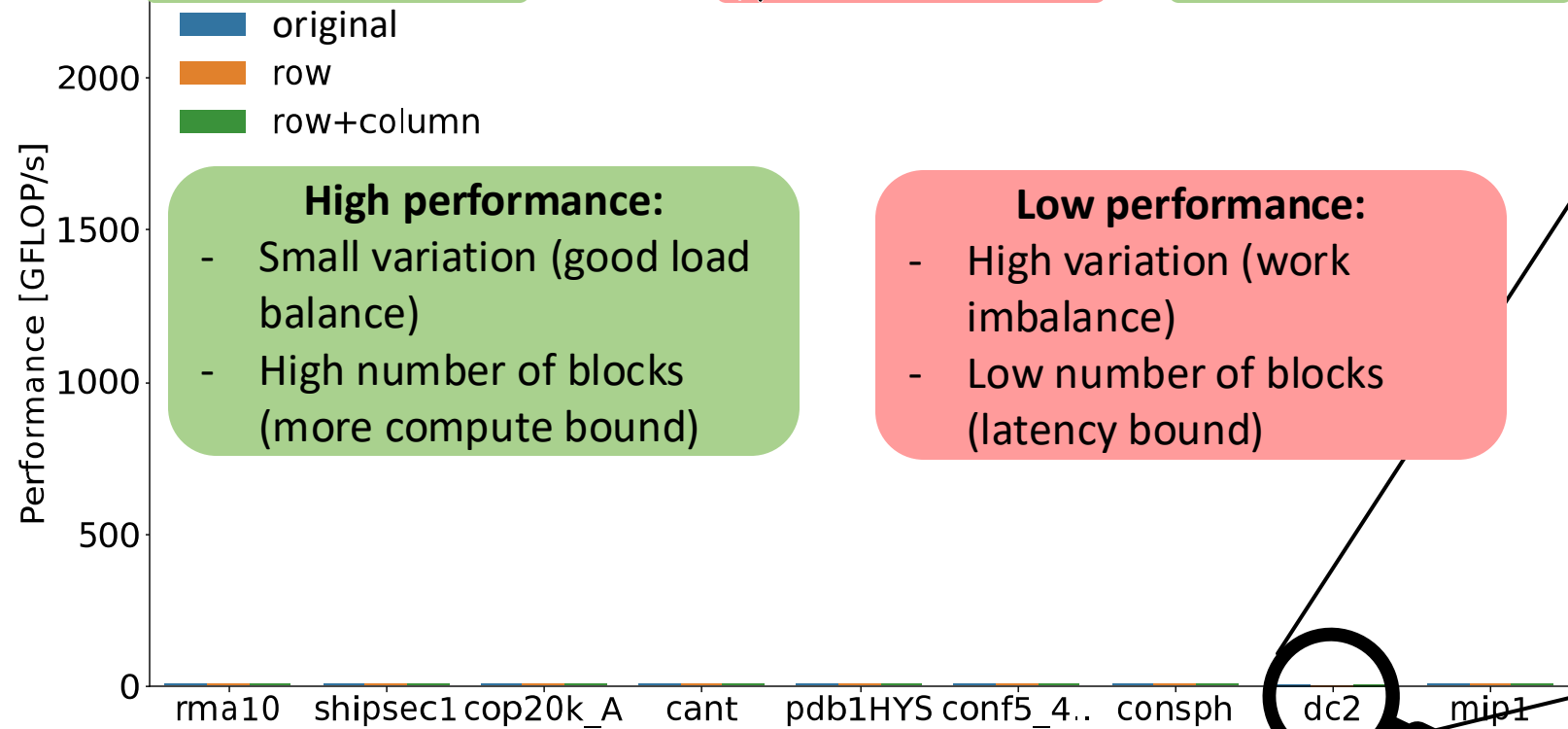
# Preprocessing reordering



👍 2.5x fewer blocks  
 👍 3x smaller standard deviation

👎 99.994% sparsity  
 👎 Coefficient of variation: 10.9

👍 1.8x fewer blocks  
 👍 8.4x smaller standard deviation

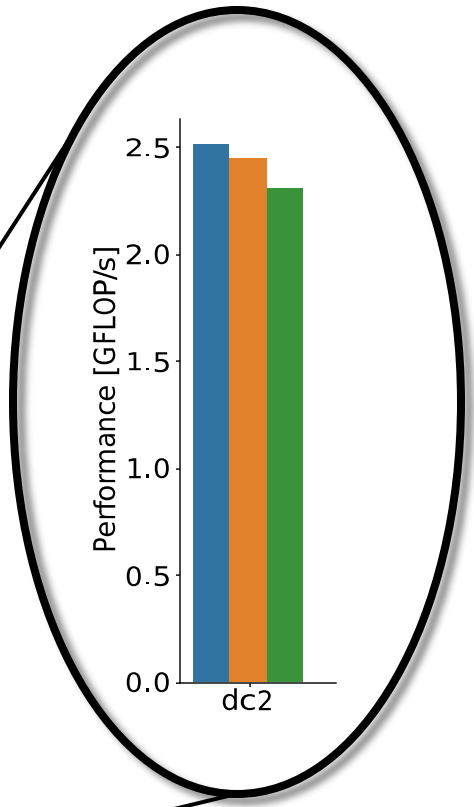


**High performance:**

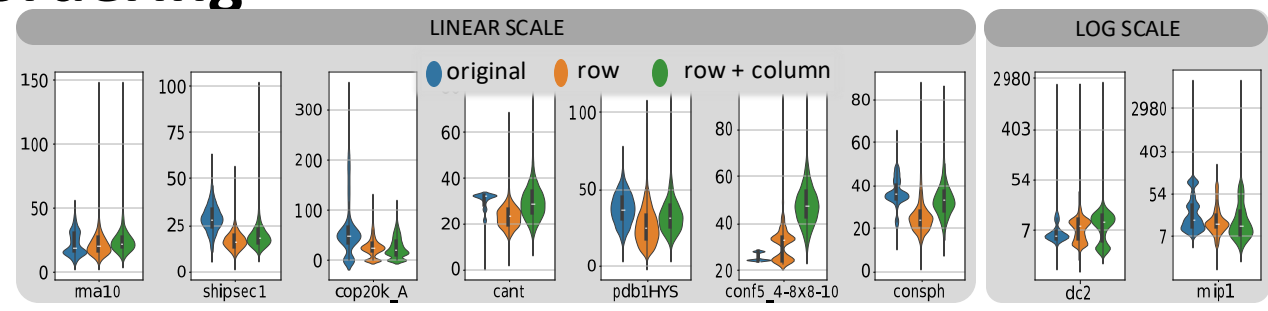
- Small variation (good load balance)
- High number of blocks (more compute bound)

**Low performance:**

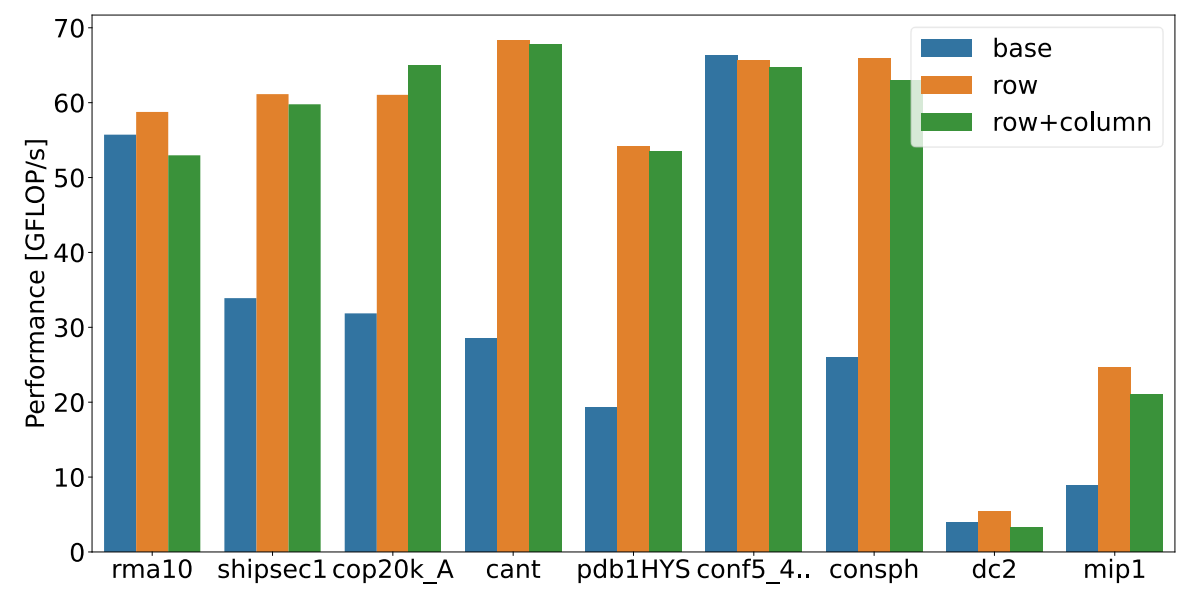
- High variation (work imbalance)
- Low number of blocks (latency bound)



# Preprocessing reordering

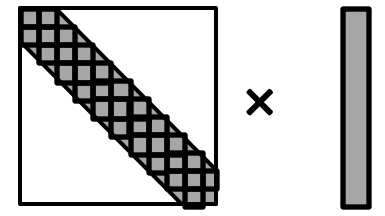


**MIPS**

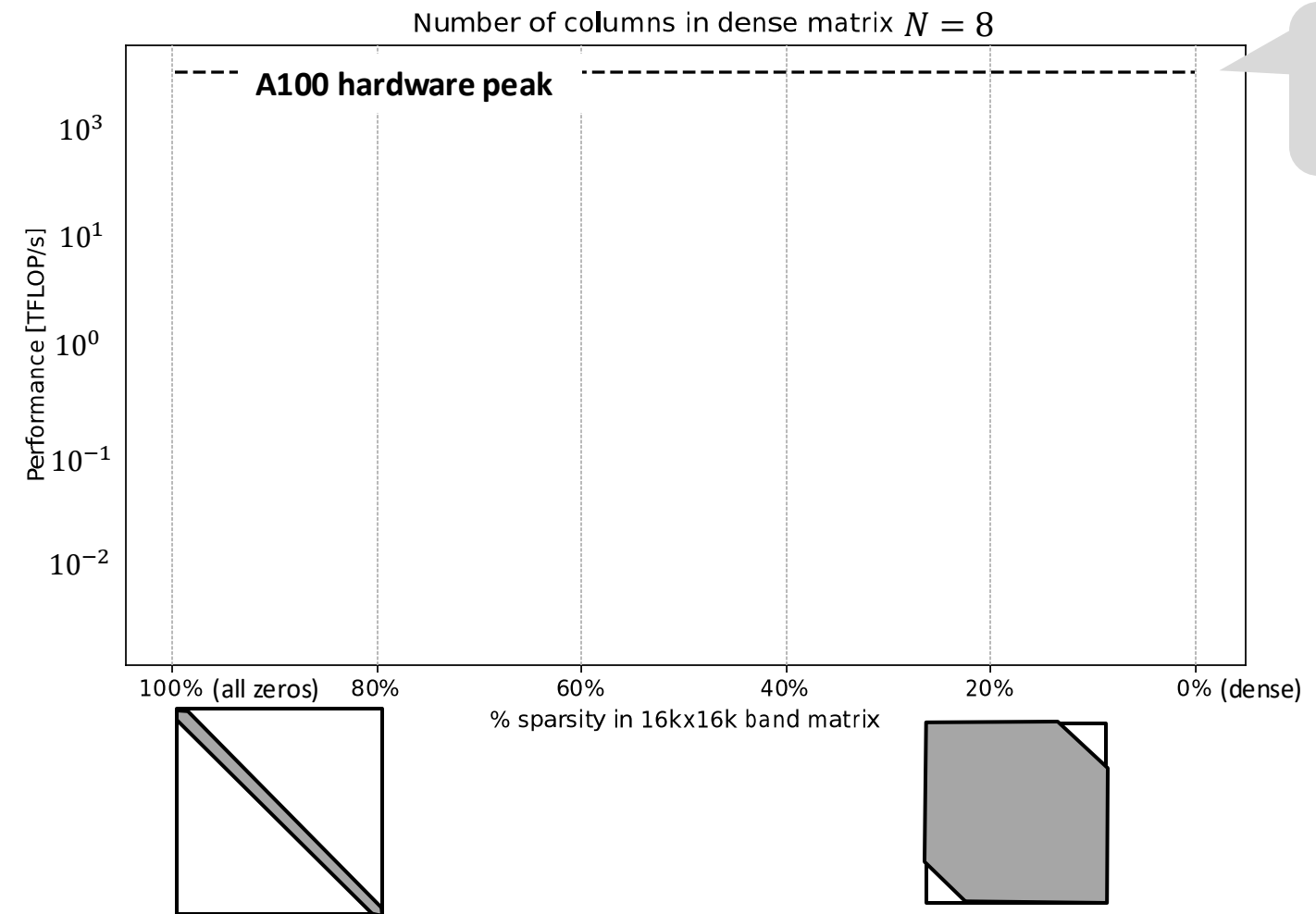


# Synthetic Matrices

16,384 x 16,384



Measure the dependence on the number of blocks and isolate the randomness of the sparse matrix structure

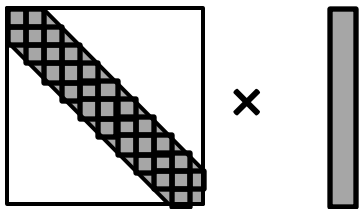


At what sparsity threshold a sparse library can outperform a highly-optimized dense library?

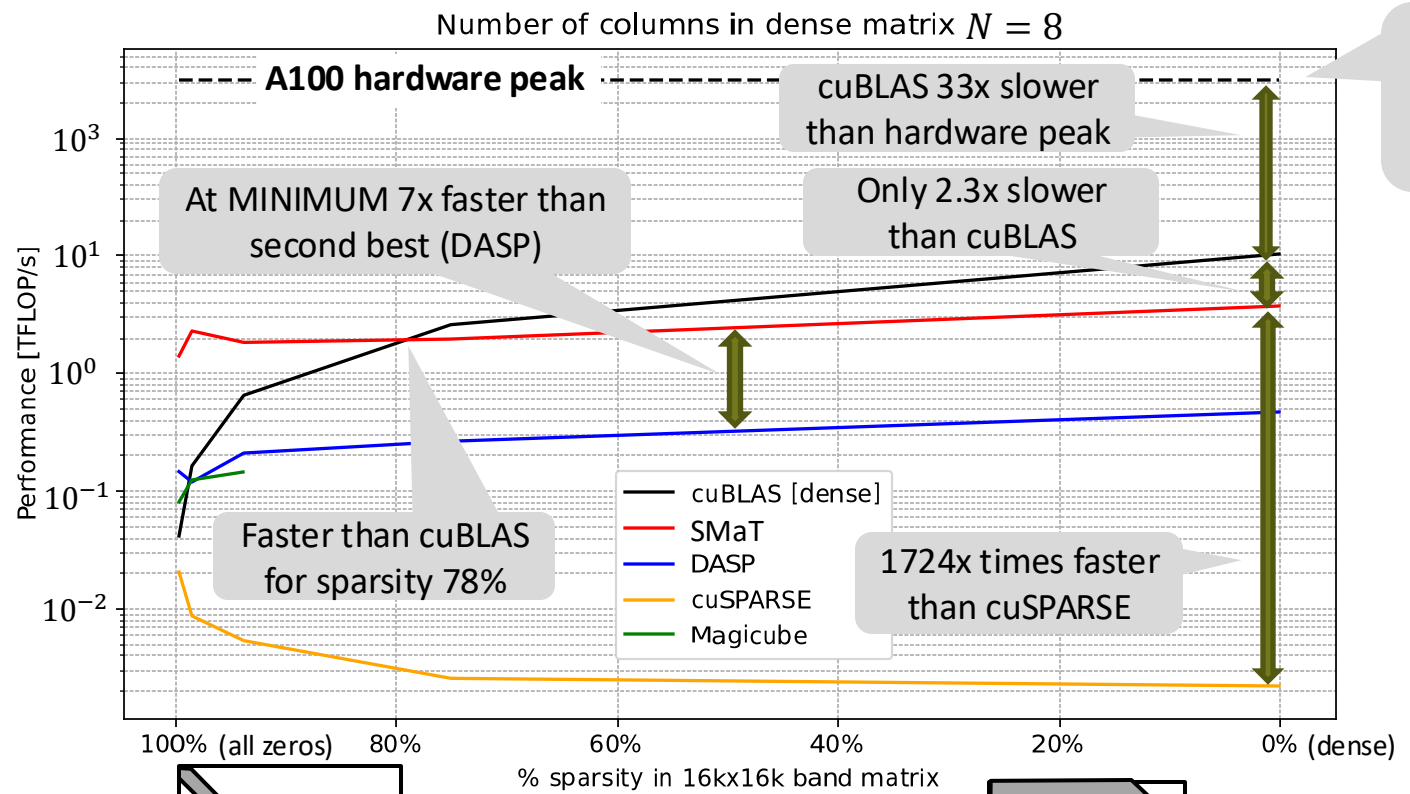


# Synthetic Matrices

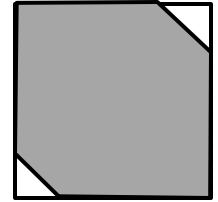
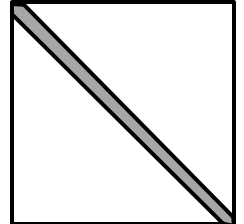
16,384 x 16,384



Measure the dependence on the number of blocks and isolate the randomness of the sparse matrix structure

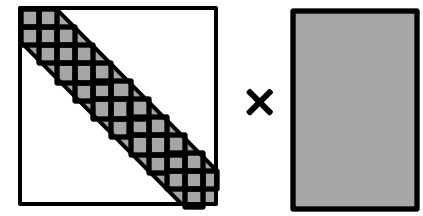


At what sparsity threshold a sparse library can outperform a highly-optimized dense library?

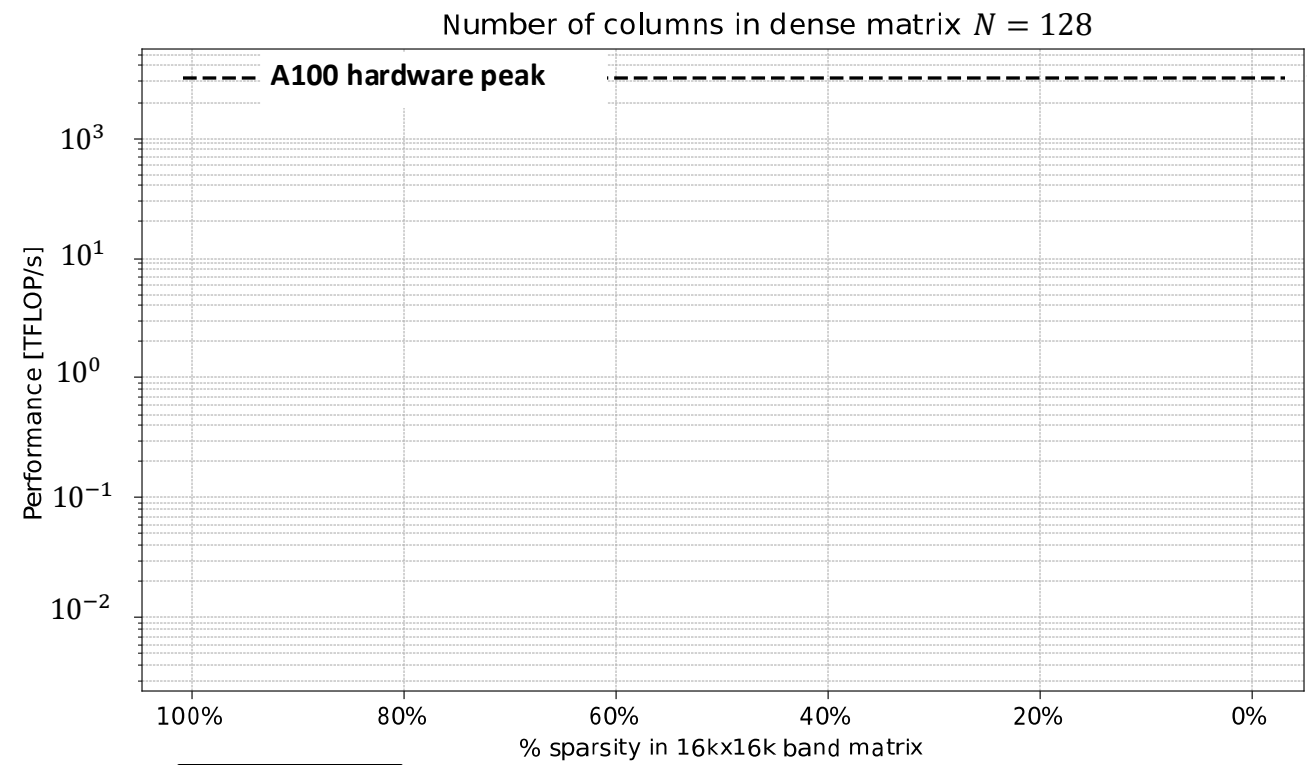


# Synthetic Matrices

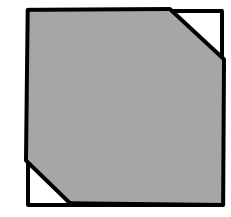
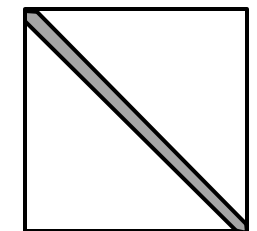
16,384 x 16,384



Measure the dependence on the number of blocks and isolate the randomness of the sparse matrix structure

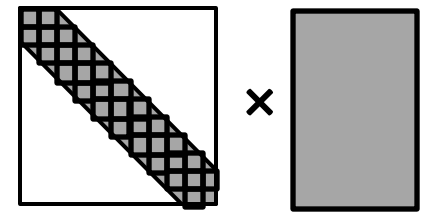


At what sparsity threshold a sparse library can outperform a highly-optimized dense library?

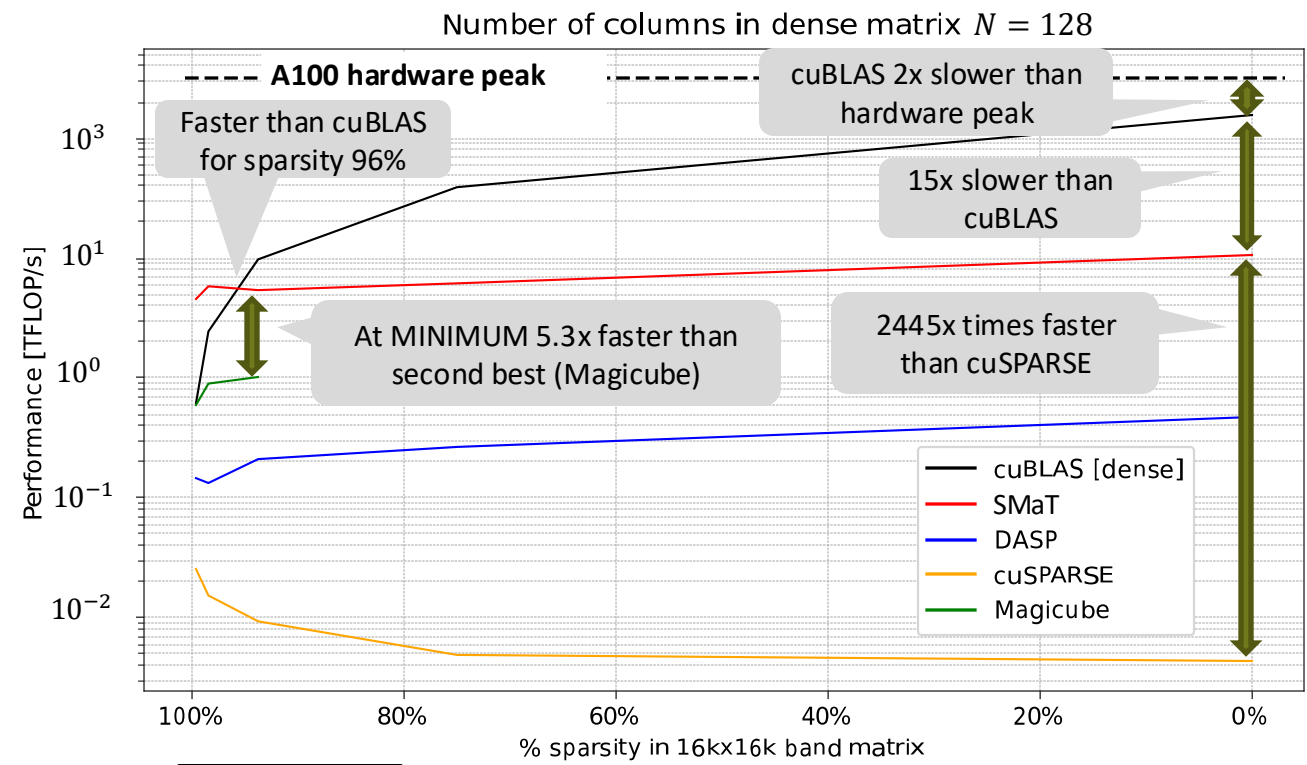


# Synthetic Matrices

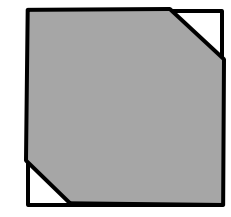
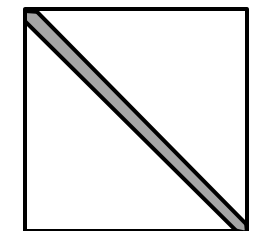
16,384 x 16,384



Measure the dependence on the number of blocks and isolate the randomness of the sparse matrix structure



At what sparsity threshold a sparse library can outperform a highly-optimized dense library?

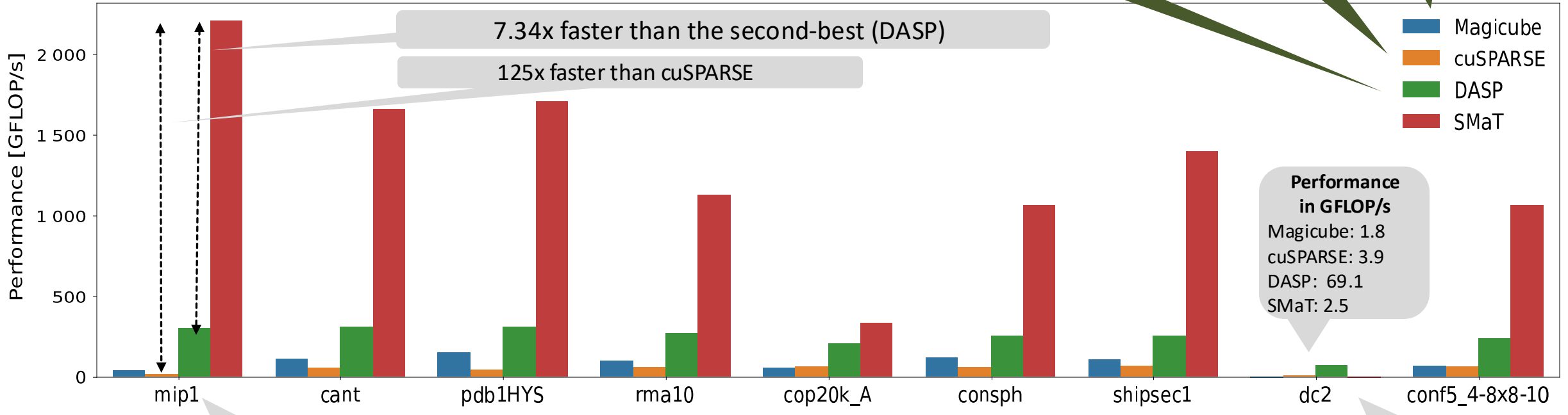


# SuiteSparse performance comparison

2.60x faster (up to 7.34x)

16.32x faster (up to 125.48x)

10.78x faster (up to 51.23x)



7.34x faster than the second-best (DASP)

125x faster than cuSPARSE

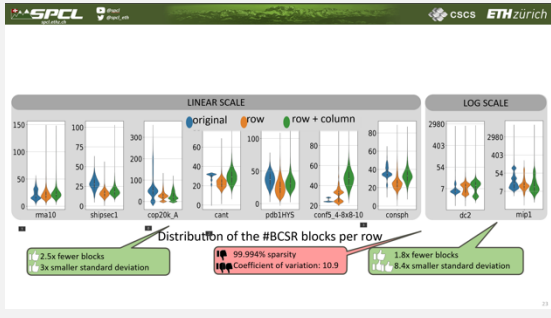
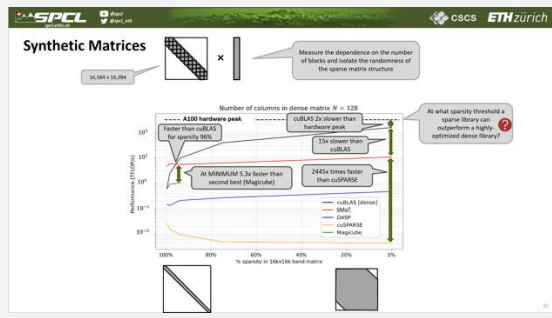
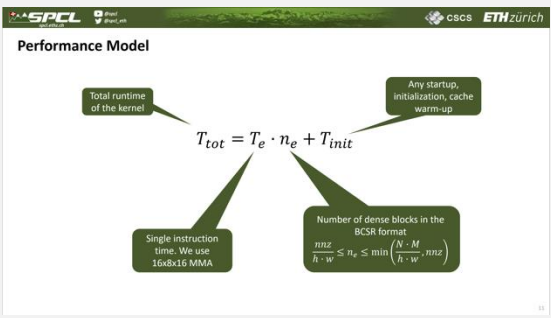
**Performance in GFLOP/s**  
 Magicube: 1.8  
 cuSPARSE: 3.9  
 DASP: 69.1  
 SMaT: 2.5

Importance of good preprocessing

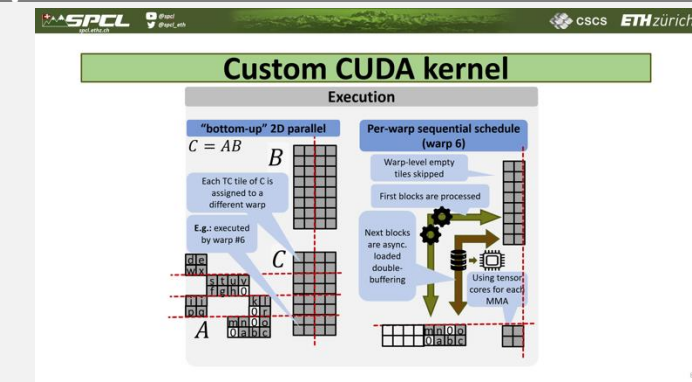
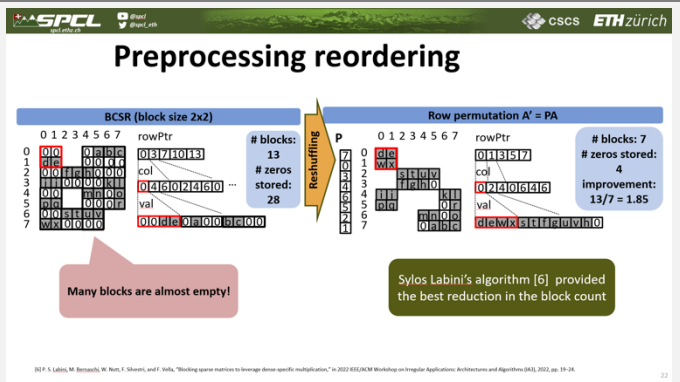
Ill-suited for SMaT's execution model with static 2D parallel schedule

SMaT is better on average **7.71x** (geometric mean) than the respective baselines

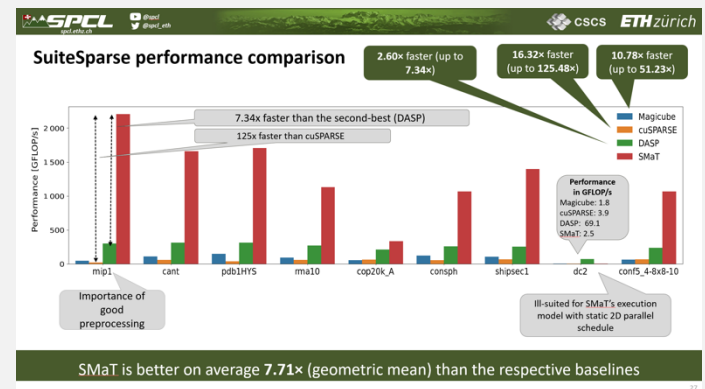
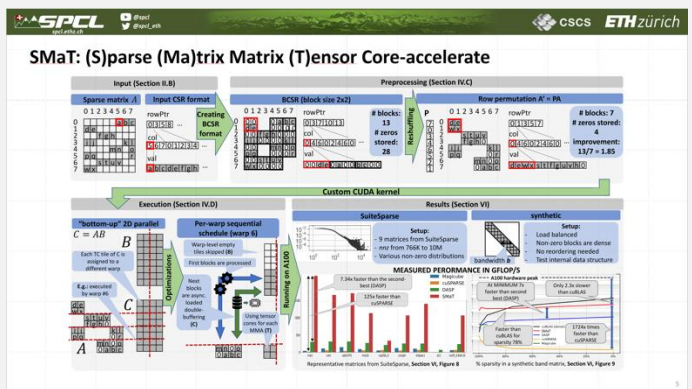
# UNDERSTANDING performance



# OPTIMIZING performance



# DELIVERING performance



More of SPCL's research:

youtube.com/@spl **210+ Talks**

twitter.com/spl\_eth **1.6K+ Followers**

github.com/spl **5.6K+ Stars**

... or [spl.ethz.ch](http://spl.ethz.ch)

