# Reordering sparse matrices
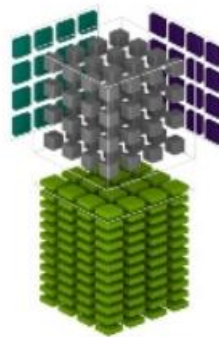# for fast multiplication on tensor architectures

Paolo Sylos Labini
Faculty of Computer Science,
Free University of Bozen-Bolzano

# Motivation



**Tensor  cores** are good at multiplying dense arrays

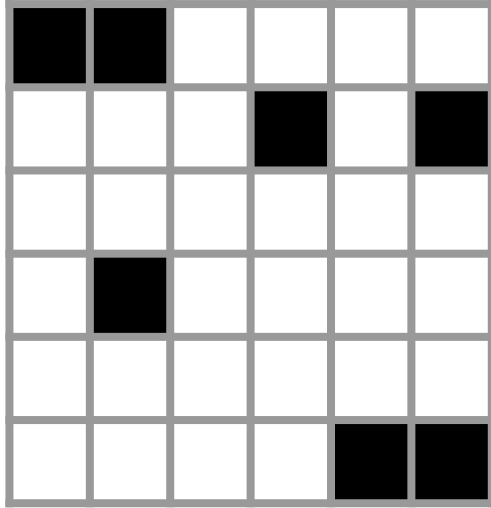But what about **sparse matrices** multiplication**?**

**Should we use sparse-specific data structures (CSR) and algorithms (cusparse)?**

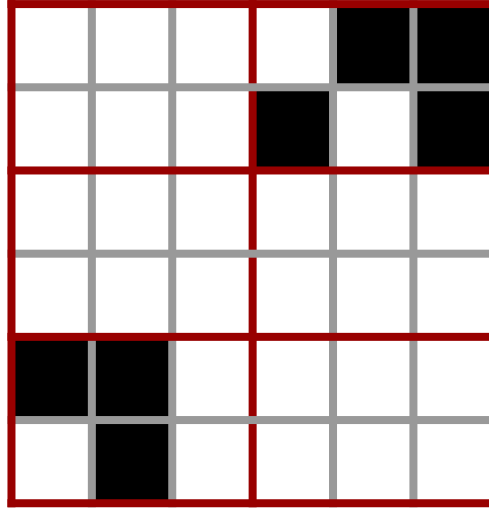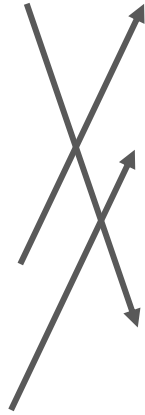**We can't benefit of the tensor architecture**

**Should we treat the sparse matrix as a dense one and use tensor architectures?**

**We end up multiplying a lot of zeros**
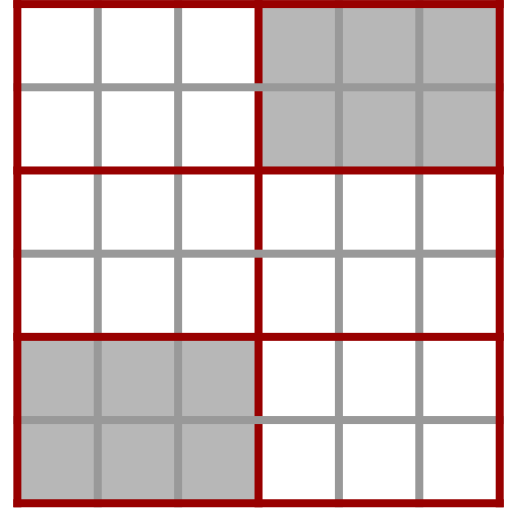
# Our solution



ORIGINAL                    REORDERED                    BLOCKED

# Applications

Reordering takes time

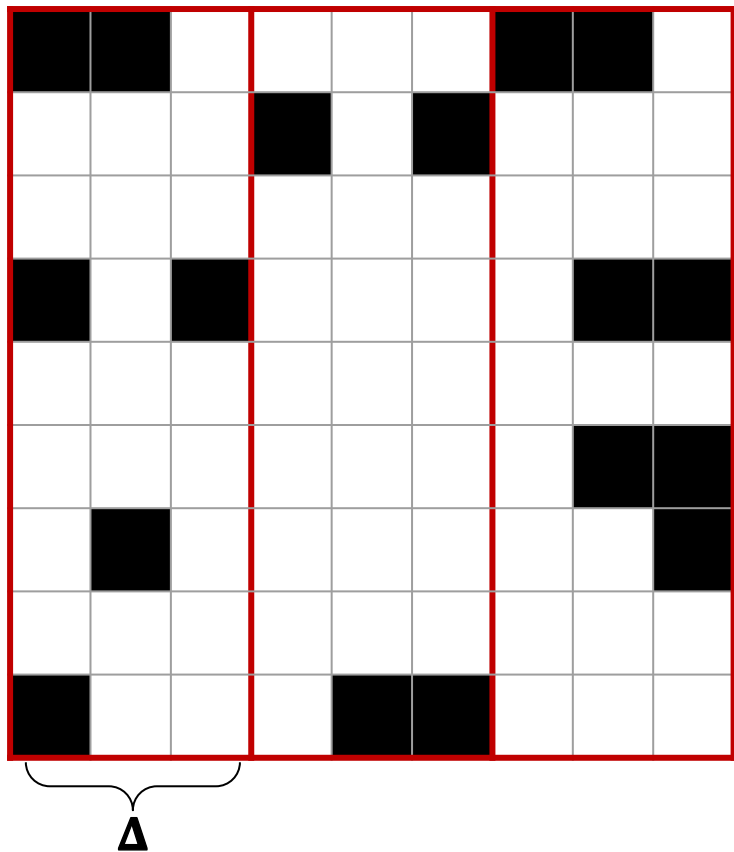If the same matrix is multiplied repeatedly, we can absorb the cost

- graphs
- neural networks
- matrix collections

# The reordering algorithm
Find similar rows and group them together in blocks

1) Partitioning the columns
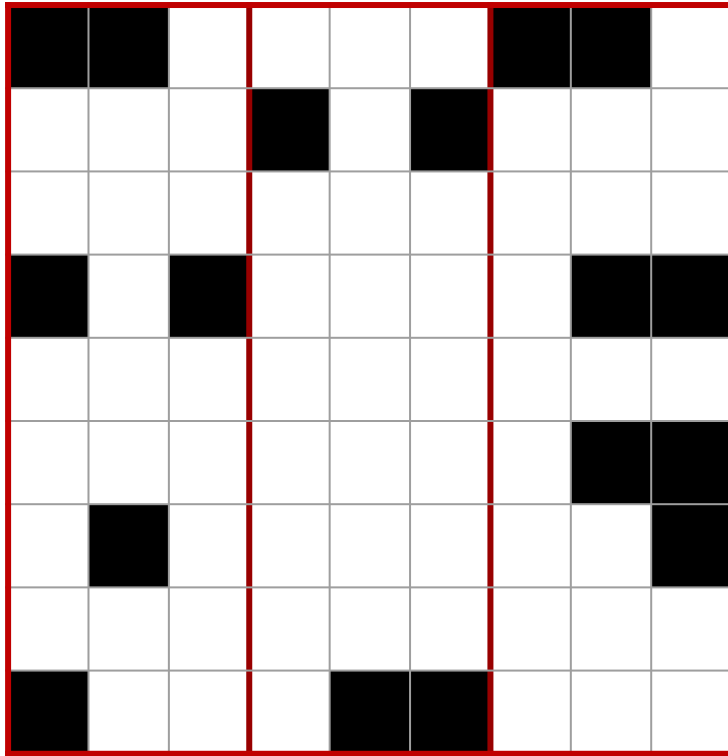
2) Grouping by hash

3) Grouping by similarity

# The reordering algorithm – partitioning the columns



Choice of **Δ** may depend on
- architecture
- previous knowledge of the matrix
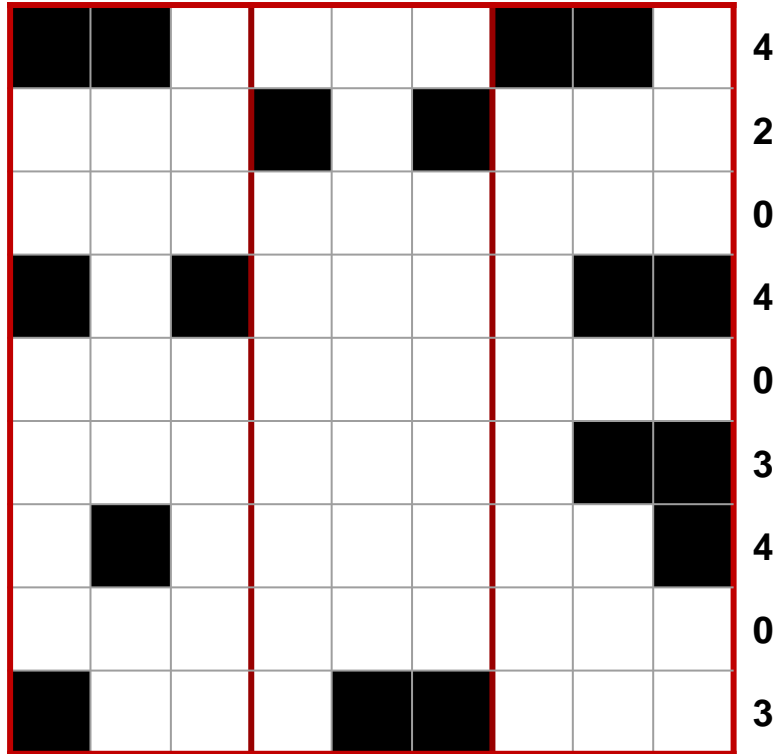- ...

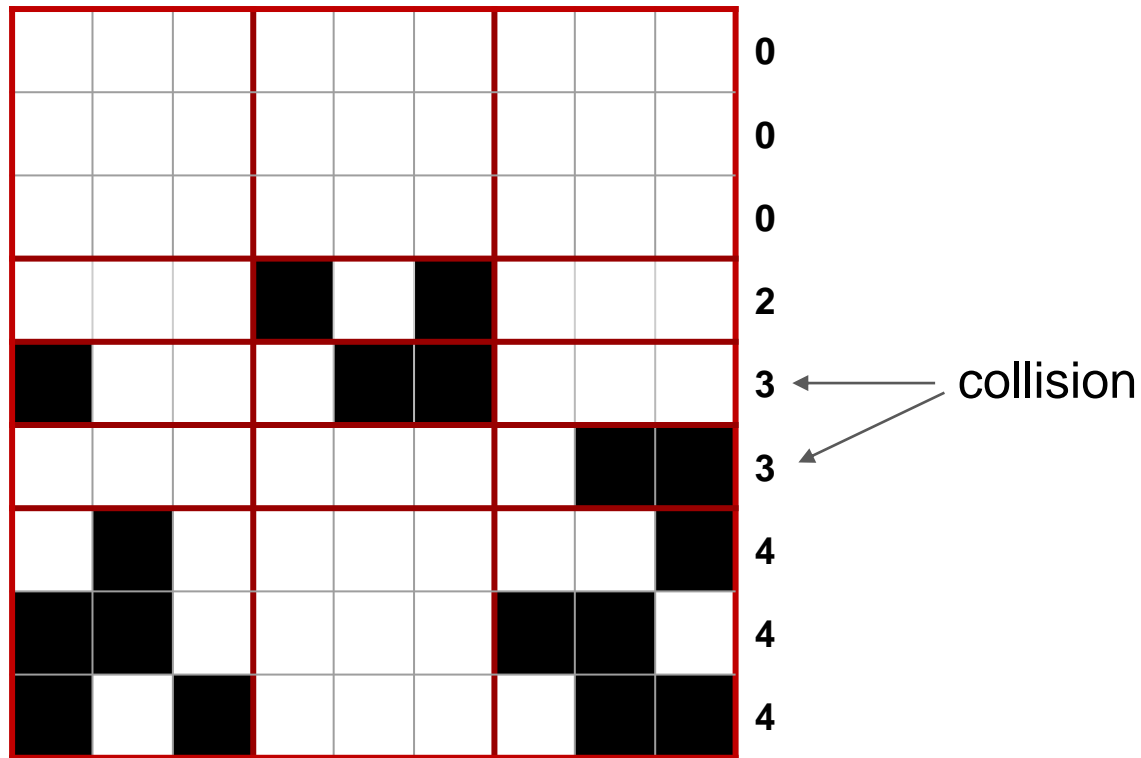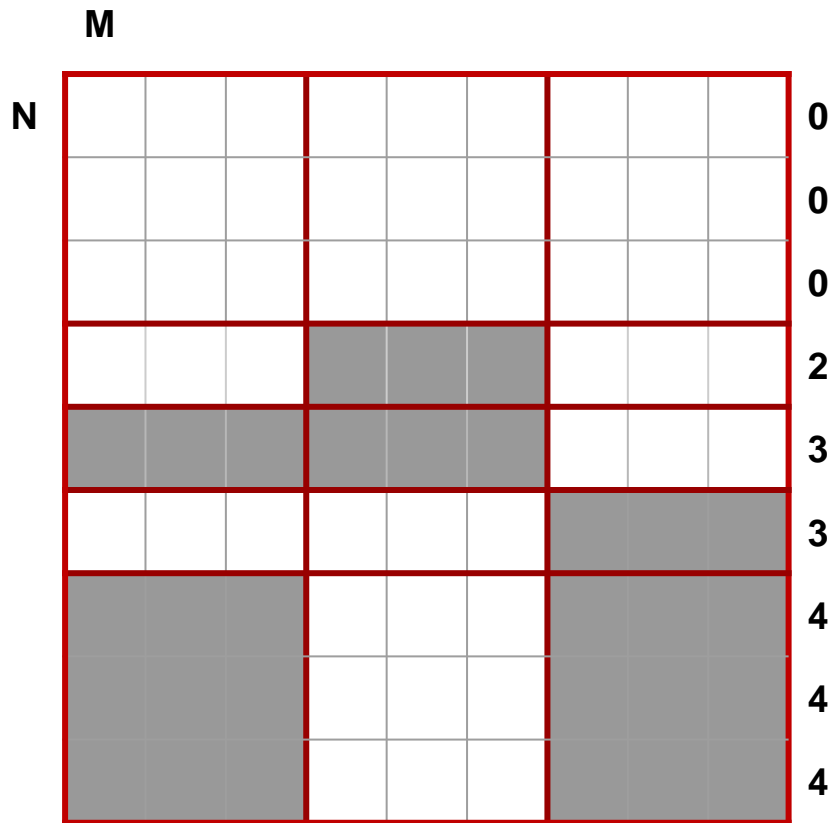# The reordering algorithm - hashing the rows



1          3

**h = 1 + 3 = 4**

# The reordering algorithm - hashing the rows

# The reordering algorithm – grouping by hash

# The reordering algorithm - grouping by hash

**M**

**N**

| | | | | |
|---|---|---|---|---|
| | | | | 0 |
| | | | | 0 |
| | | | | 0 |
| | | | | 2 |
| | | | | 3 |
| | | | | 3 |
| | | | | 4 |
| | | | | 4 |
| | | | | 4 |

1) **hashing**: **O(NxM)**
   for CSR matrix with K entries: O(K)

1) **sorting:  O(N log N)**

1) **grouping:** approx. NxM

# The reordering algorithm - grouping by cosine



**For each** pair of groups A, B

**if**    cos **θ**(A,B)  ≥  ε

**merge** A and B

A

| 0 | 0 | 1 |
|---|---|---|

B

| 3 | 0 | 3 |
|---|---|---|

M

N

# The reordering algorithm - grouping by cosine

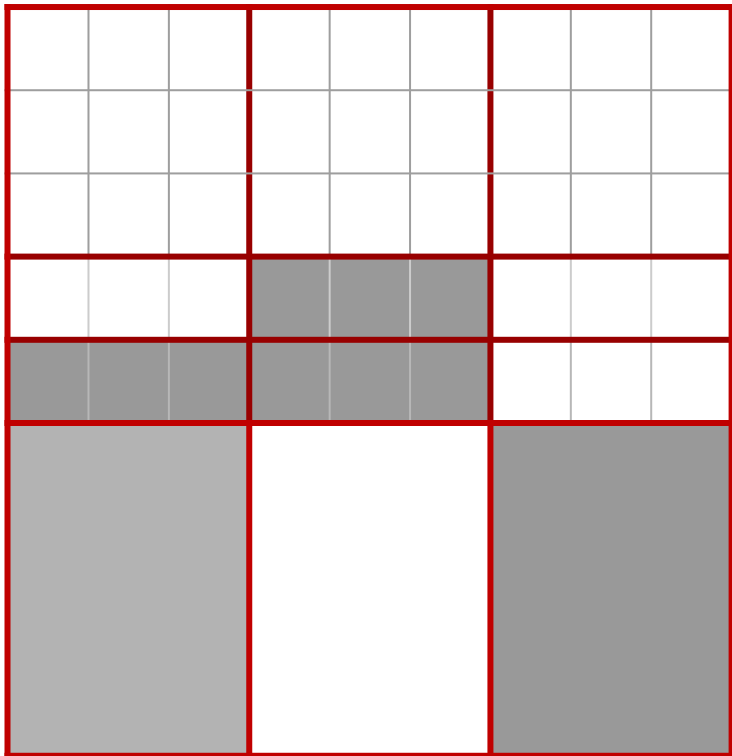**M**

**N**

**For each** pair of groups A, B

**if**   cos **θ**(A,B)  ≥  ε

**merge** A and B

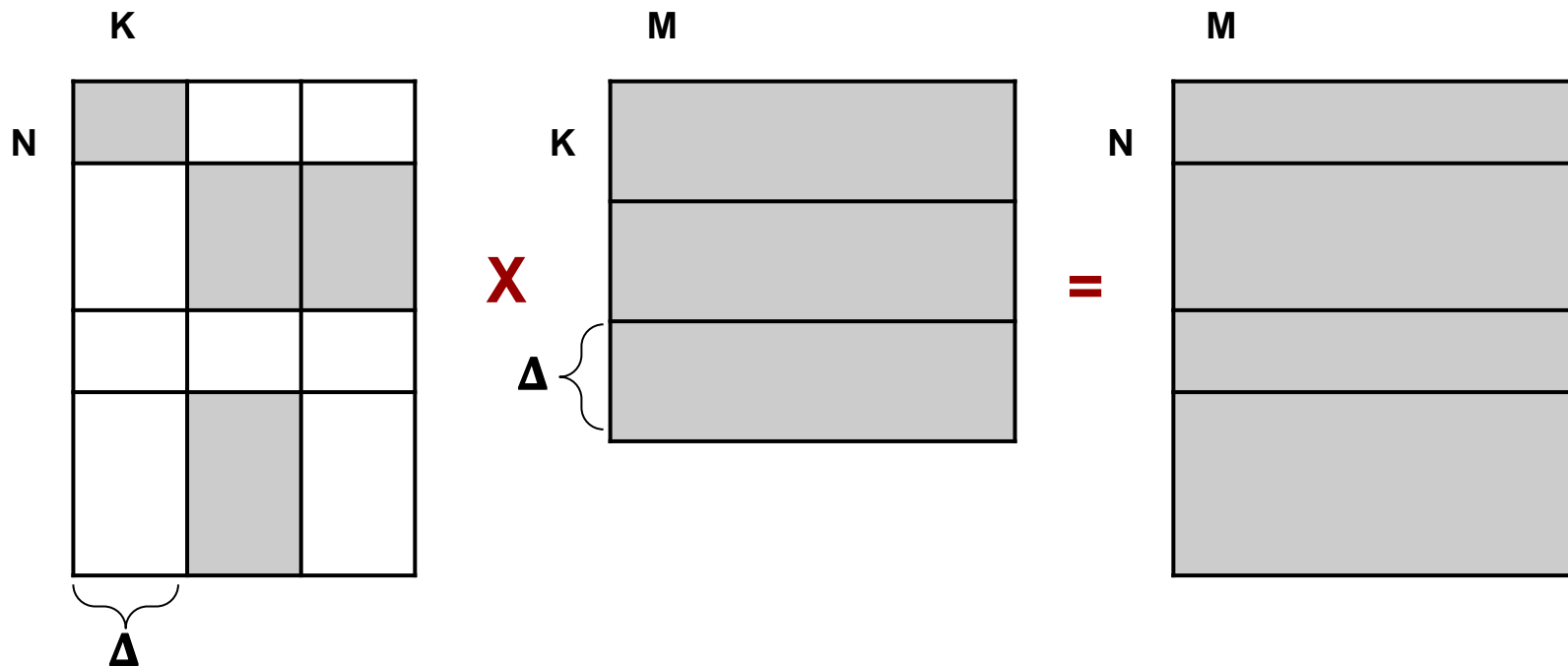# The reordering algorithm - grouping by cosine

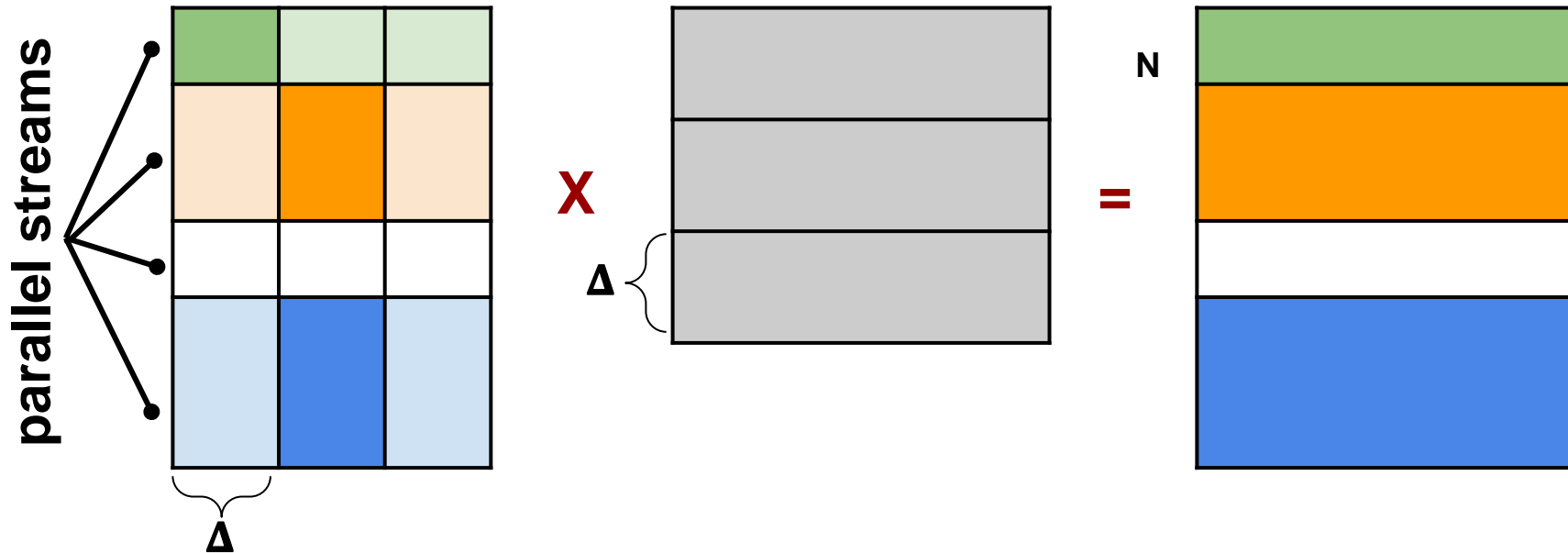**M**

**N**



**O(N*N*M)**

But with some optimization:
- Only groups are compared, not rows
- Only patterns are compared, not entries

Still, too expensive!

# Parallel multiplication
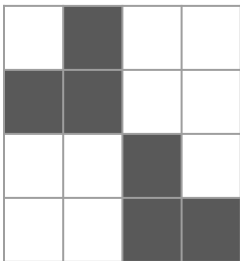
# Parallel multiplication

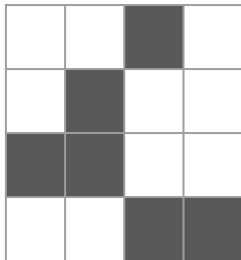# Reorder - experiments with synthetic matrices
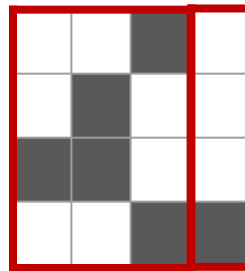
Synthetic block matrix          Scrambled          Partitioned

**Δ = 2**

**μ = 75%**

**η = 50%**

**Δ' = 3**

**Δ = original block size**
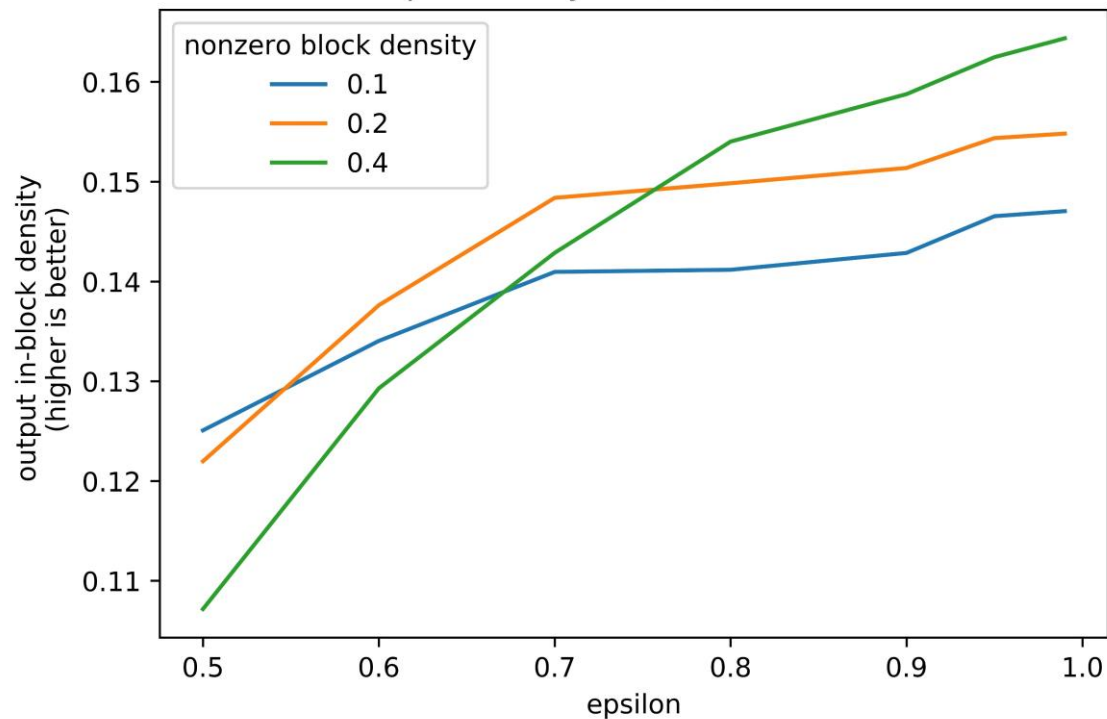
**μ = original in-block density**

**η = original nonzero block density**

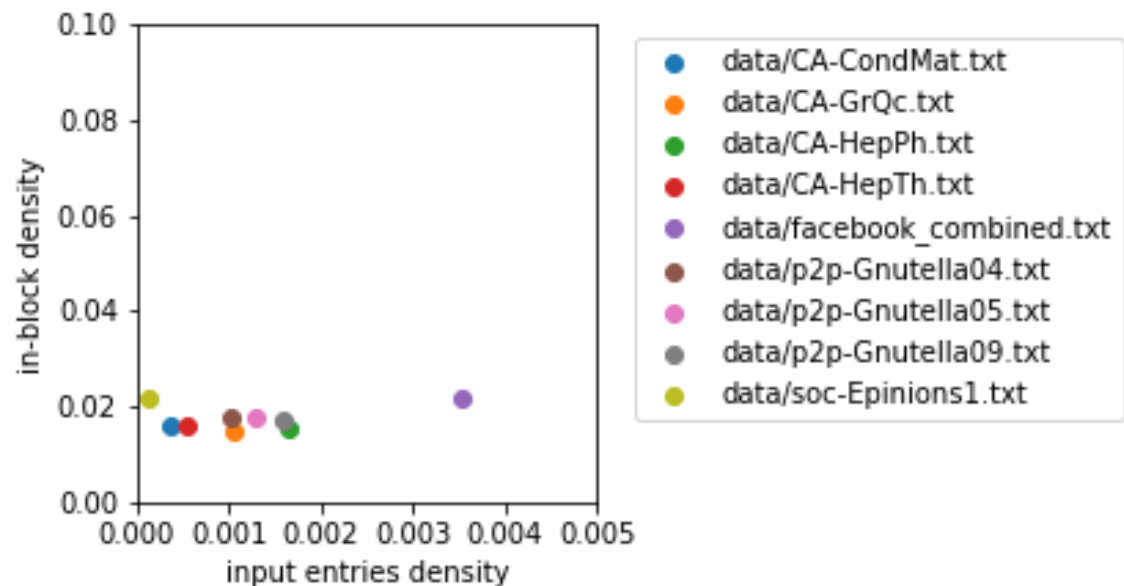# Reorder - experiments with synthetic matrices



algorithm block size = 35
input block size = 64
input density inside blocks = 0.2
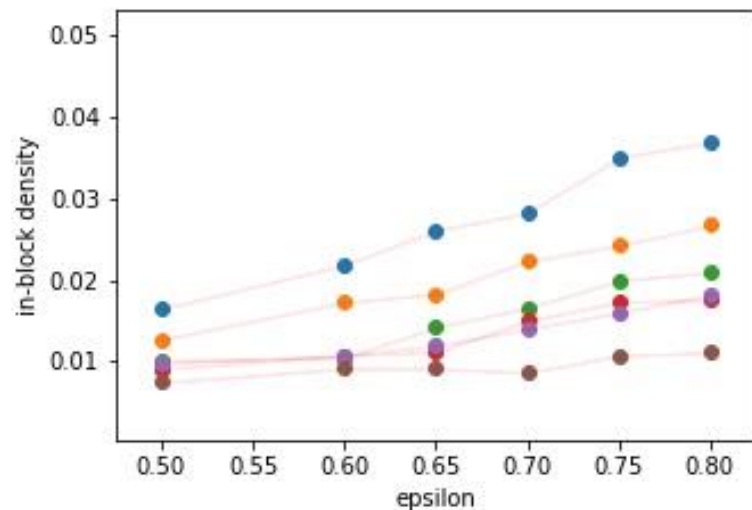
# Reorder - experiments with real matrices



algo_block_size:64
scramble:1
epsilon:0.6

Legend:
- data/CA-CondMat.txt
- data/CA-GrQc.txt
- data/CA-HepPh.txt
- data/CA-HepTh.txt
- data/facebook_combined.txt
- data/p2p-Gnutella04.txt
- data/p2p-Gnutella05.txt
- data/p2p-Gnutella09.txt
- data/soc-Epinions1.txt

x-axis: input entries density
y-axis: in-block density

# Reorder - experiments with real matrices



scramble:1
input_source:data/facebook_combined.txt

scramble:1
input_source:data/facebook_combined.txt

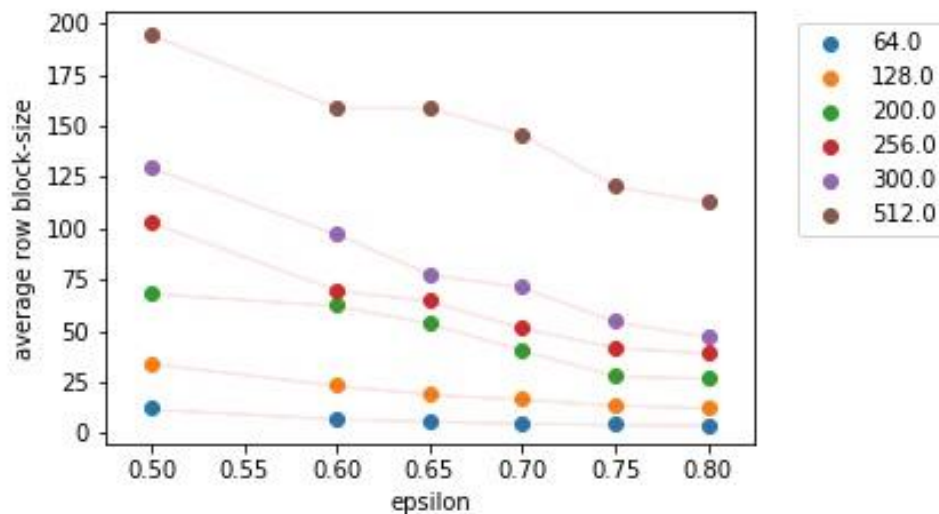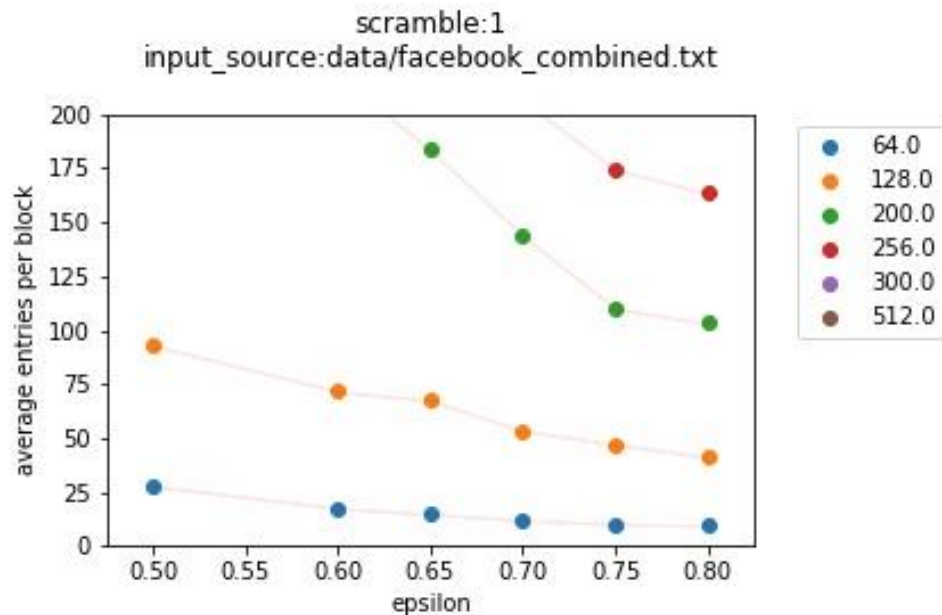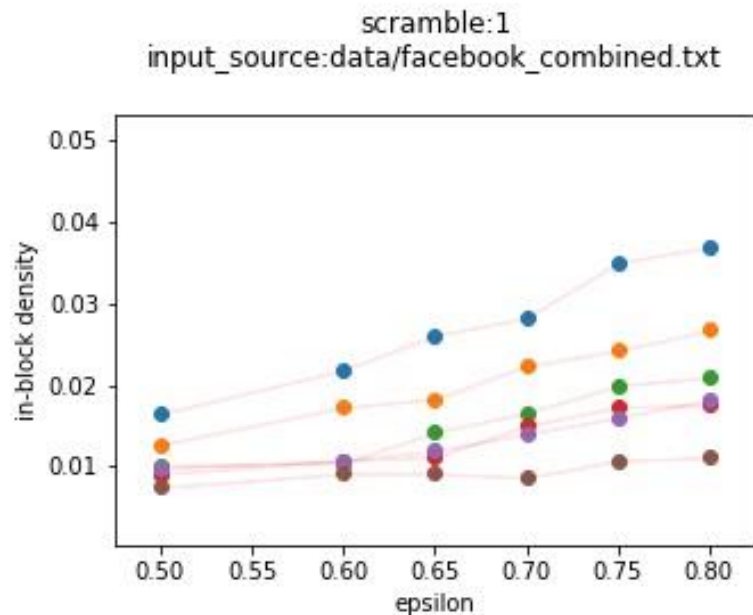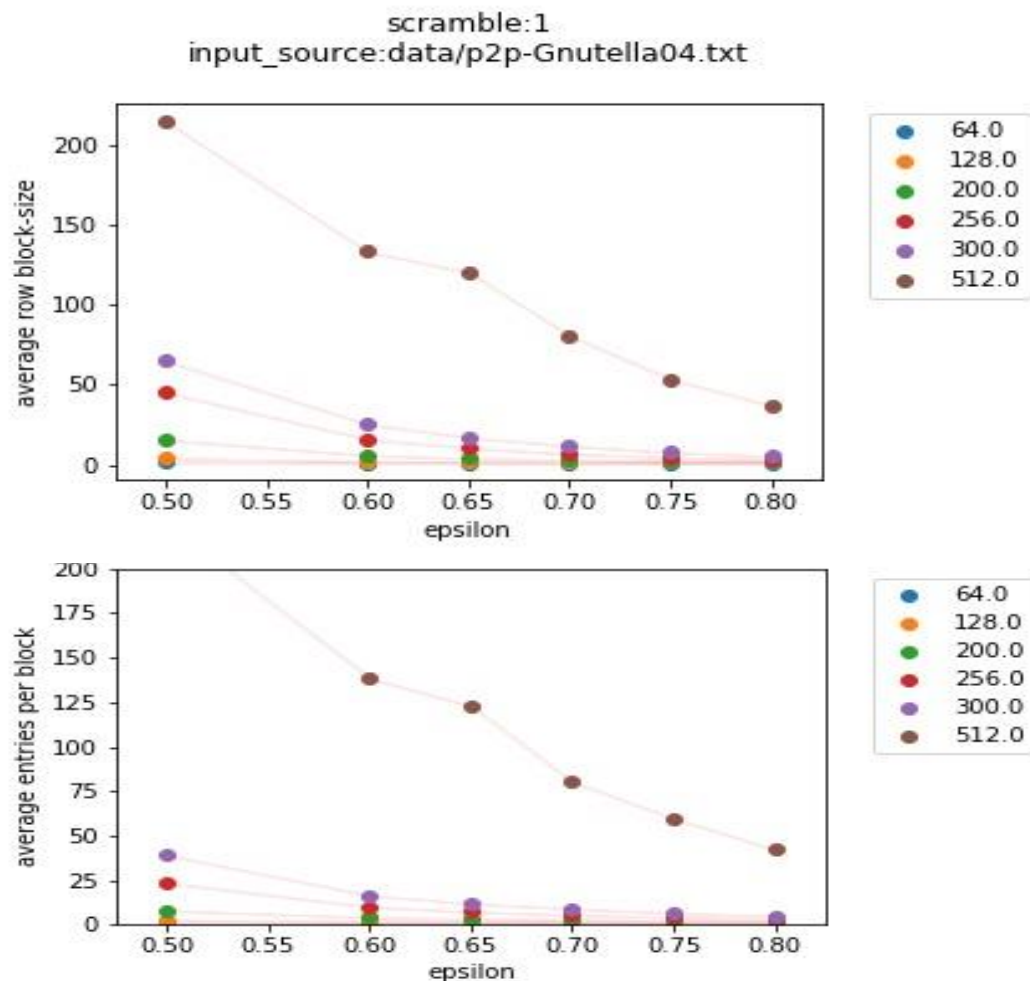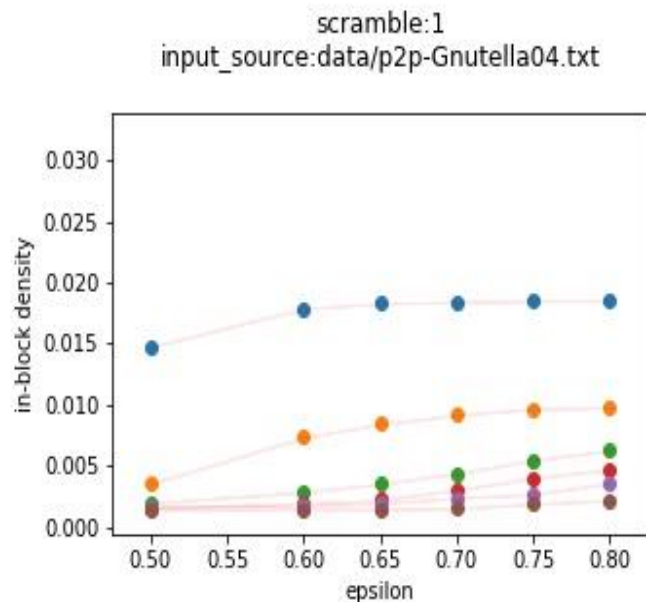# Reorder - experiments with real matrices



scramble:1
input_source:data/facebook_combined.txt

scramble:1
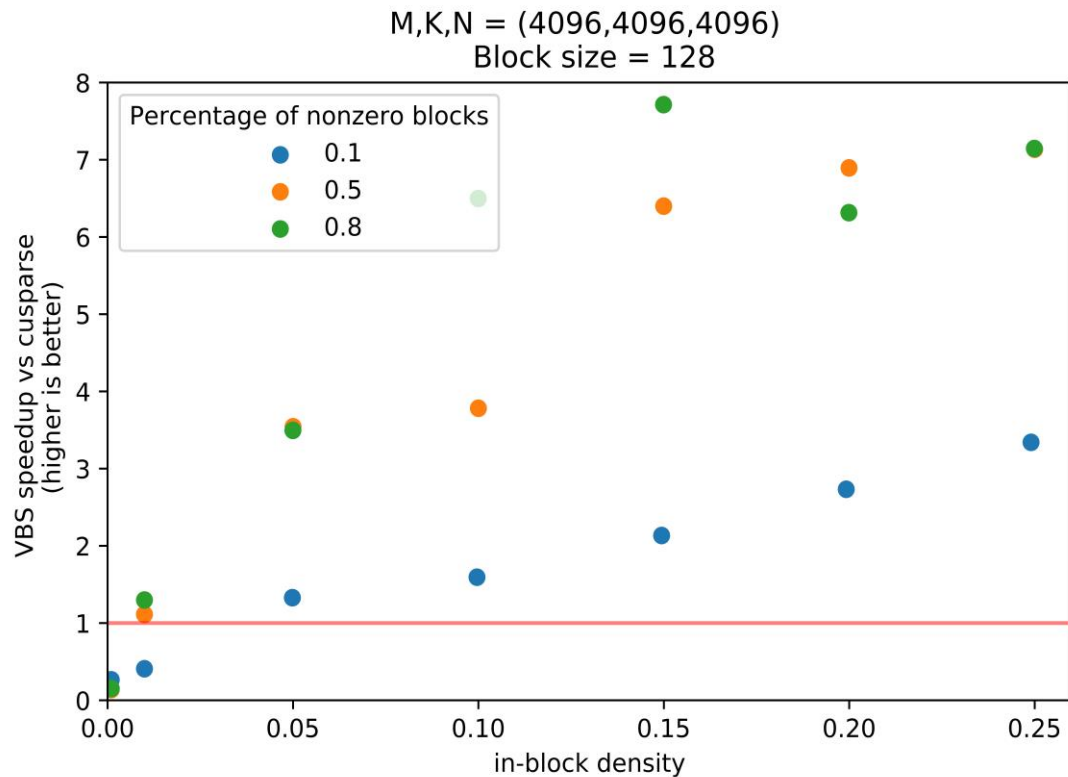input_source:data/facebook_combined.txt

# Reorder



scramble:1
input_source:data/p2p-Gnutella04.txt



scramble:1
input_source:data/p2p-Gnutella04.txt

# Multiplication - experiments

M,K,N = (4096,4096,4096)
Block size = 128

# Multiplication - experiments



M,K,N = (4096,4096,16384)
Block size = 128

# Multiplication - experiments



M,K,N = (4096,2048,16384)
Block size = 128

# Multiplication - experiments

# Future directions – Better, faster reordering

- Parallel implementation
- Local Sensitive Hashing

# Future directions – LSH

**Faster reordering algorithm through local sensitive hashing and nearest neighbours search**
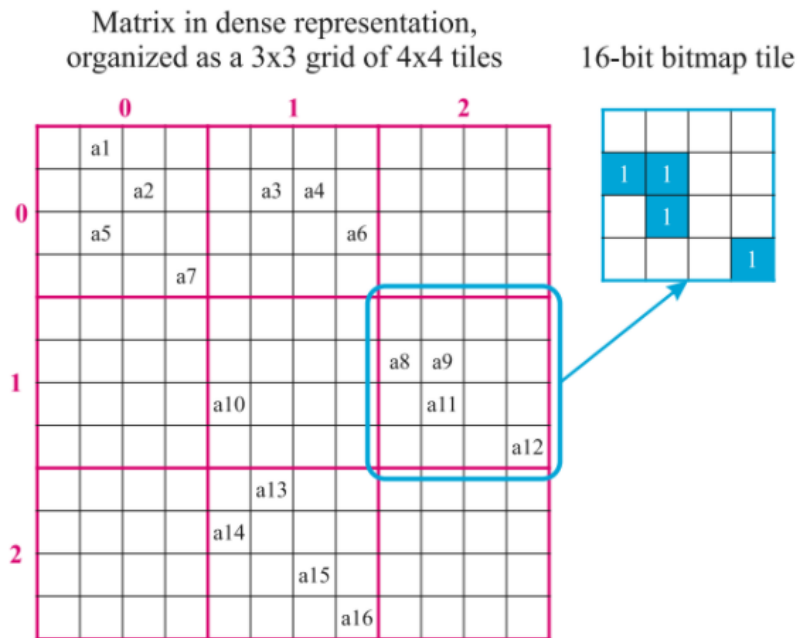
MinHash (random permutations): approximate Jaccard similarity

SimHash (random projections): approximate cosine similarity

# Future directions - Compare to related work

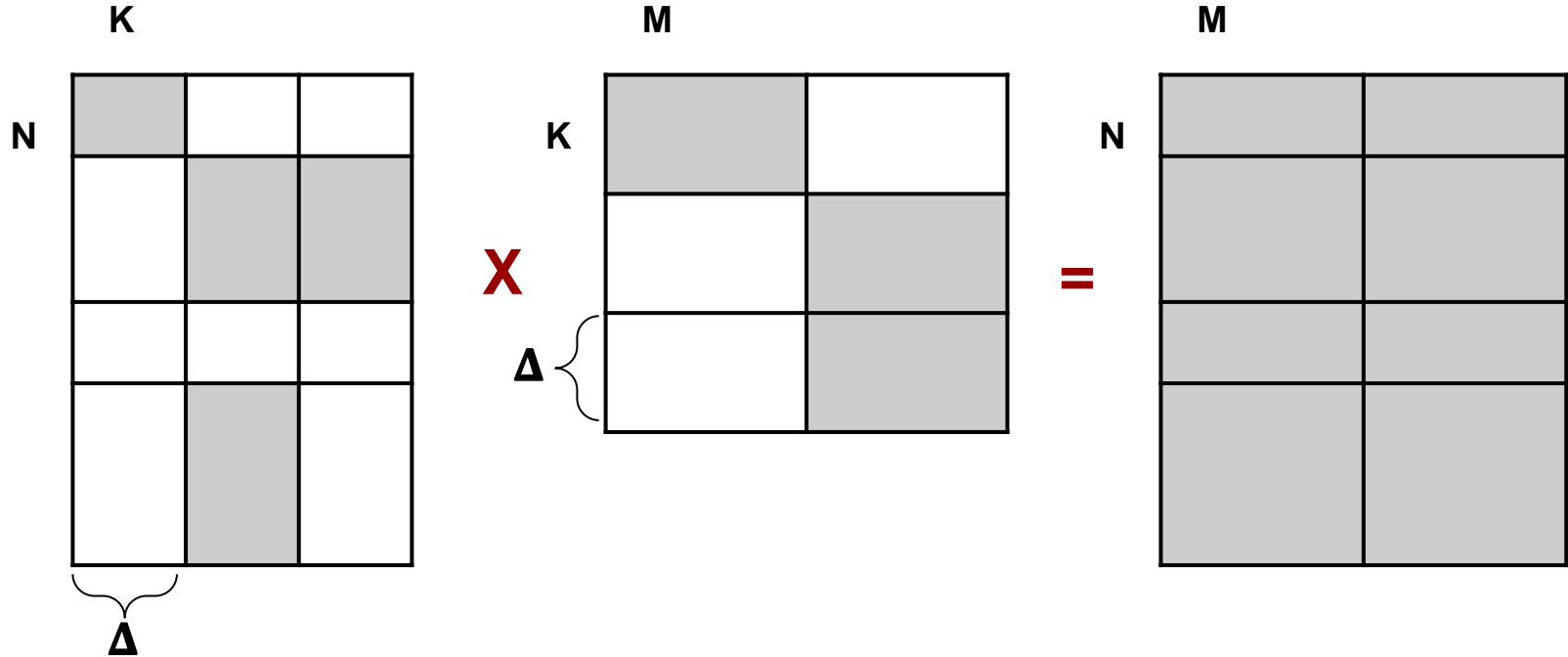Accelerating sparse matrix–matrix multiplication with GPU Tensor Cores[☆]

Orestis Zachariadis [a,*], Nitin Satpute [a], Juan Gómez-Luna [b], Joaquíı

Matrix in dense representation, organized as a 3x3 grid of 4x4 tiles

16-bit bitmap tile

# Future directions - Tensor-core implementation

- Optimal size of block partition?
- Best way to parallelize?
- Focused experiments

# Future directions - Sparse-Sparse multiplication



Problem: zeros gets multiplied by each other
Number of unnecessary operations is squared