

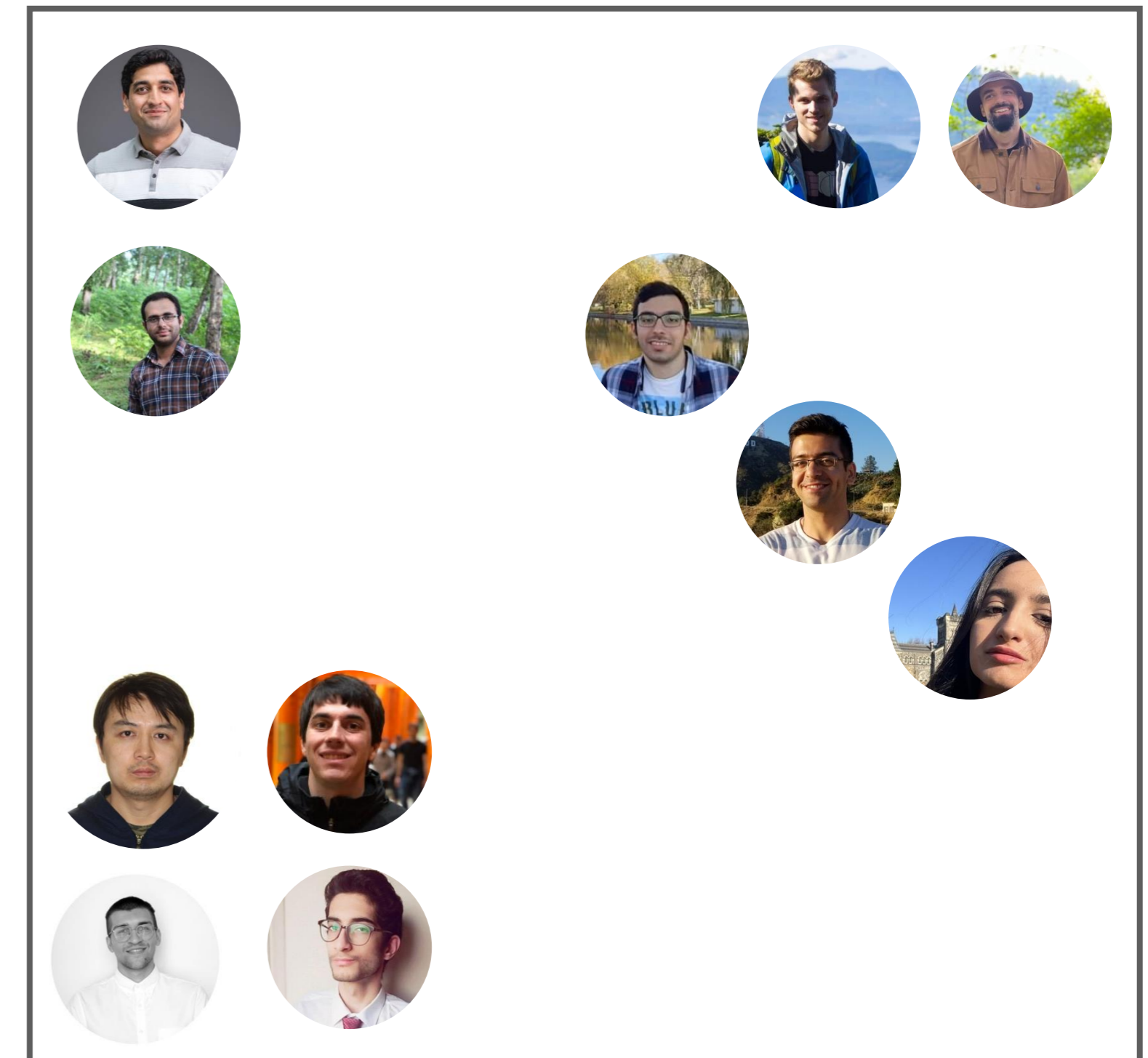
The Power of Less: Harnessing Sparsity for Performance Optimization

Maryam Mehri Dehnavi

University of Toronto

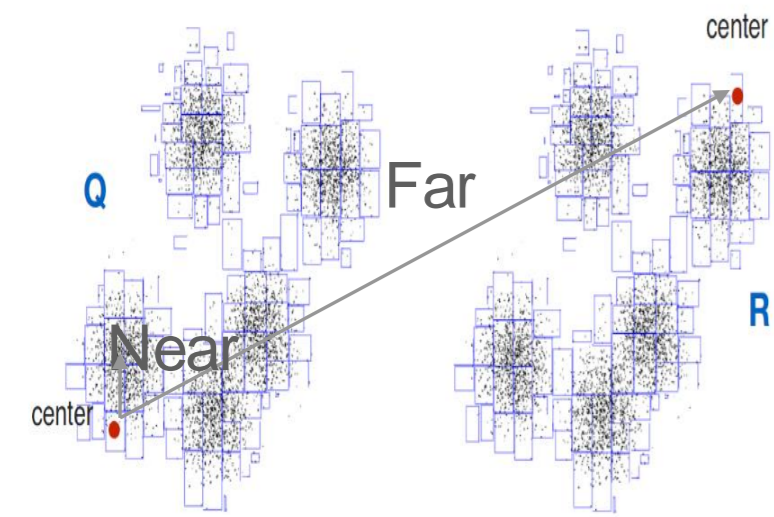
Microsoft Research

<http://www.paramathic.com>

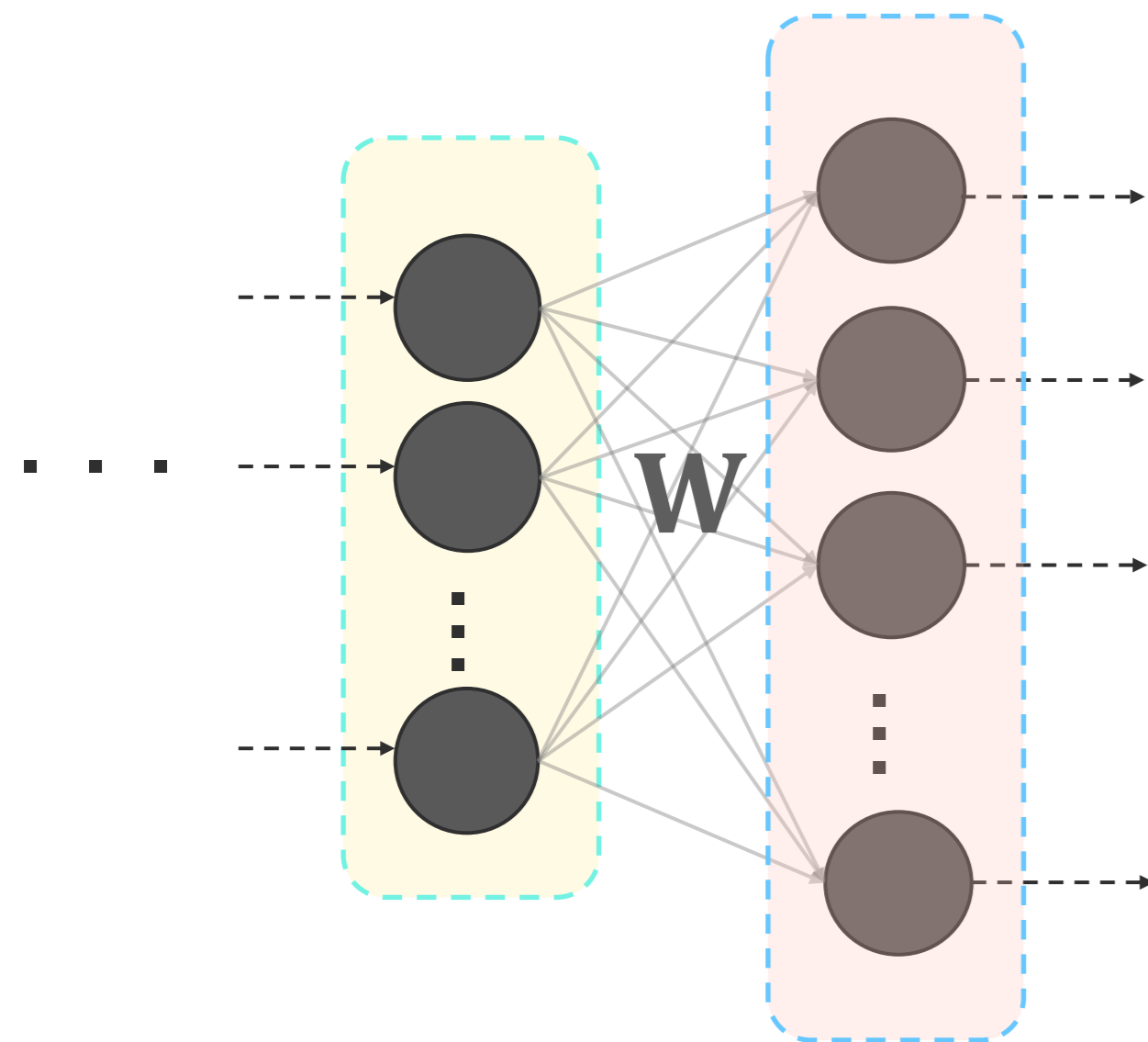


Sparsity is Everywhere

Machine Learning

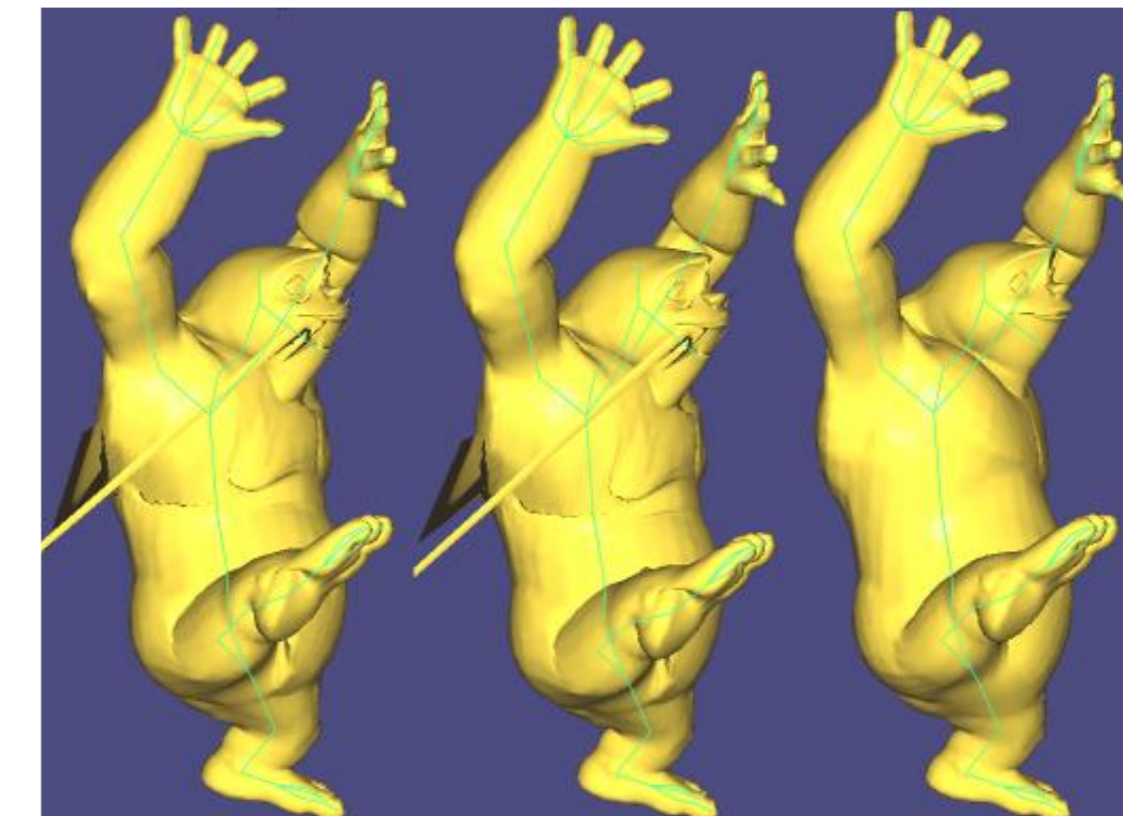


Clustering

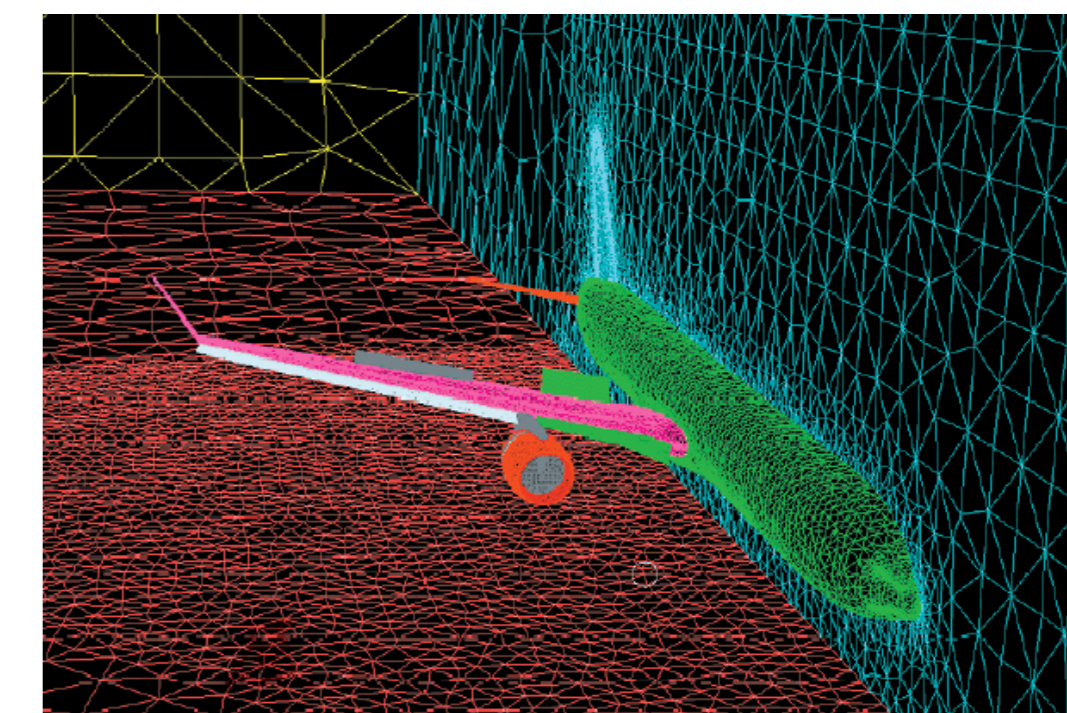
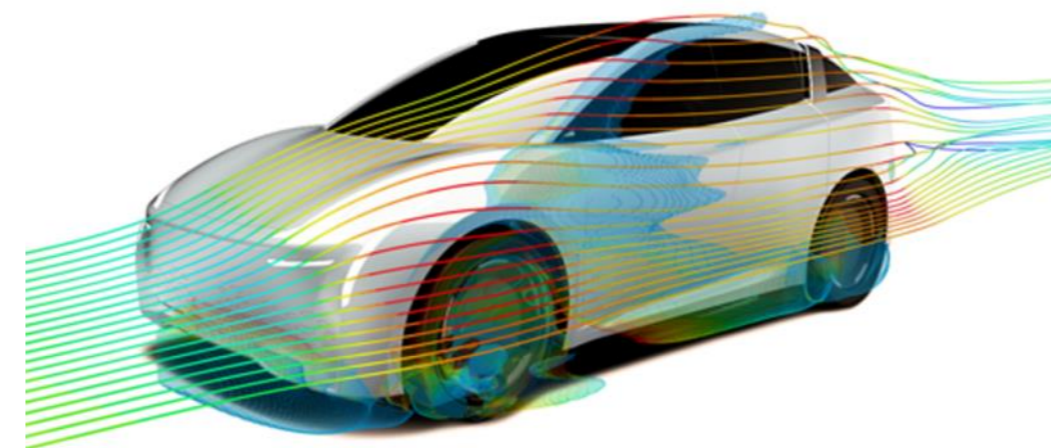


Deep Neural Networks

Science and Engineering



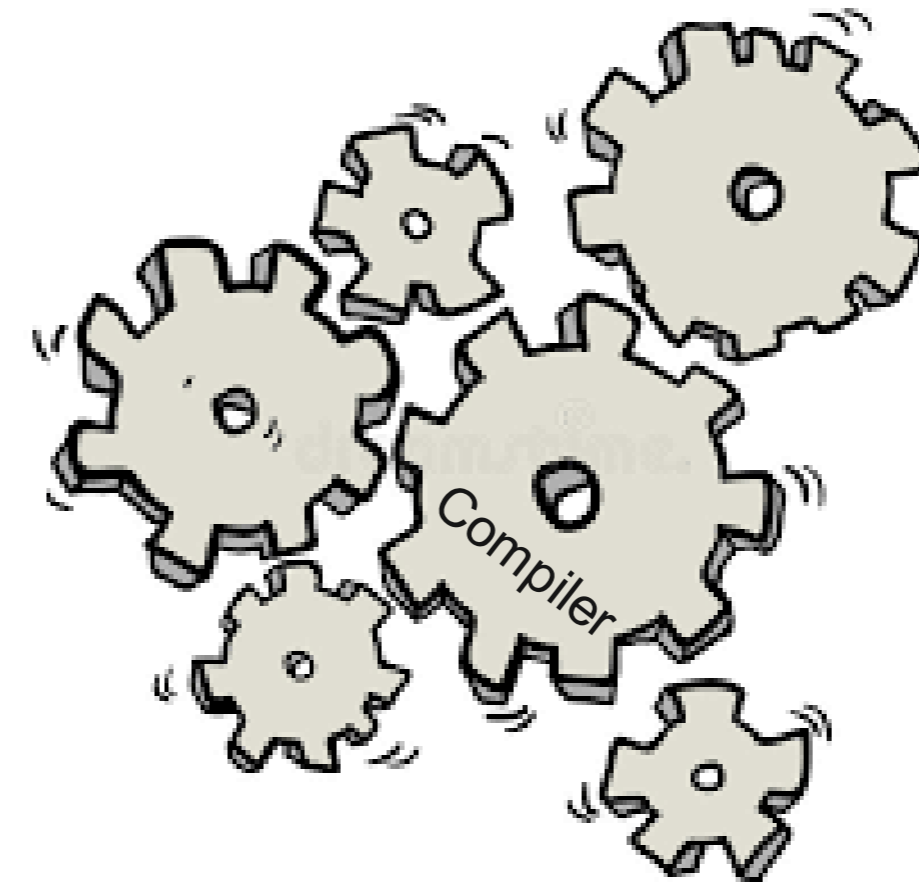
Graphics Simulations



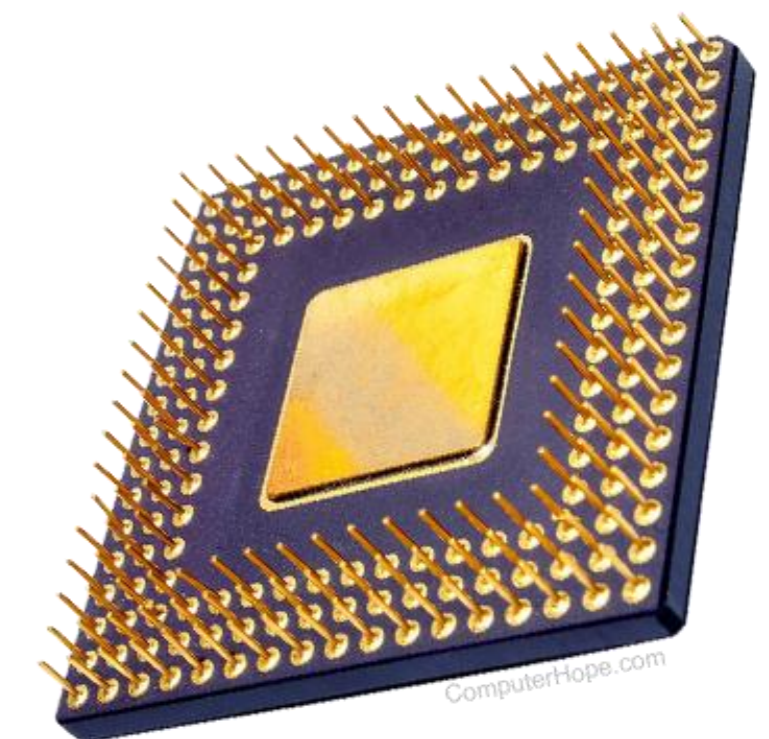
Fluid Dynamics

The World is Built for Dense!

- Hardware utilization optimizations: prefetching, branch predictors, registers, caches, ...
- Compiler optimizations: tiling, unrolling, vectorization, parallelization, ...
- Libraries: BLAS, LINPACK, LAPACK

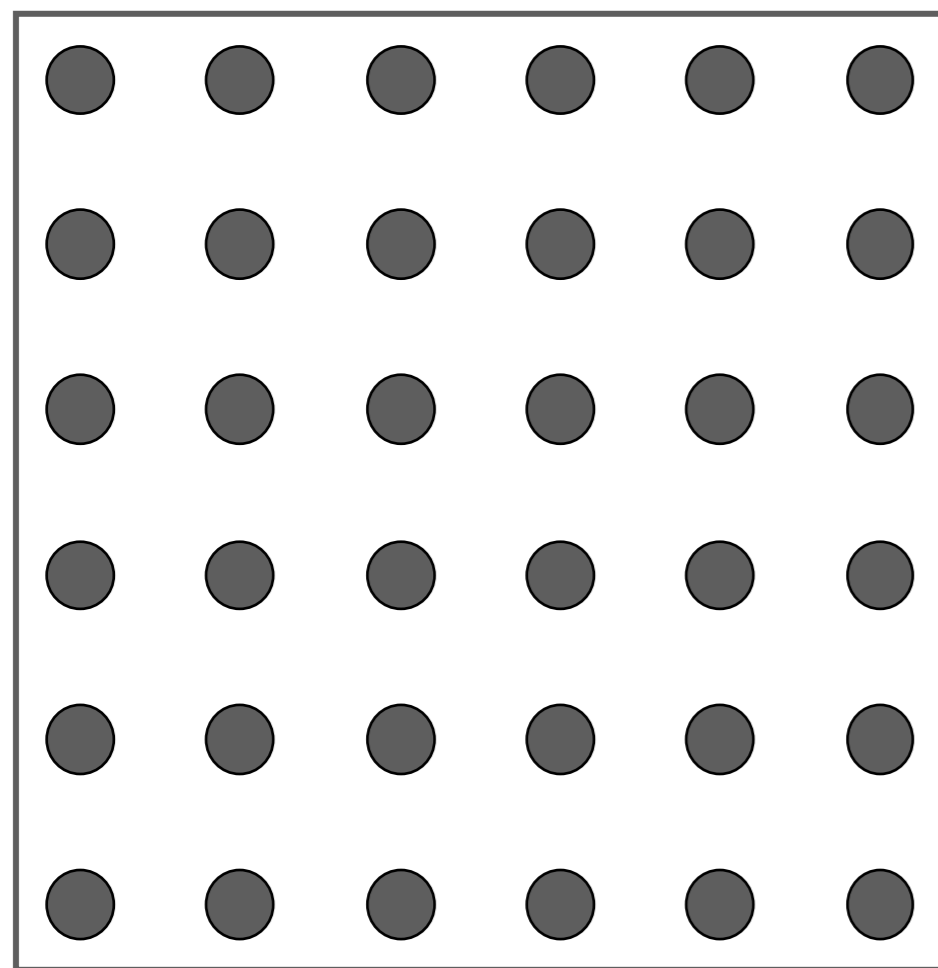


L A P A C K
L -A P -A C -K
L A P A -C -K
L -A P -A -C K
L A -P -A C K
L -A -P A C -K

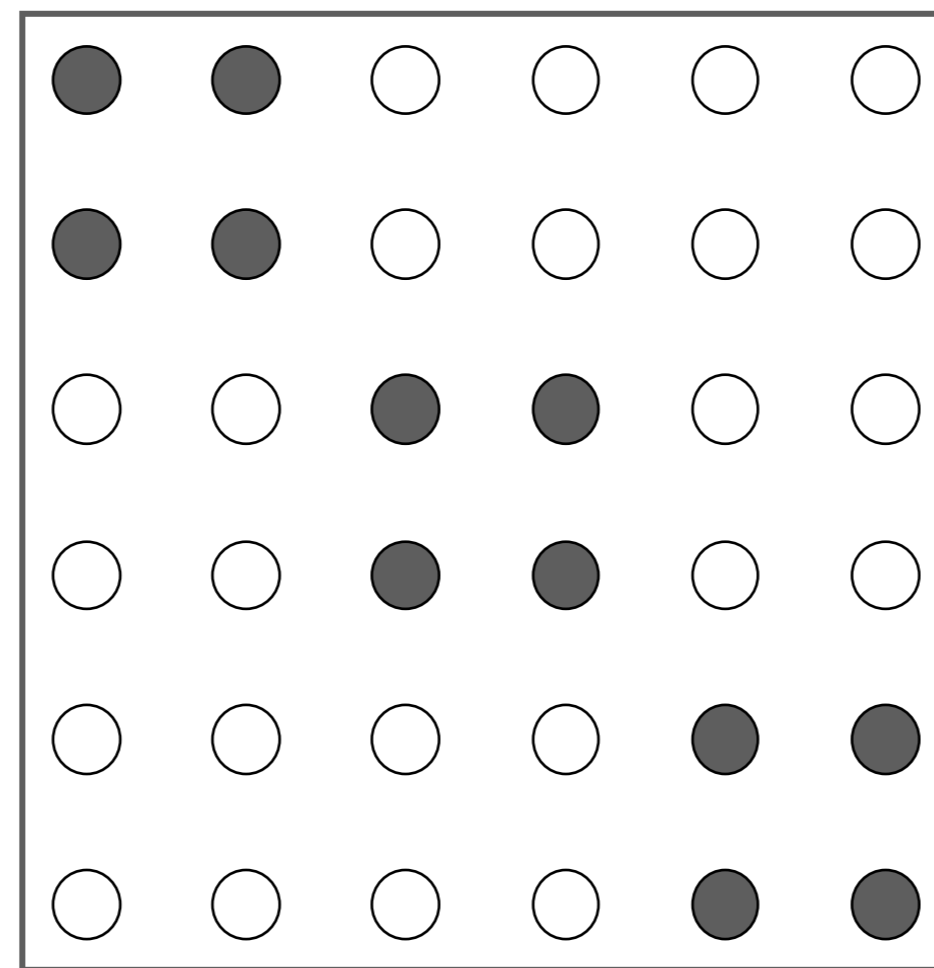


Sparsity Pattern and Ratio

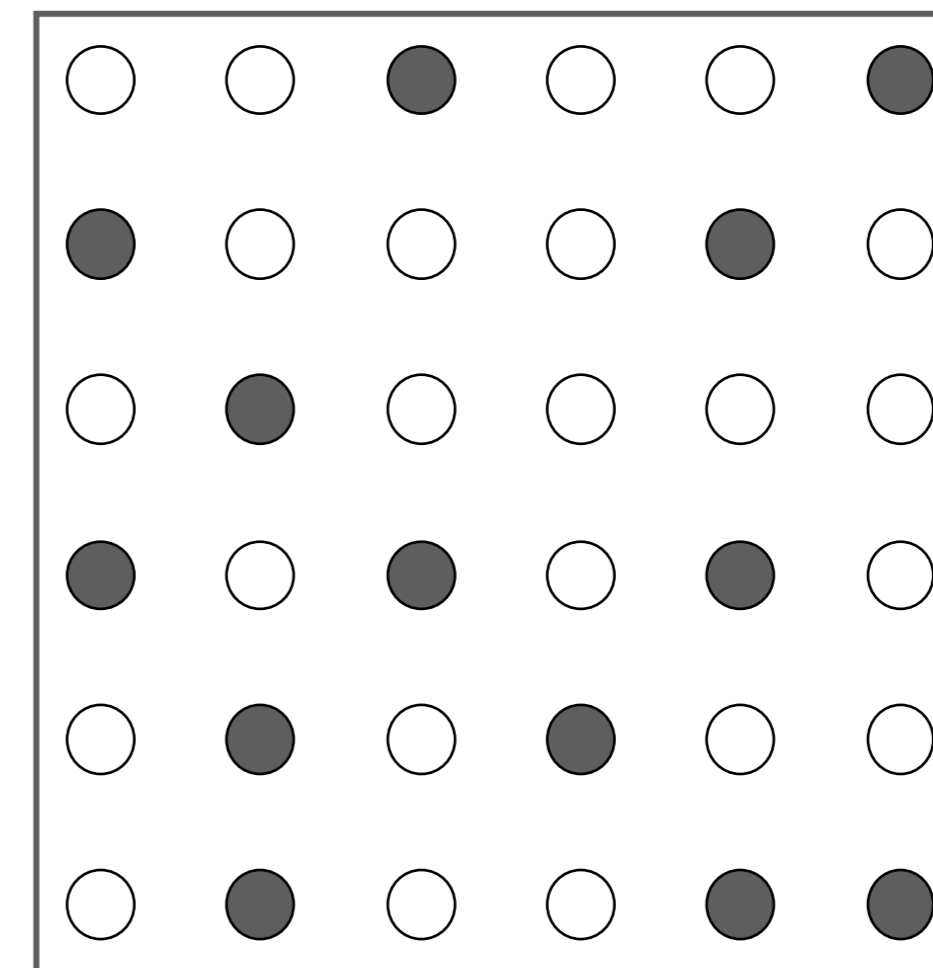
The pattern and location of non-zeros in a sparse matrix is called the **sparsity pattern** and the ratio of non-zeros over all the elements in the matrix is **density ratio**.



Dense A

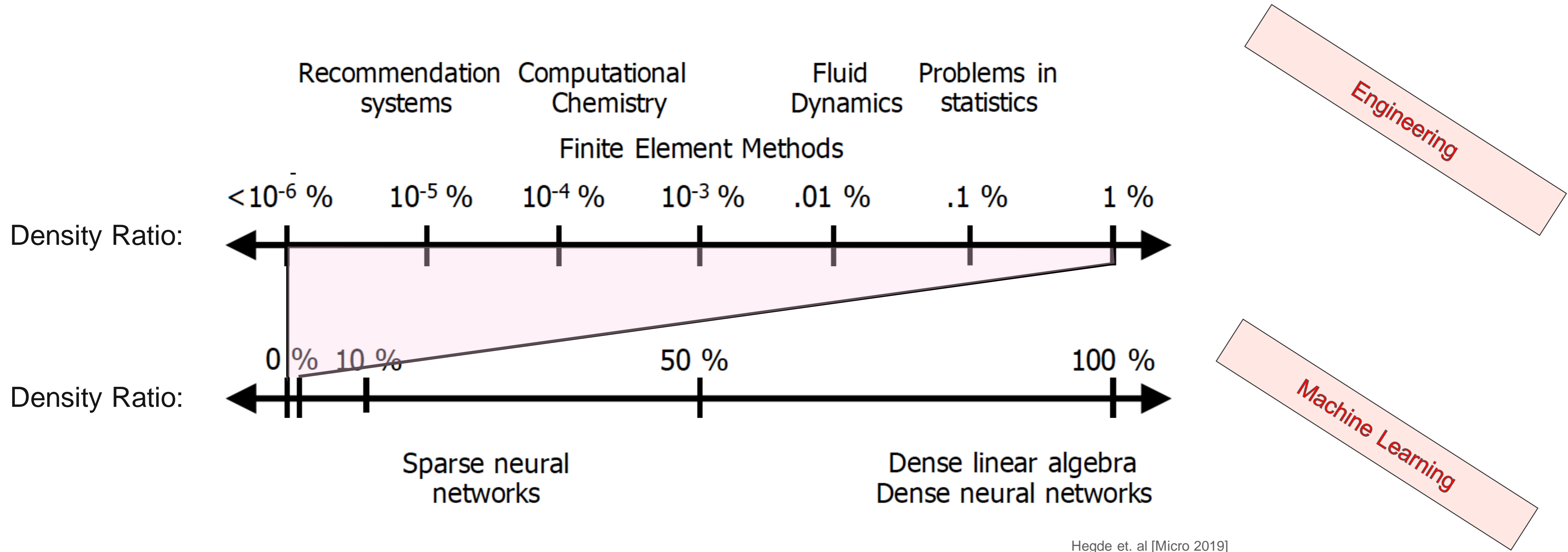


Sparse B



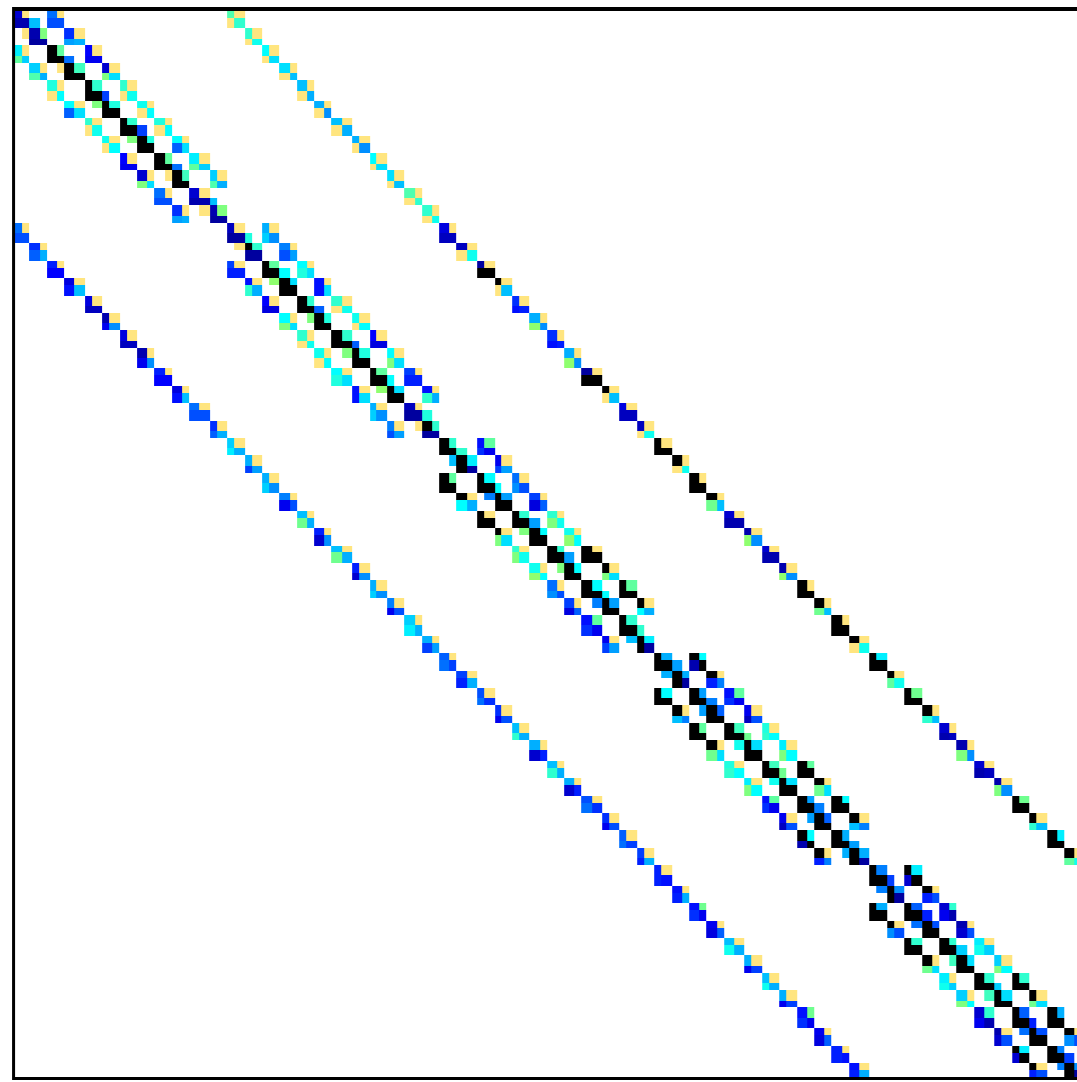
Sparse C

Density Ratios in Machine Learning vs. Engineering

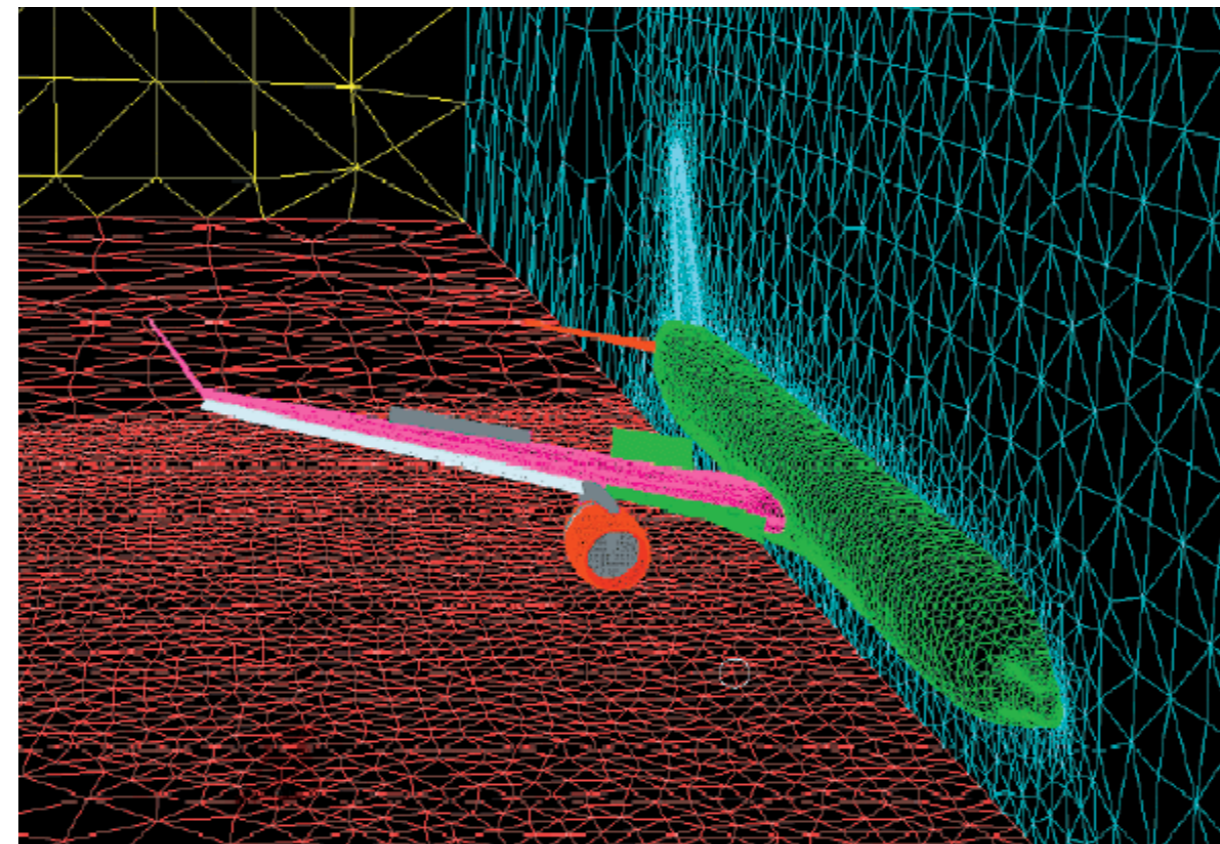


On average DNN matrices are 13x less sparse compared to scientific matrices.

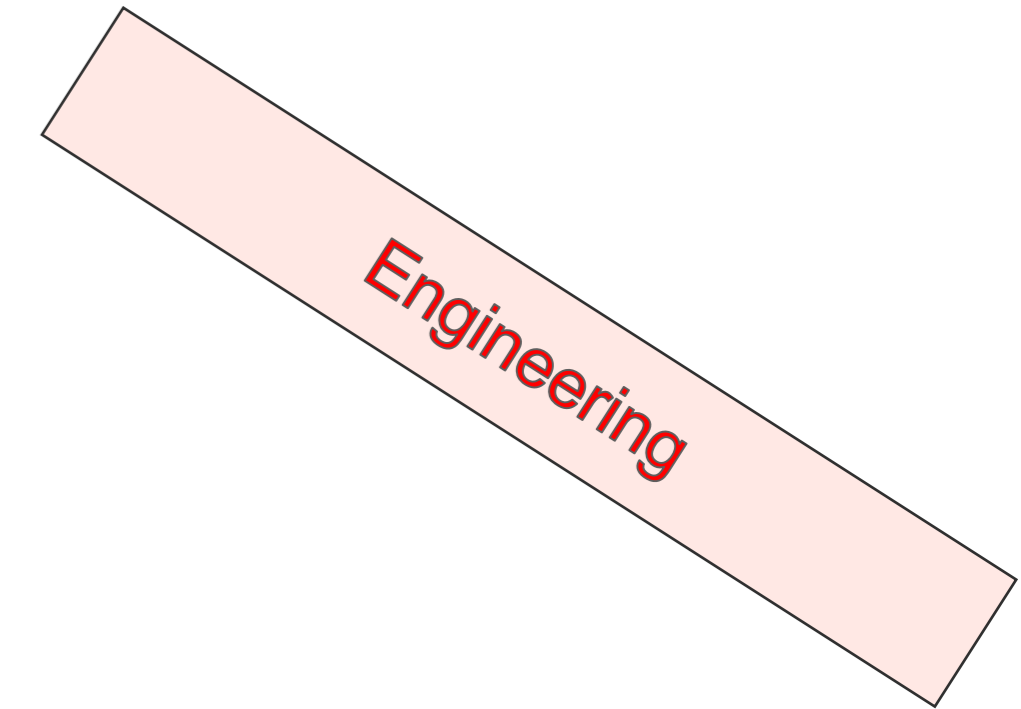
Sparsity Patterns in Machine Learning vs. Engineering



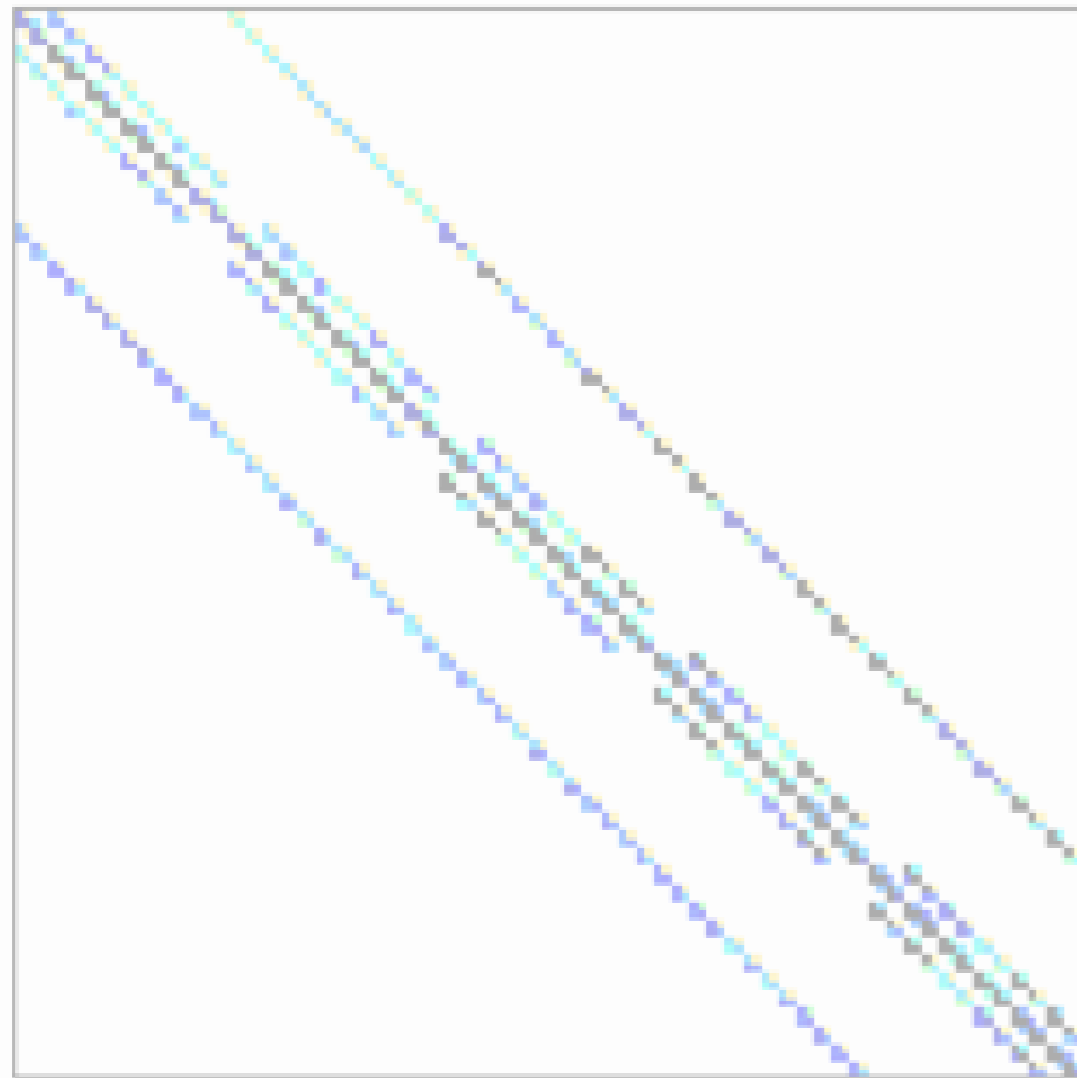
Density Ratio: 2%



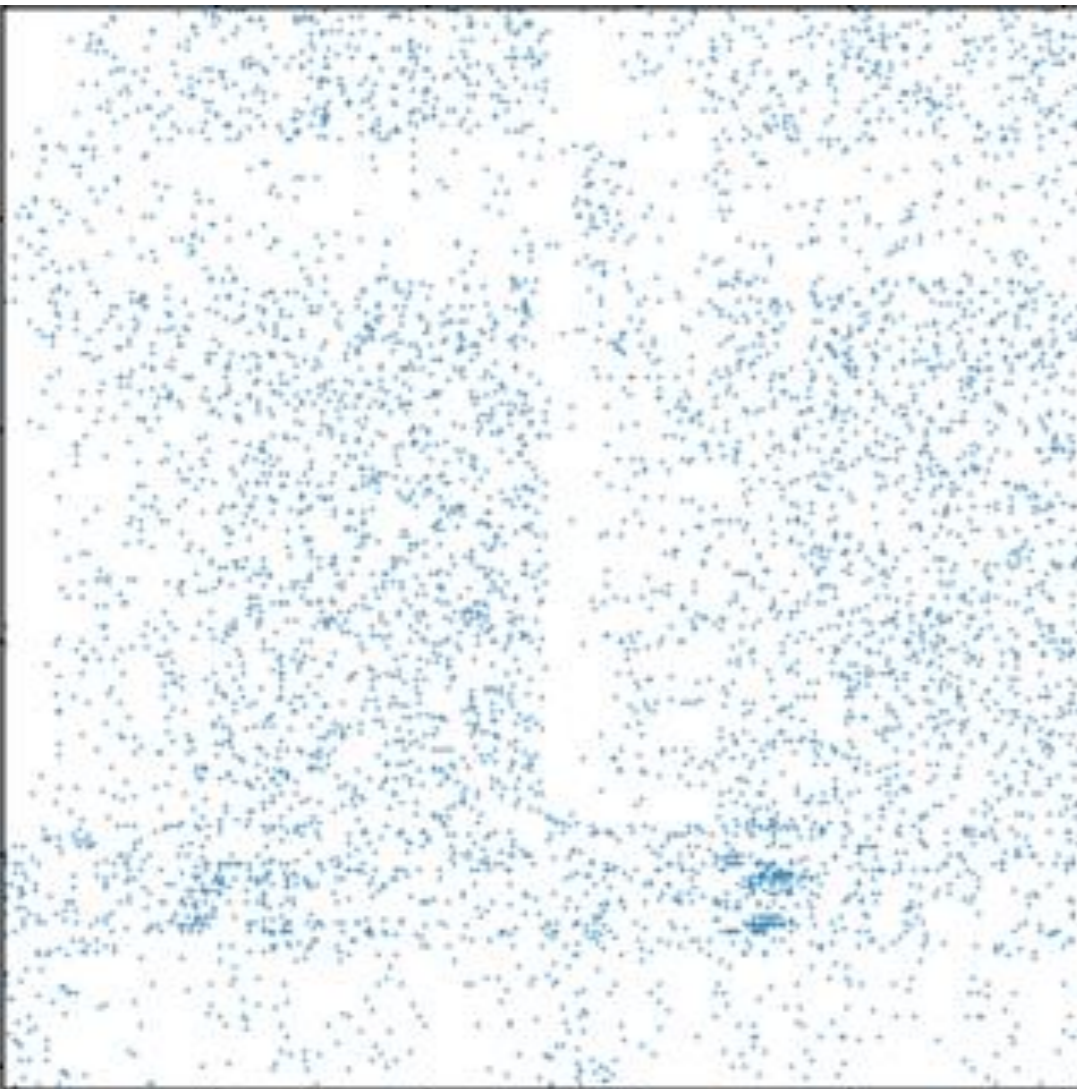
Grid structure, physical properties, etc.



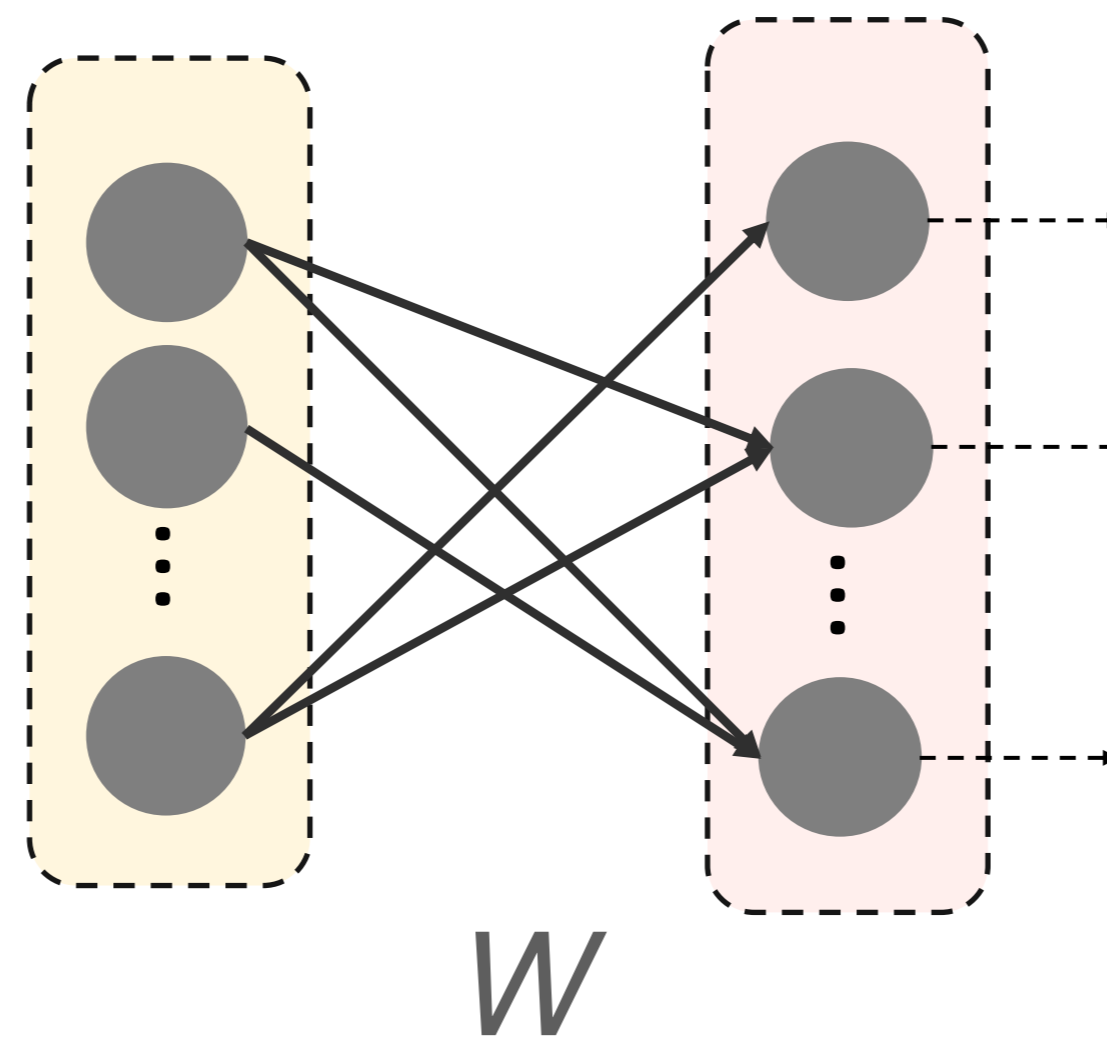
Sparsity Patterns in Machine Learning vs. Engineering



Computational Fluid Dynamics Problem
Density Ratio: 98%



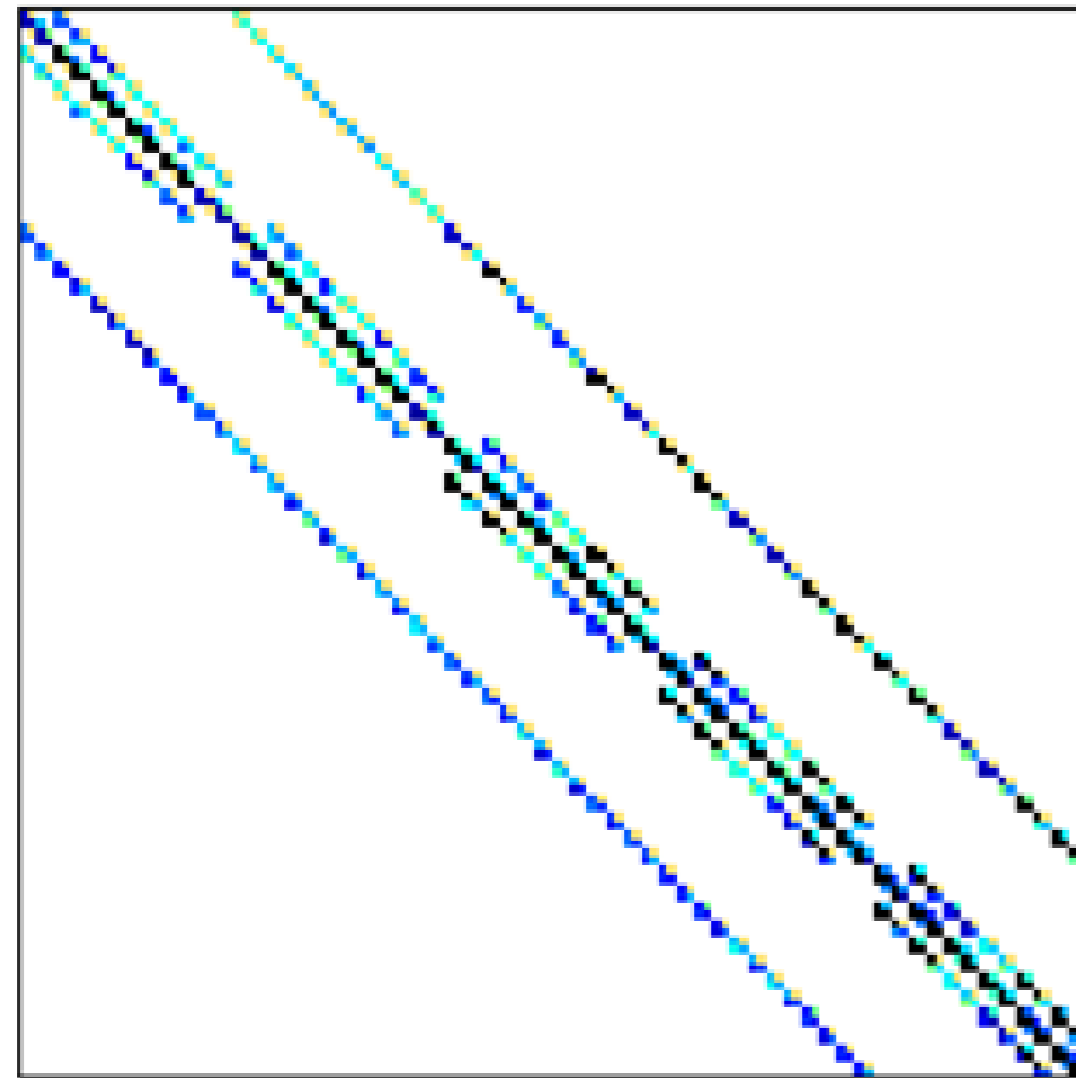
Density Ratio: 2%



Engineering

Machine Learning

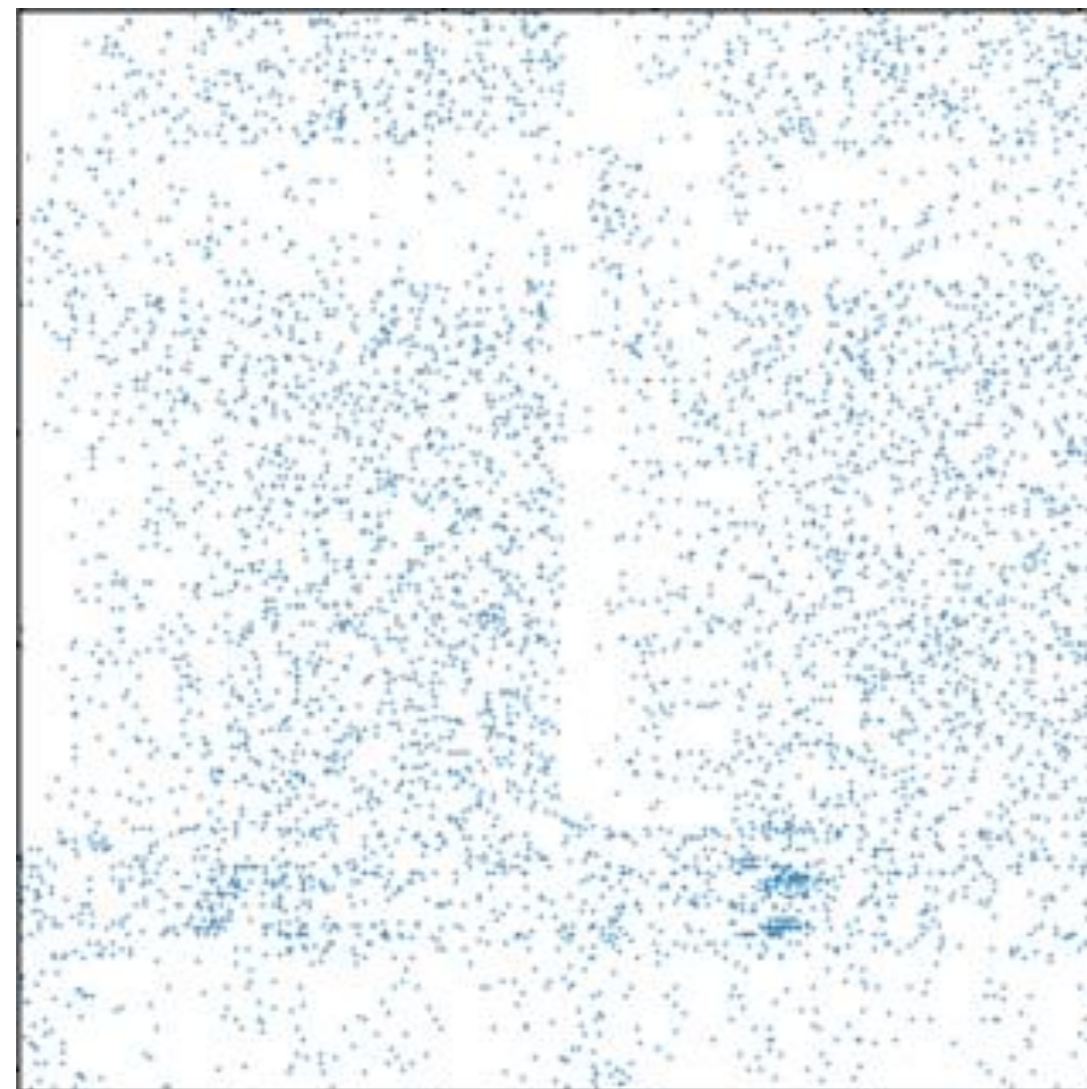
Sparsity Patterns in Machine Learning vs. Engineering



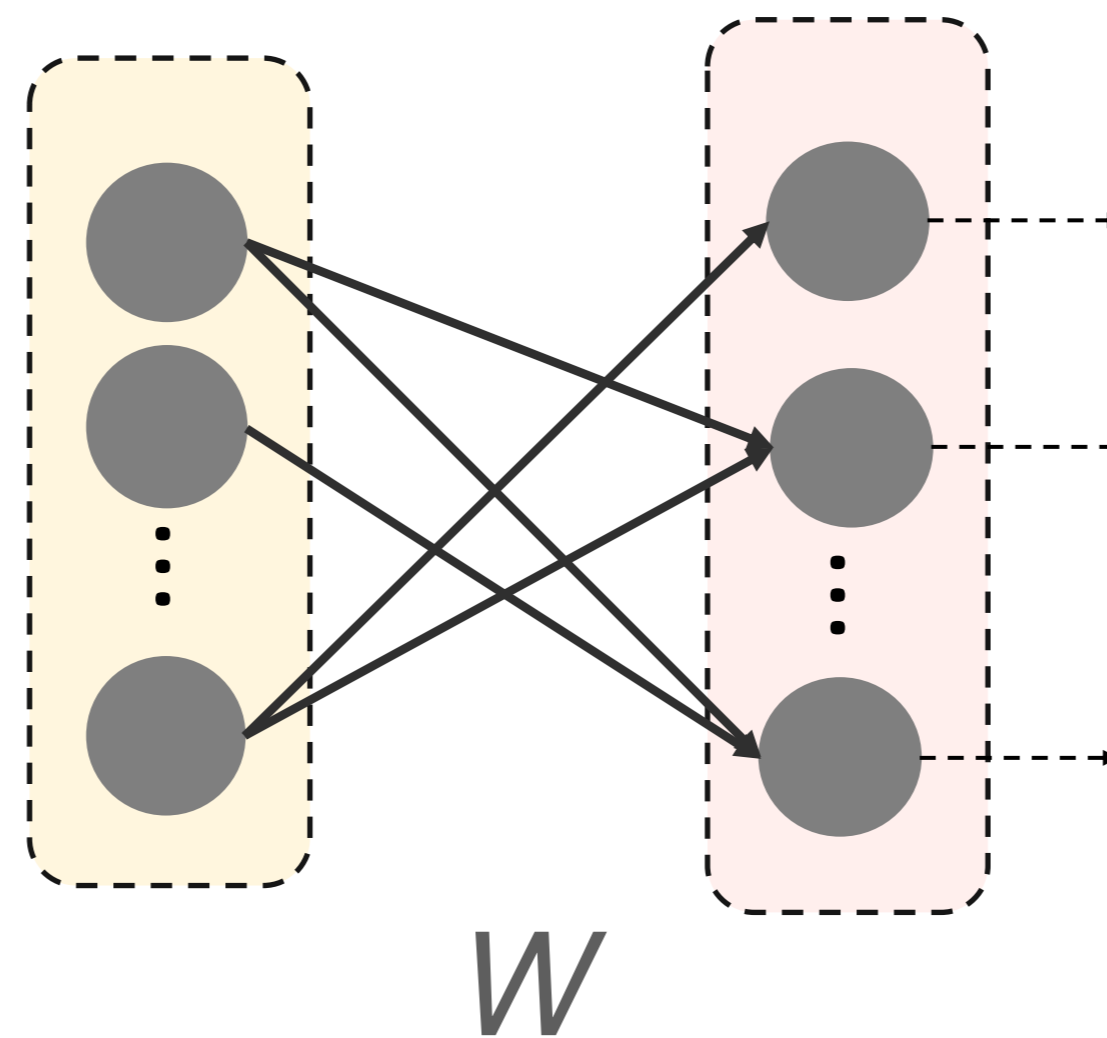
Computational Fluid Dynamics Problem
Density Ratio: 98%

Sparsity is often static or moderately changes

Engineering



Density Ratio: 2%



Weight sparsity static during inference

Machine Learning

There is no One-Size-Fits-All Approach!

Thesis

Sparsity is often static or changes moderately in most simulations.

Key contribution

Inspect sparsity patterns to **automatically generate highly-optimized code for sparsity.**

Part 1: Sparse Matrix Multiplication in ML & Engineering

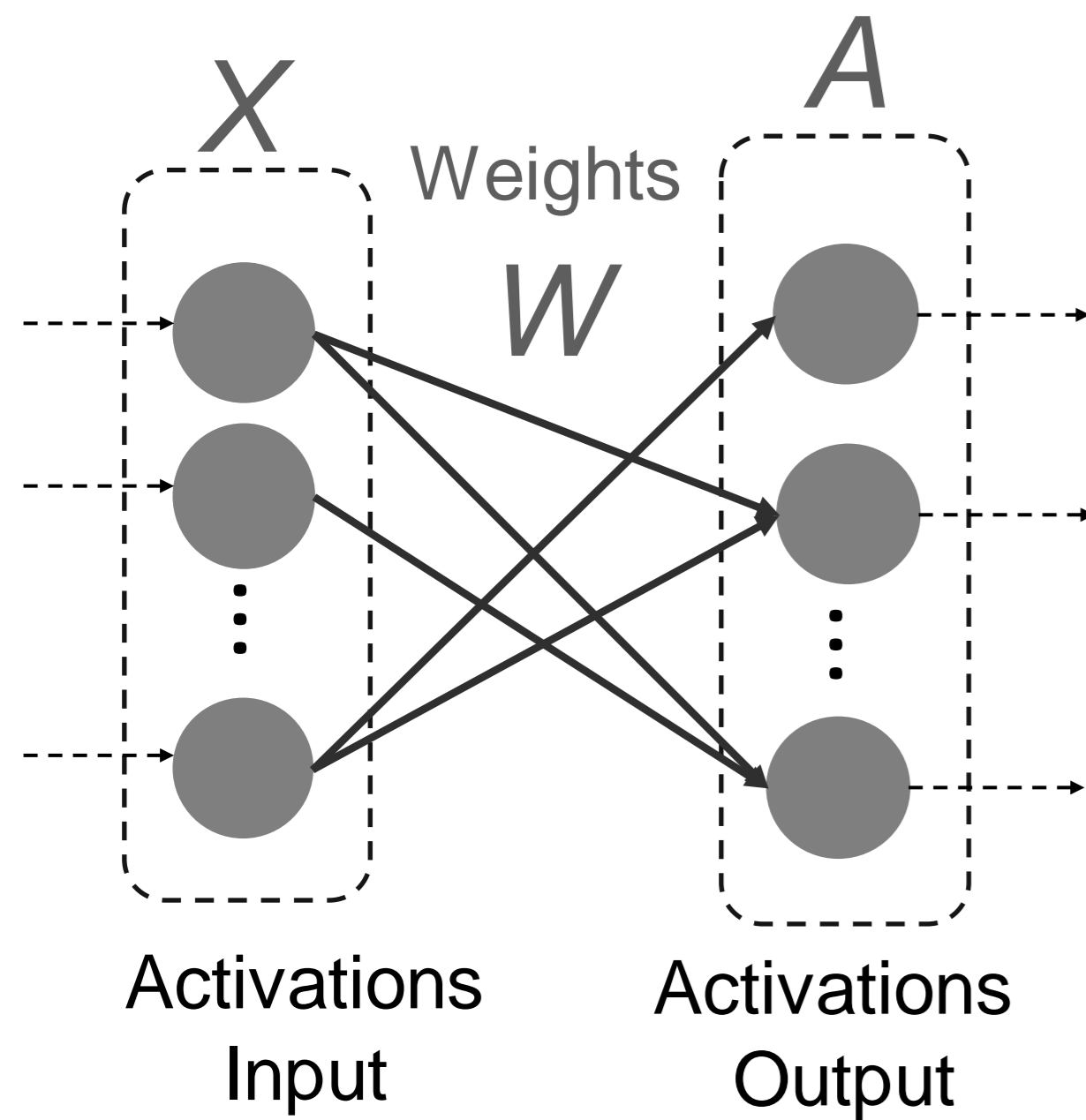
```
L, ... ← factorization(W)
while (not converged) {
...
Solve:  $Lx_{k+1} = b$ 
 $r_{k+1} = r_k - \alpha Wx_{k+1}$ 
....
}
```

Part 1

Matrix Multiplication routines

Engineering

Machine Learning



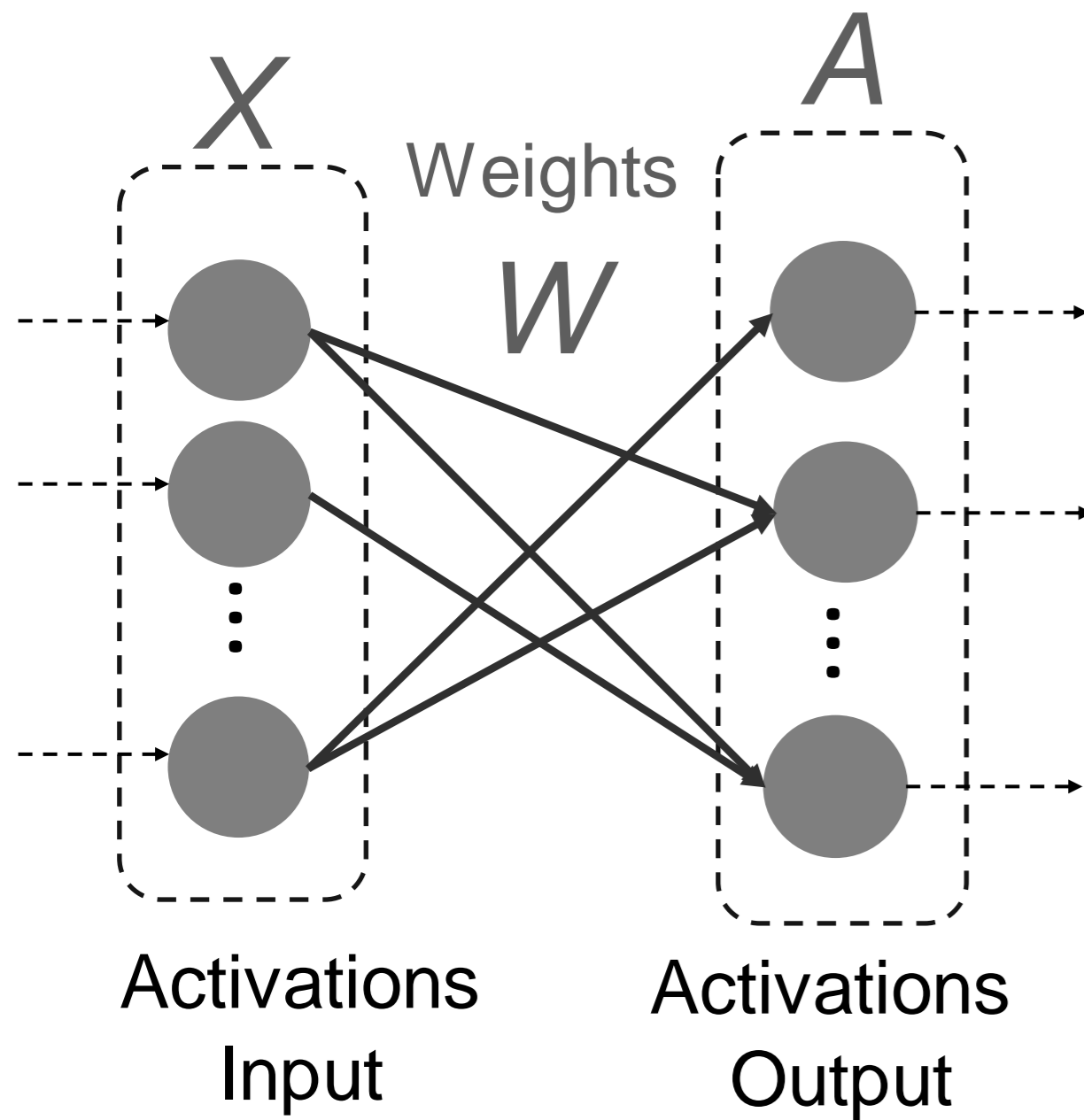
$$A = WX$$

Part 1: Sparse Matrix Multiplication in ML & Engineering

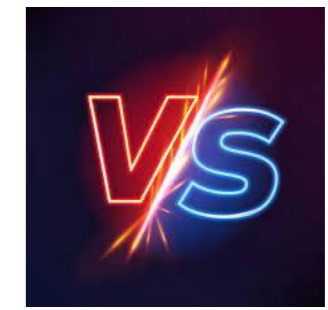
```
L, ... ← factorization(W)
while (not converged) {
...
Solve:  $Lx_{k+1} = b$ 
 $r_{k+1} = r_k - \alpha Wx_{k+1}$ 
....
}
```

Part 1

Matrix Multiplication routines



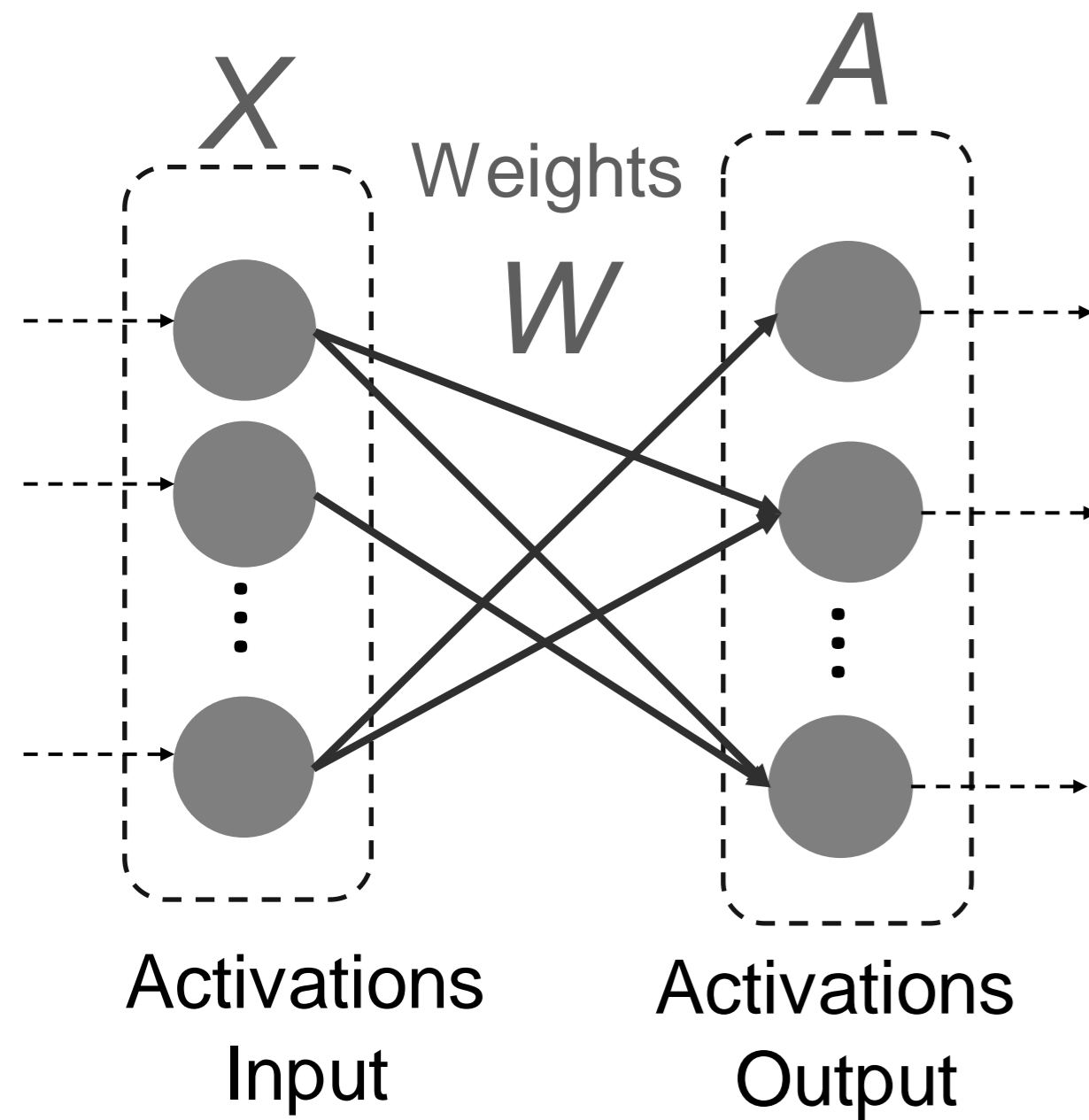
Engineering



Machine Learning

Part 1: Sparse Matrix Multiplication in ML & Engineering

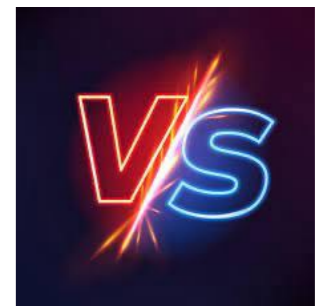
```
L, ... ← factorization(W)
while (not converged) {
...
Solve:  $Lx_{k+1} = b$ 
 $r_{k+1} = r_k - \alpha Wx_{k+1}$ 
....
}
```



$$A = WX$$

Part 1

Matrix Multiplication routines



Engineering

Machine Learning

Inspect memory access patterns
for spatial and temporal **locality**

Part 2: Sparse Solvers in Engineering

```
L, ... ← factorization(W)
while (not converged) {
...
  Solve:  $Lx_{k+1} = b$ 
   $r_{k+1} = r_k - \alpha Wx_{k+1}$ 
....
}
```

Engineering

PART 2

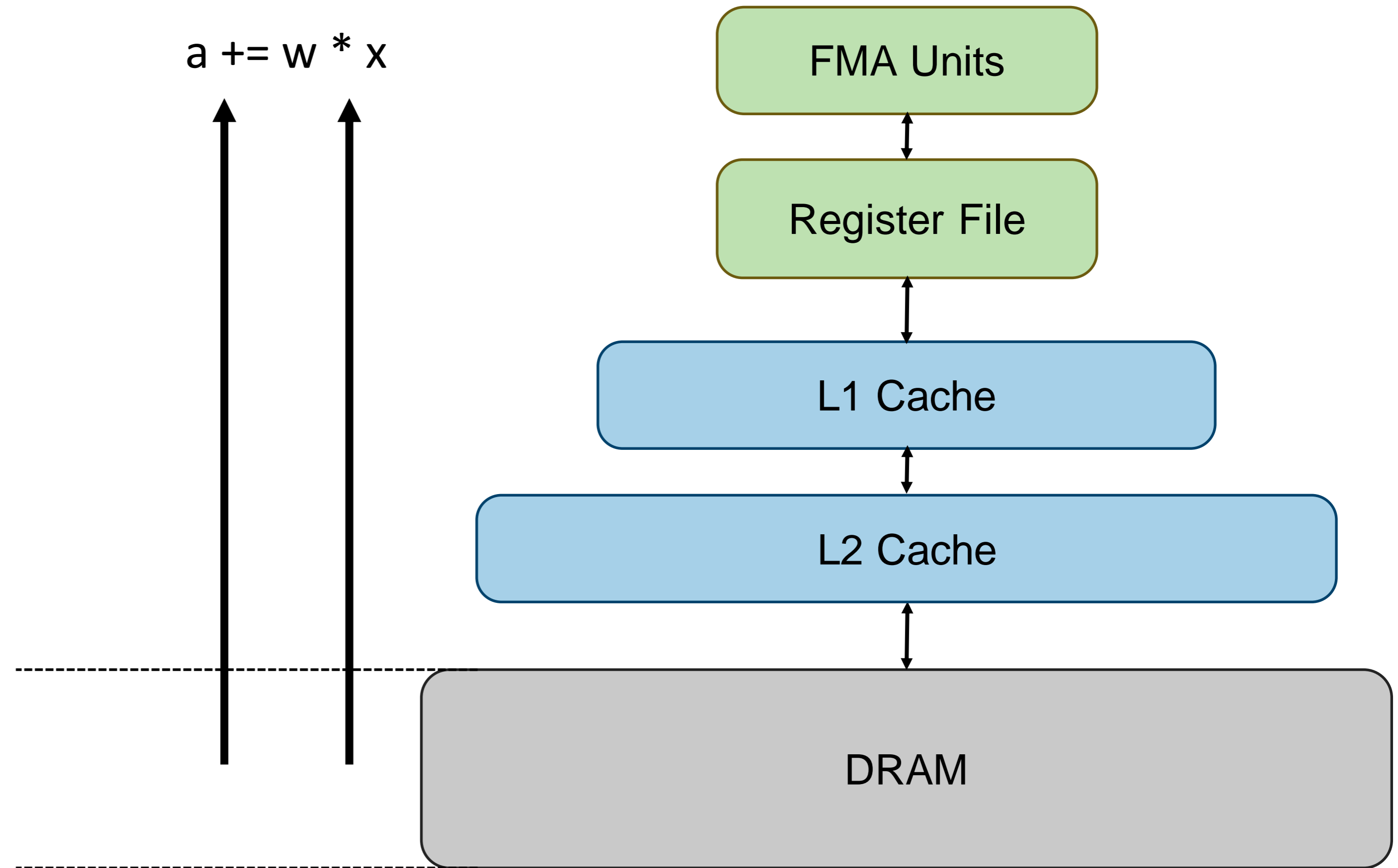
Solver routines

LU factorization, forward-backward solvers, etc.

Inspect data dependence patterns for **pruning** and **parallelism**, etc.

CPU Memory Hierarchy

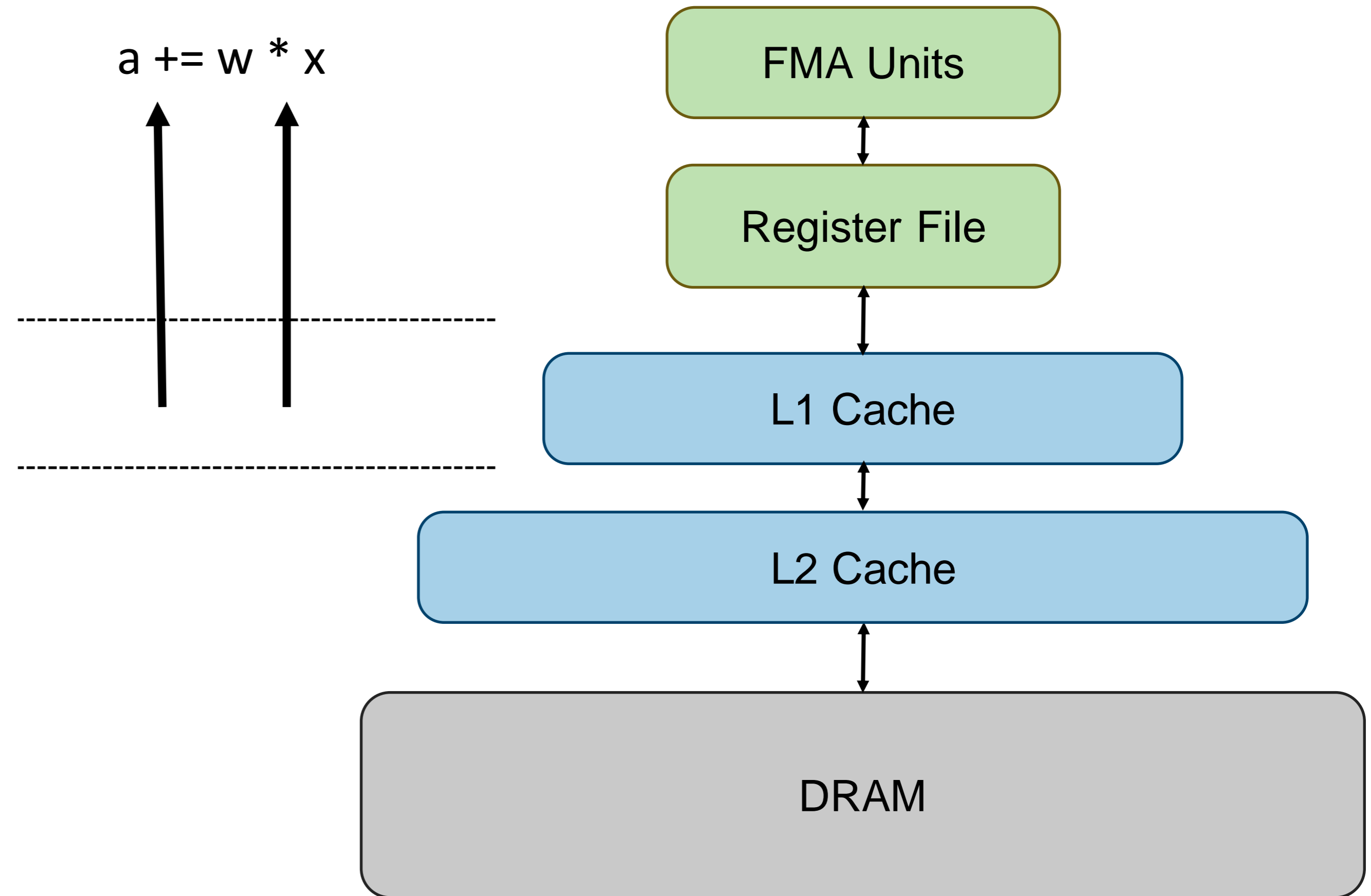
To compute $A = Wx$ values must be moved up the memory hierarchy.



Temporal Locality

To compute $A = Wx$ values must be moved up the memory hierarchy.

Temporal Locality: Reuse data while in fast memory.

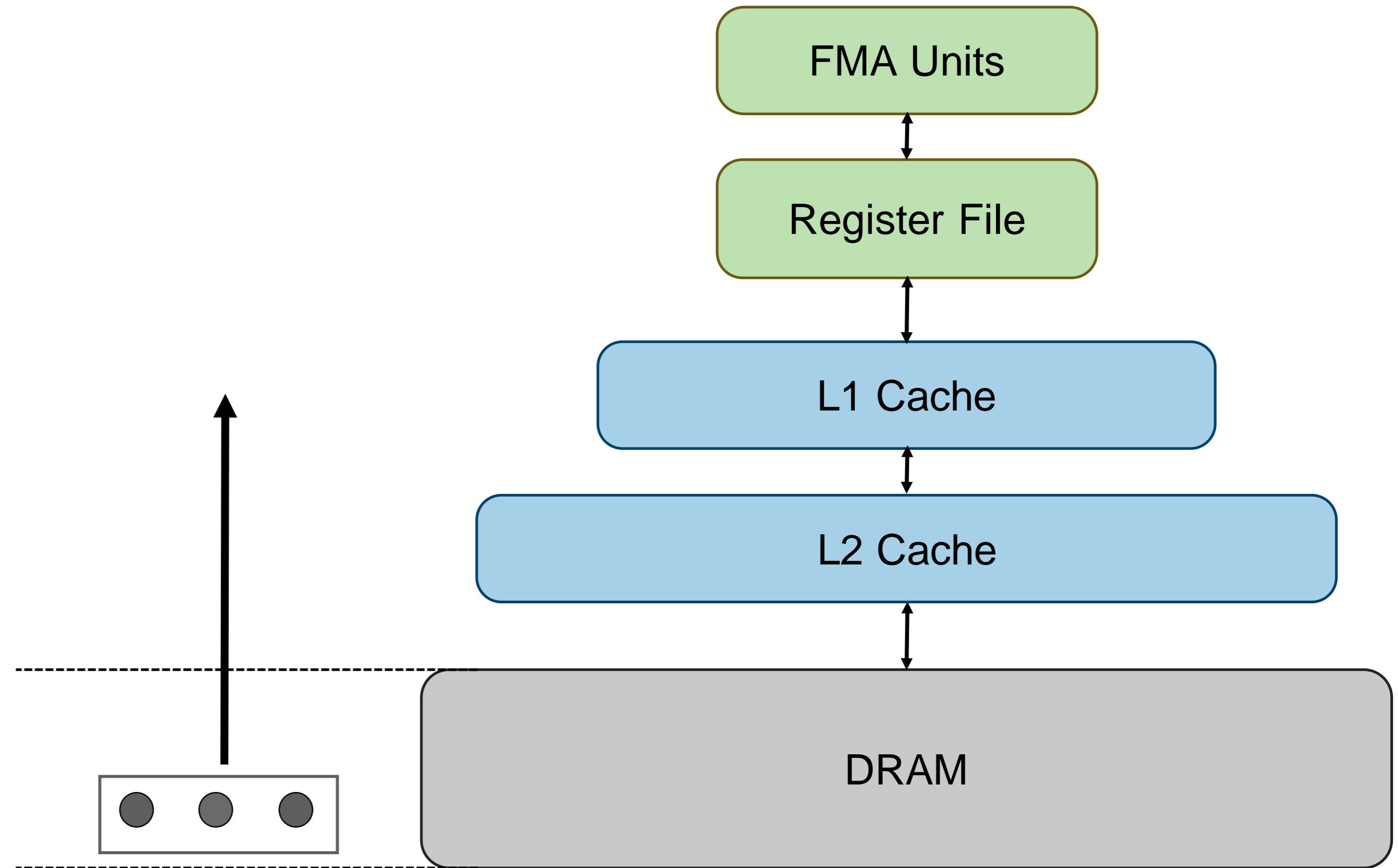


Spatial Locality

To compute $A = Wx$ values must be moved up the memory hierarchy.

Temporal Locality: Reuse data while in fast memory.

Spatial Locality: Access close-by data.



What's next!

Instruction reordering



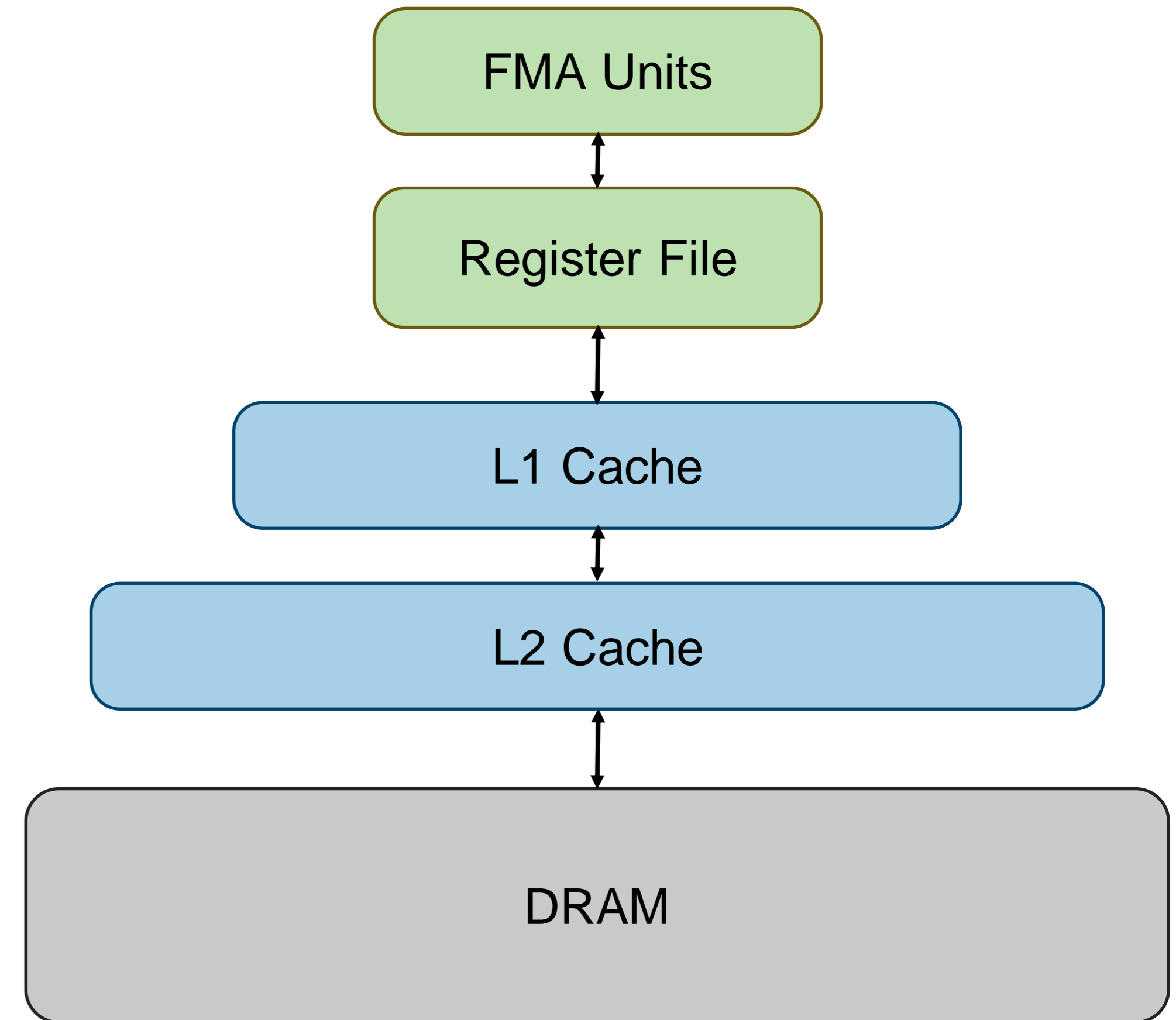
Engineering problems ✓

Machine Learning ✗

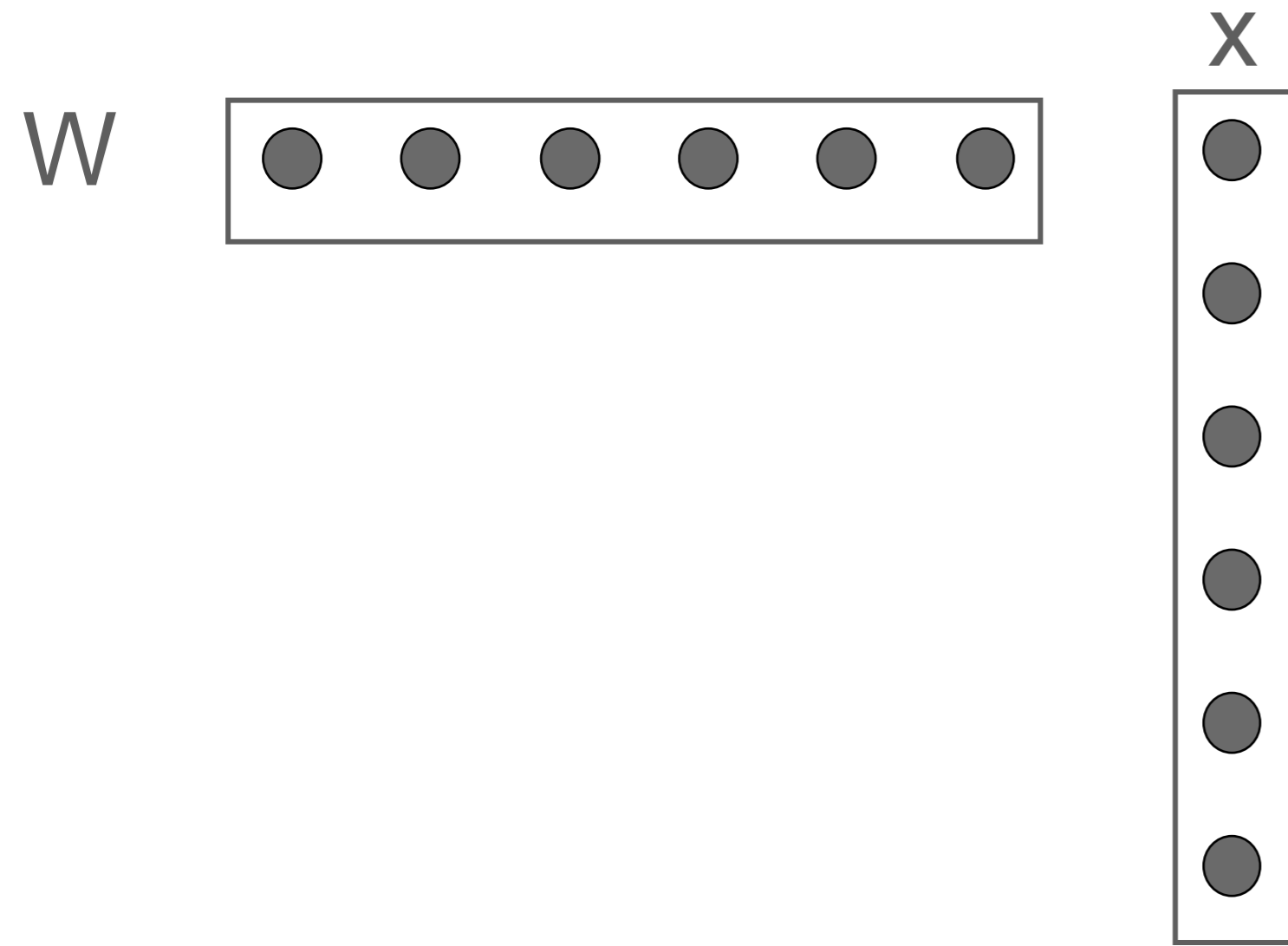
Sparse Register Tiling



Machine Learning ✓



Strided Memory Accesses and Locality



for $i = 0 \dots 5$

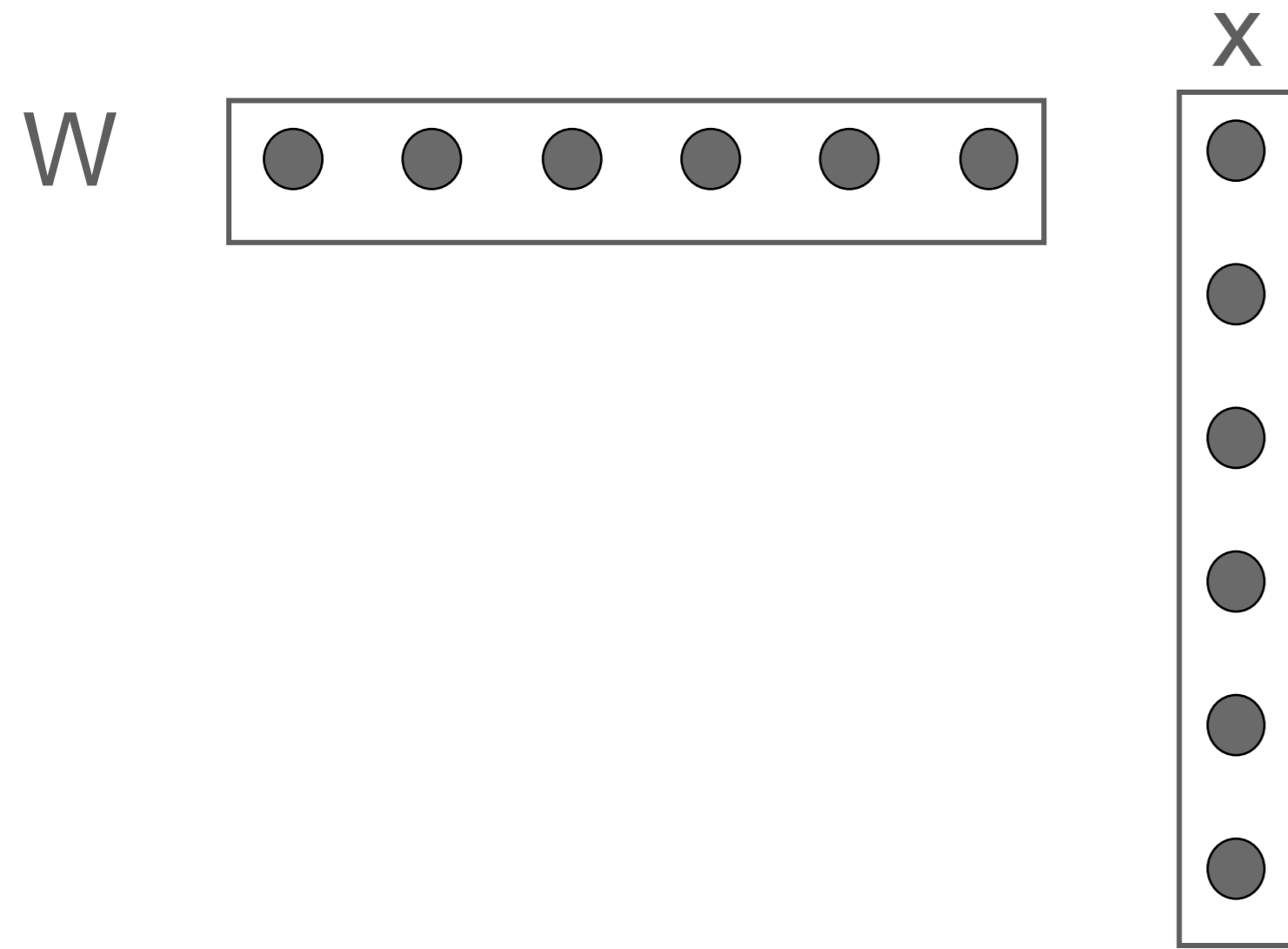
$a[1] += W[i] * x[1 * i].$

Zero-Strided Access
(Temporal locality)

Unit-Strided Access
(Spatial locality)



Strided Memory Accesses and Locality

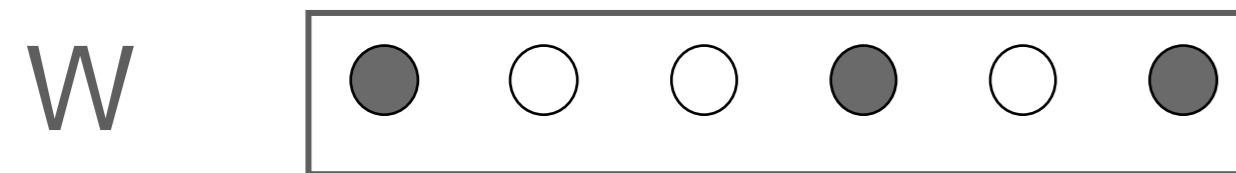


for i = 0...5

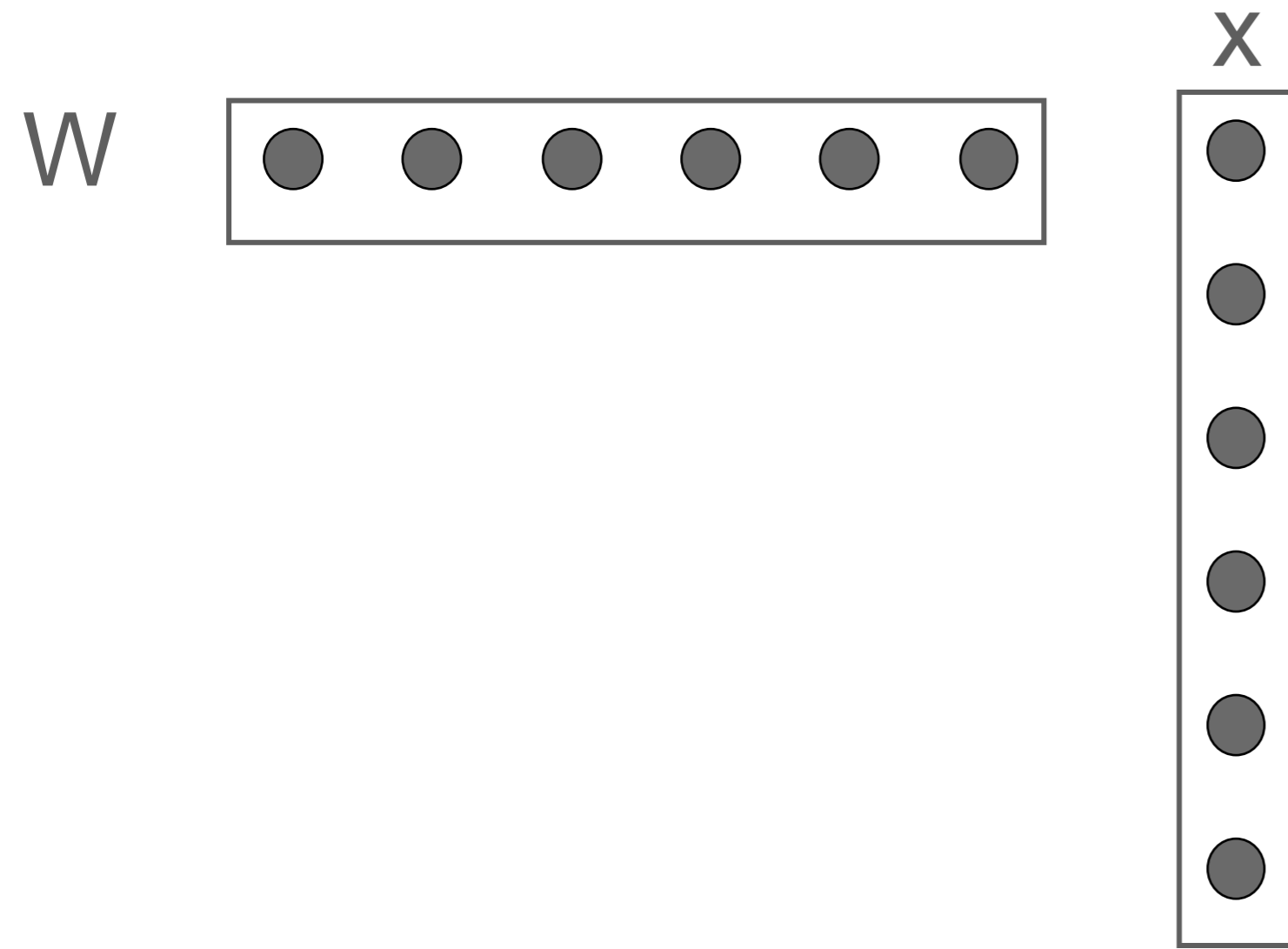
$a[1] += W[i] * x[1 * i].$

Zero-Strided Access
(Temporal locality)

Unit-Strided Access
(Spatial locality)



Strided Memory Accesses and Locality



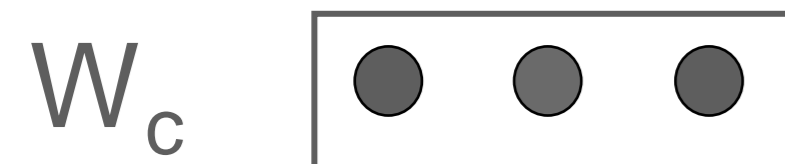
```
for i = 0...5  
  a[1] += W[i] * x[1 * i].
```

W[i]
Zero-Strided Access
(Temporal locality)

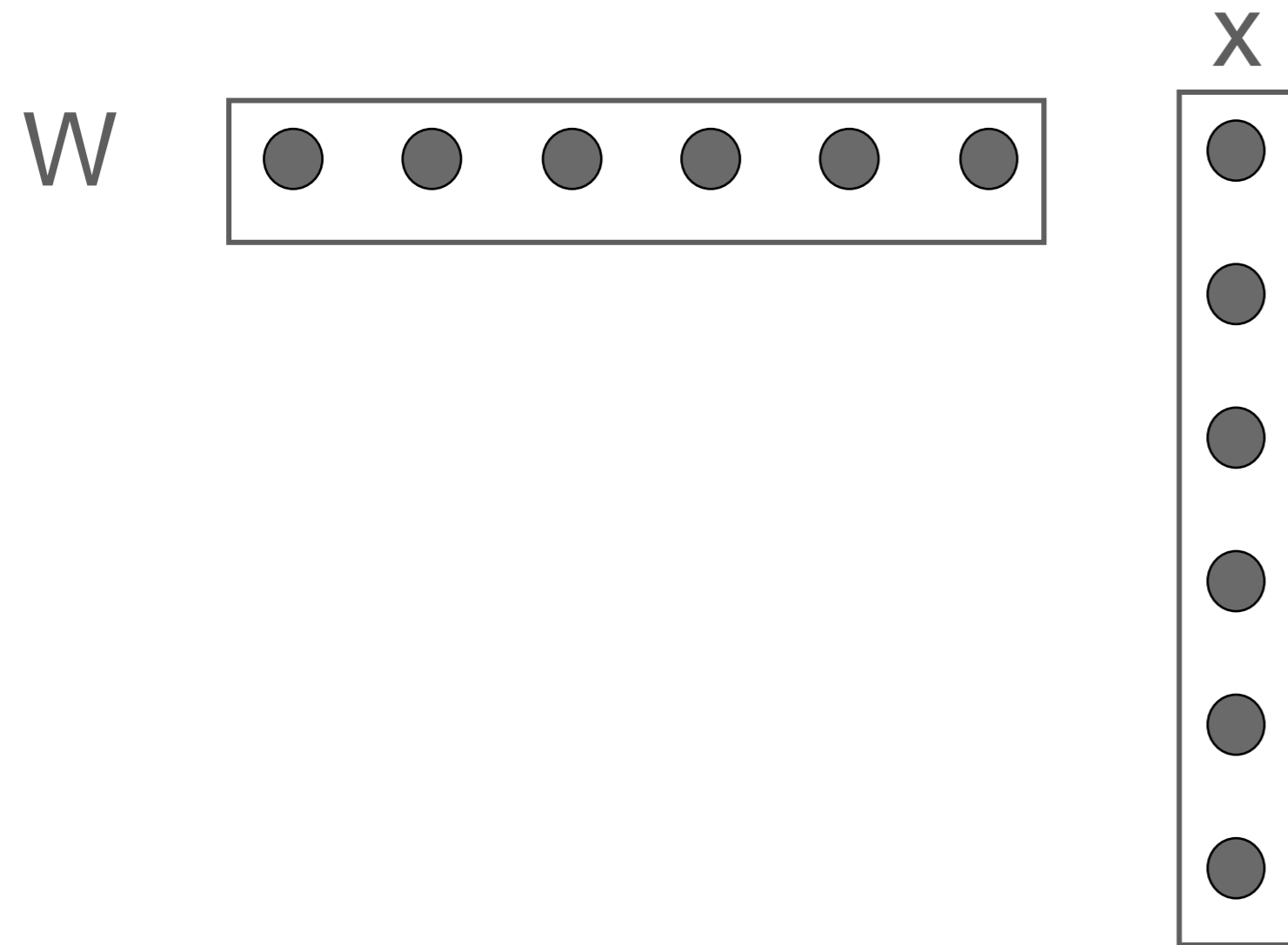
x[1 * i]
Unit-Strided Access
(Spatial locality)



Column 1 3 6



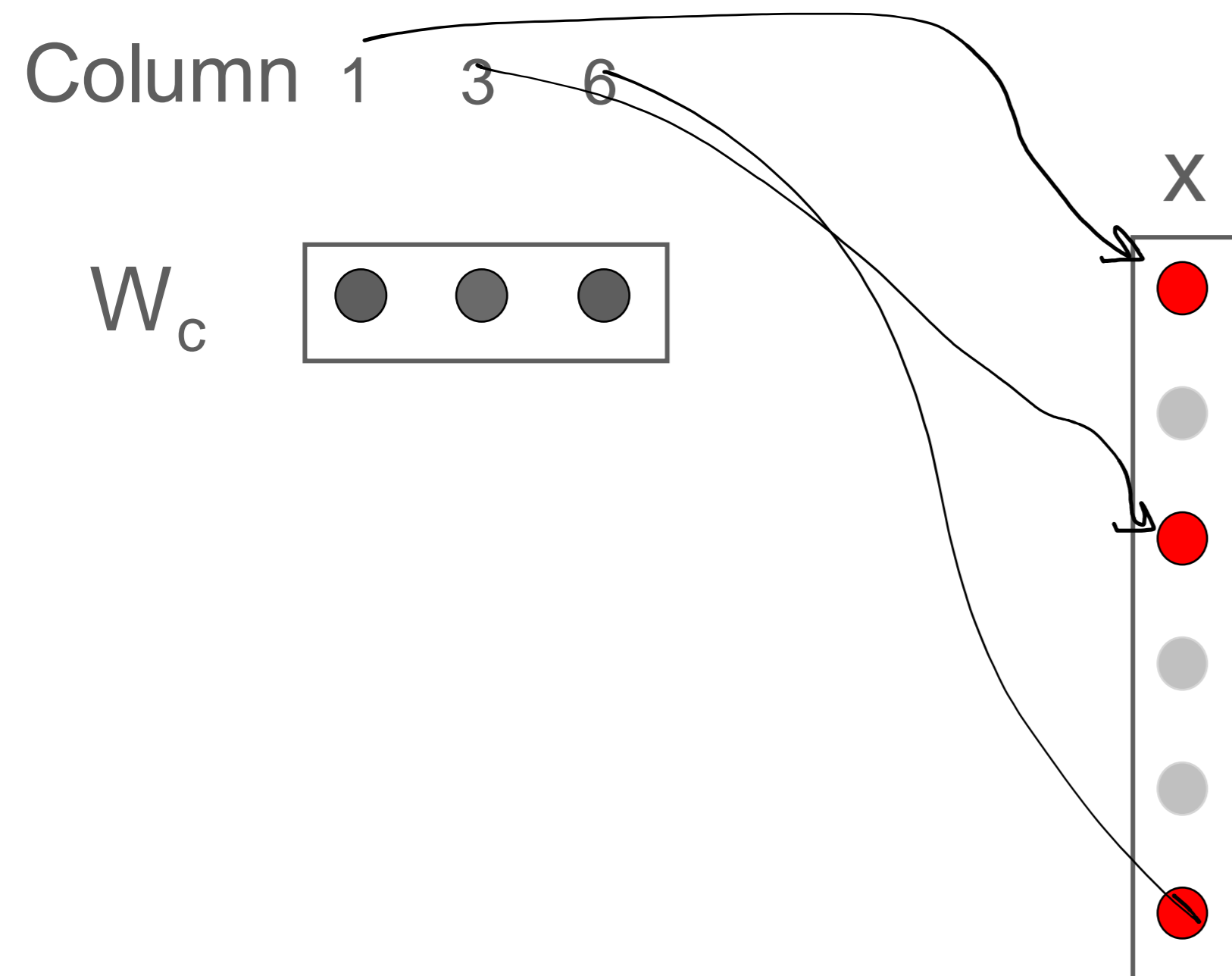
Strided Memory Accesses and Locality



```
for i = 0...5  
  a[1] += W[i] * x[1 * i].
```

Zero-Strided Access
(Temporal locality)

Unit-Strided Access
(Spatial locality)



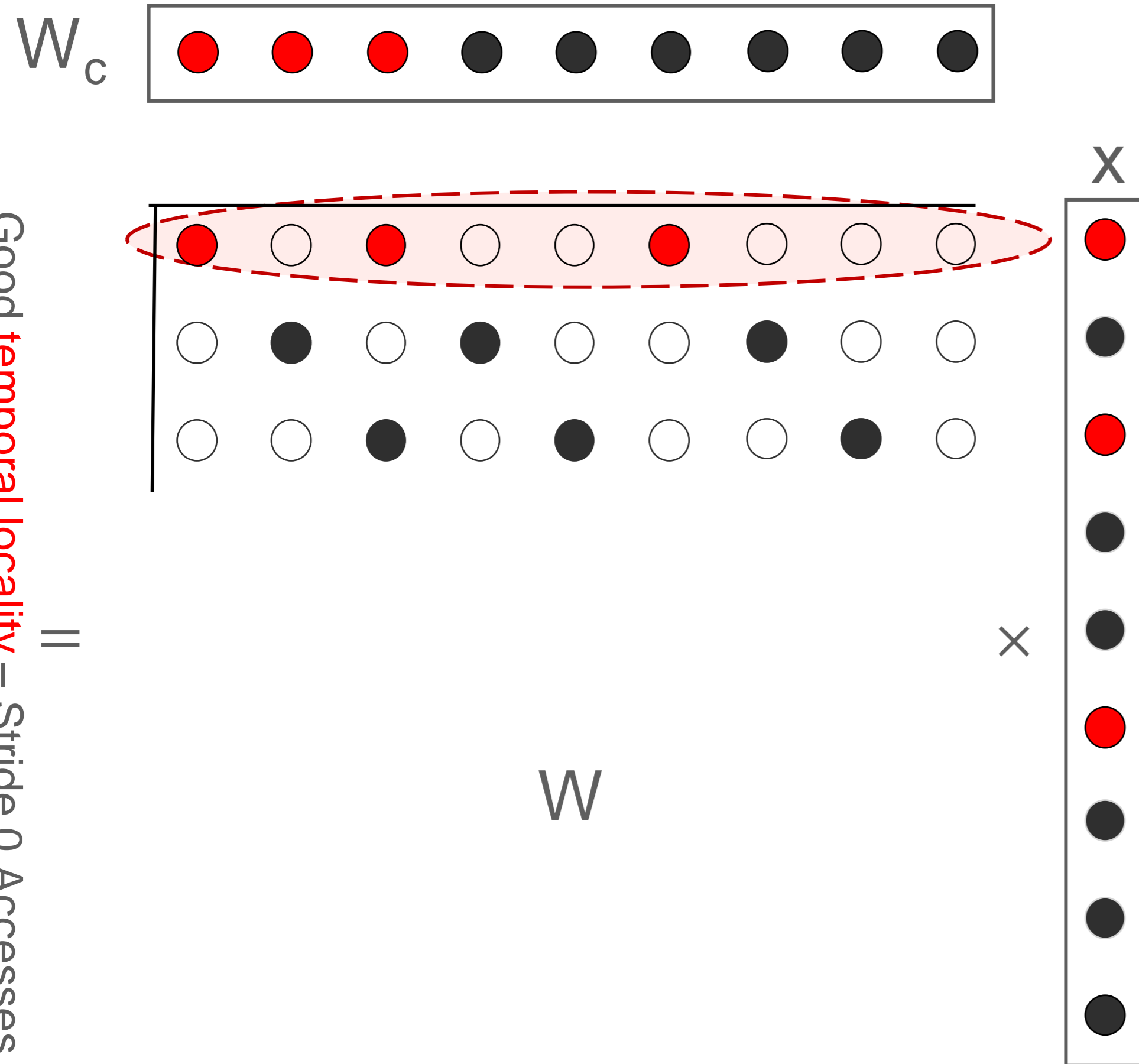
```
for i = 0...2  
  a[1] += Wc[i] * x[column[i]]
```

Unknown stride Access



Strided Accesses in SpMV

Good **spatial locality** – Stride 1 Accesses



Instruction Order

```

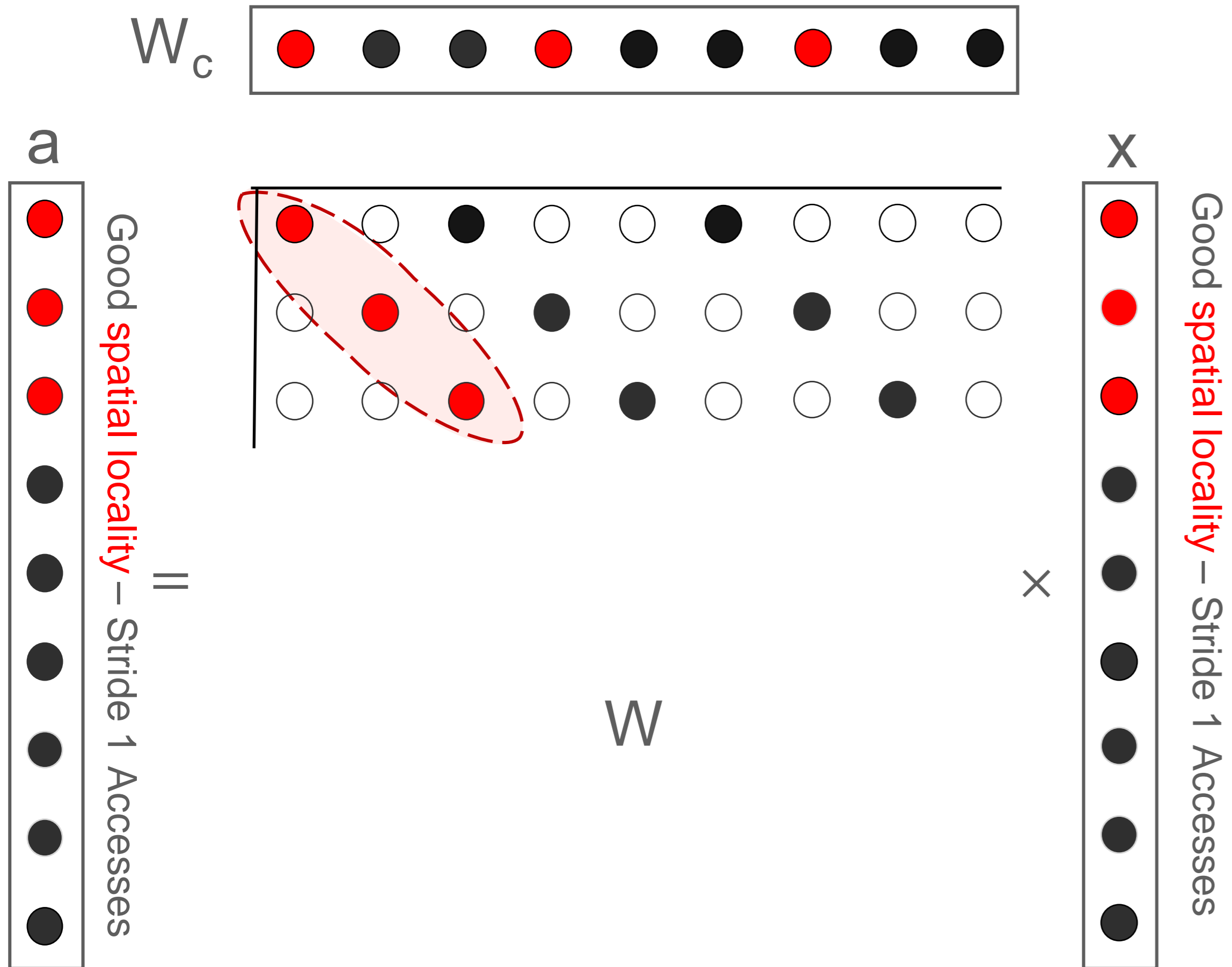
a[0] +=  $W_c[0]$  * x[0];
a[0] +=  $W_c[1]$  * x[2];
a[0] +=  $W_c[2]$  * x[5];
a[1] +=  $W_c[3]$  * x[1];
a[1] +=  $W_c[4]$  * x[3];
a[1] +=  $W_c[5]$  * x[6];
a[2] +=  $W_c[6]$  * x[2];
a[2] +=  $W_c[7]$  * x[4];
a[2] +=  $W_c[8]$  * x[7];
    
```

Code

```

for i = 0...3 {
    a[i] +=  $W_c[i*3]$  * x[i];
    a[i] +=  $W_c[i*3+1]$  * x[i+2];
    a[i] +=  $W_c[i*3+2]$  * x[i+5];
}
    
```

Strided Accesses in the New Code



Instruction Order

```

a[0] += Wc[0] * x[0];
a[1] += Wc[3] * x[1];
a[2] += Wc[6] * x[2];
a[0] += Wc[1] * x[2];
a[1] += Wc[5] * x[6];
a[2] += Wc[7] * x[4];
a[0] += Wc[2] * x[5];
a[1] += Wc[4] * x[3];
a[2] += Wc[8] * x[7];
    
```

Code

```

for i = 0...3
    a[i] += Wc[i*3] * x[i];
for i = 0...3
    a[i] += Wc[i*3+1] * x[i+2];
for i = 0...3
    a[i] += Wc[i*3+2] * x[i+5];
    
```


Which Code to Generate?

```
for i = 0...3
    a[i] += Wc[i*3] * x[i];
for i = 0...3
    a[i] += Wc[i*3+1] * x[i+2];
for i = 0...3
    a[i] += Wc[i*3+2] * x[i+5];
```



```
for i = 0...3 {
    a[i] += Wc[i*3] * x[i];
    a[i] += Wc[i*3+1] * x[i+2];
    a[i] += Wc[i*3+2] * x[i+5];
}
```

Strided accesses to a and x

Spatial locality

Spatial locality

Strided accesses to W_c and a

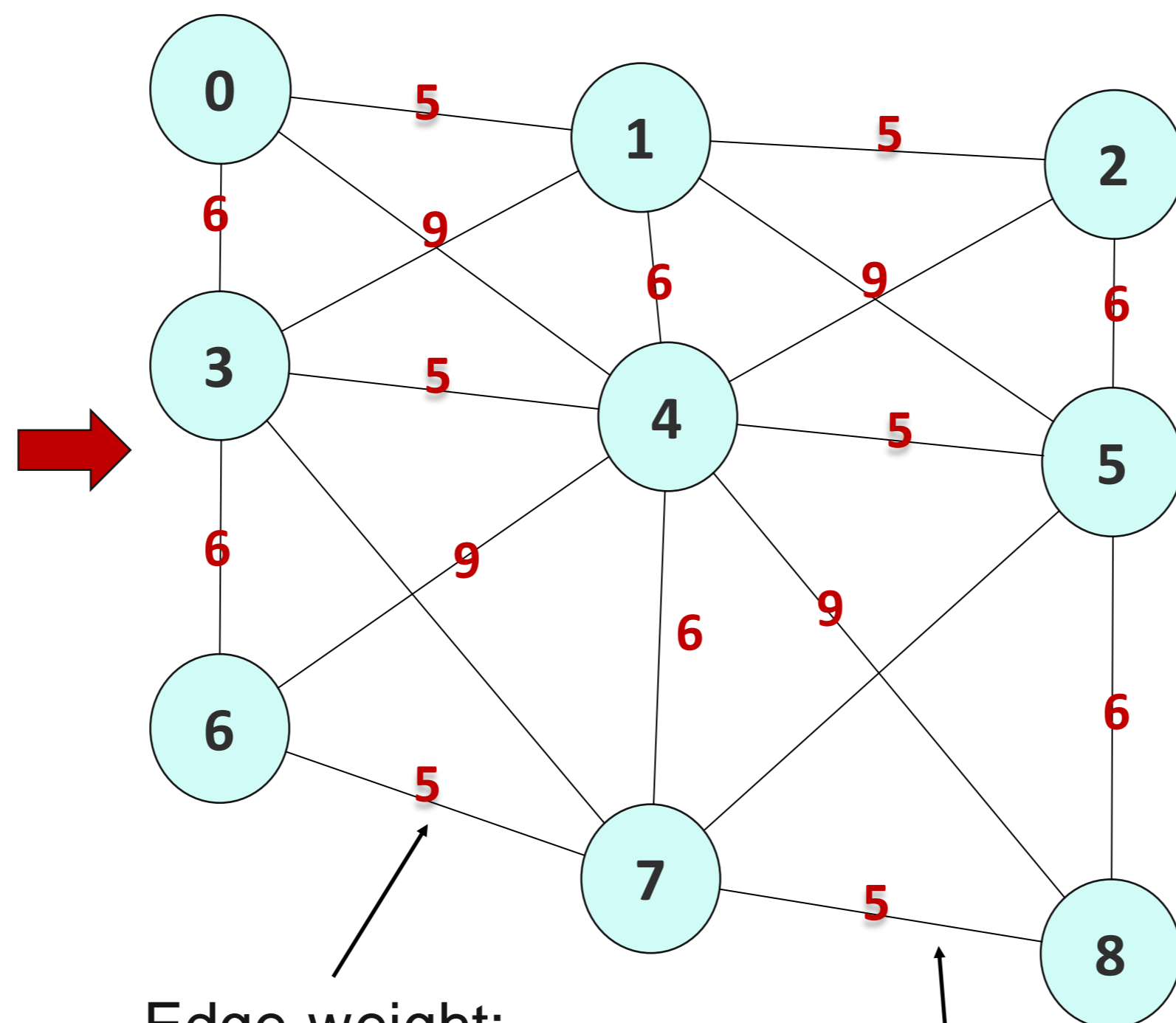
Spatial
locality

Temporal
locality

Locality-aware Codelet Mining (LCM)

Inspection

```
a[0] += Wc[0] * x[0];  
a[0] += Wc[1] * x[2];  
a[0] += Wc[2] * x[5];  
a[1] += Wc[3] * x[1];  
a[1] += Wc[4] * x[3];  
a[1] += Wc[5] * x[6];  
a[2] += Wc[6] * x[2];  
a[2] += Wc[7] * x[4];  
a[2] += Wc[8] * x[7];
```



Minimum edge covering problem:
Minimum weight-set of edges guides permutation

+

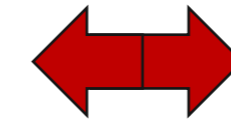
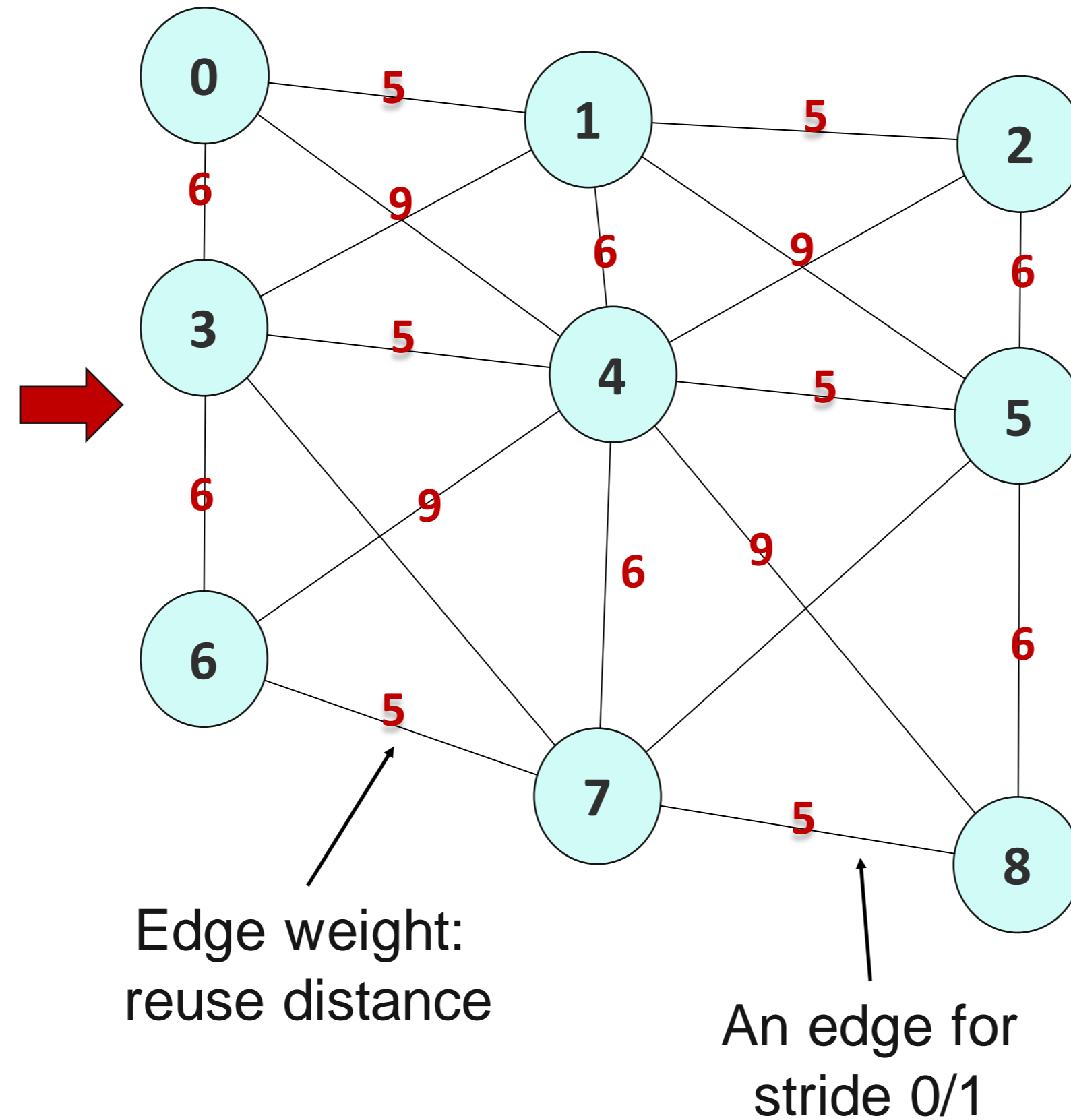
A codelet cost model



Locality-aware Codelet Mining (LCM)

```

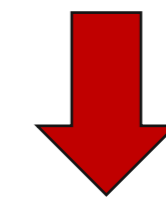
a[0] += Wc[0] * x[0];
a[0] += Wc[1] * x[2];
a[0] += Wc[2] * x[5];
a[1] += Wc[3] * x[1];
a[1] += Wc[4] * x[3];
a[1] += Wc[5] * x[6];
a[2] += Wc[6] * x[2];
a[2] += Wc[7] * x[4];
a[2] += Wc[8] * x[7];
    
```



Minimum edge covering problem:
 Minimum weight-set of edges guides permutation

+

A codelet cost model



Generate strided and partially-strided codelets.

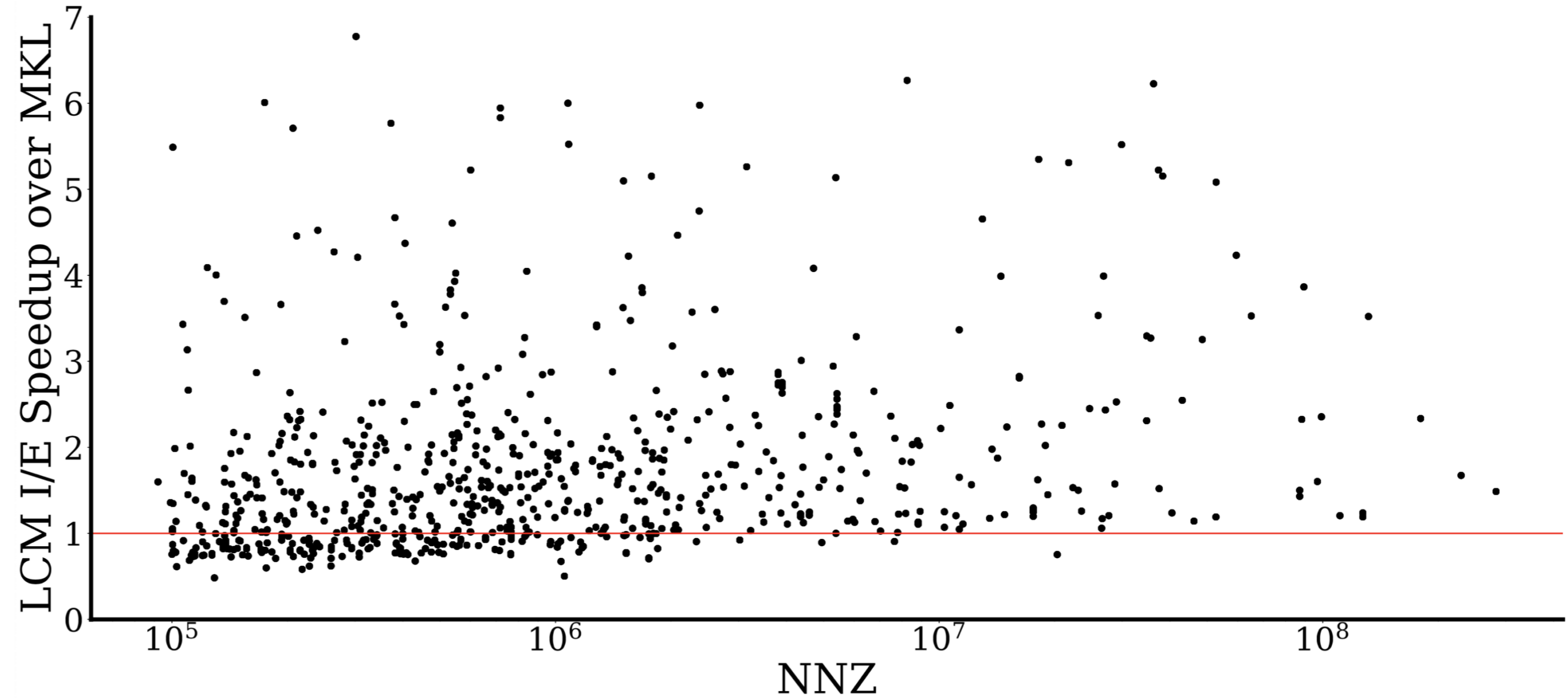
Inspection

Code Gen

LCM on Science and Engineering Problems for SpMM

Testbed: 789 matrices from the SuiteSparse¹ repository on an Intel(R) Xeon(R) Gold 5115 CPU (20 cores, 64GB main memory).

 **Engineering**



LCM is faster than Intel **MKL**² with an average speedup of **1.75x** for *Science and Engineering* problems.

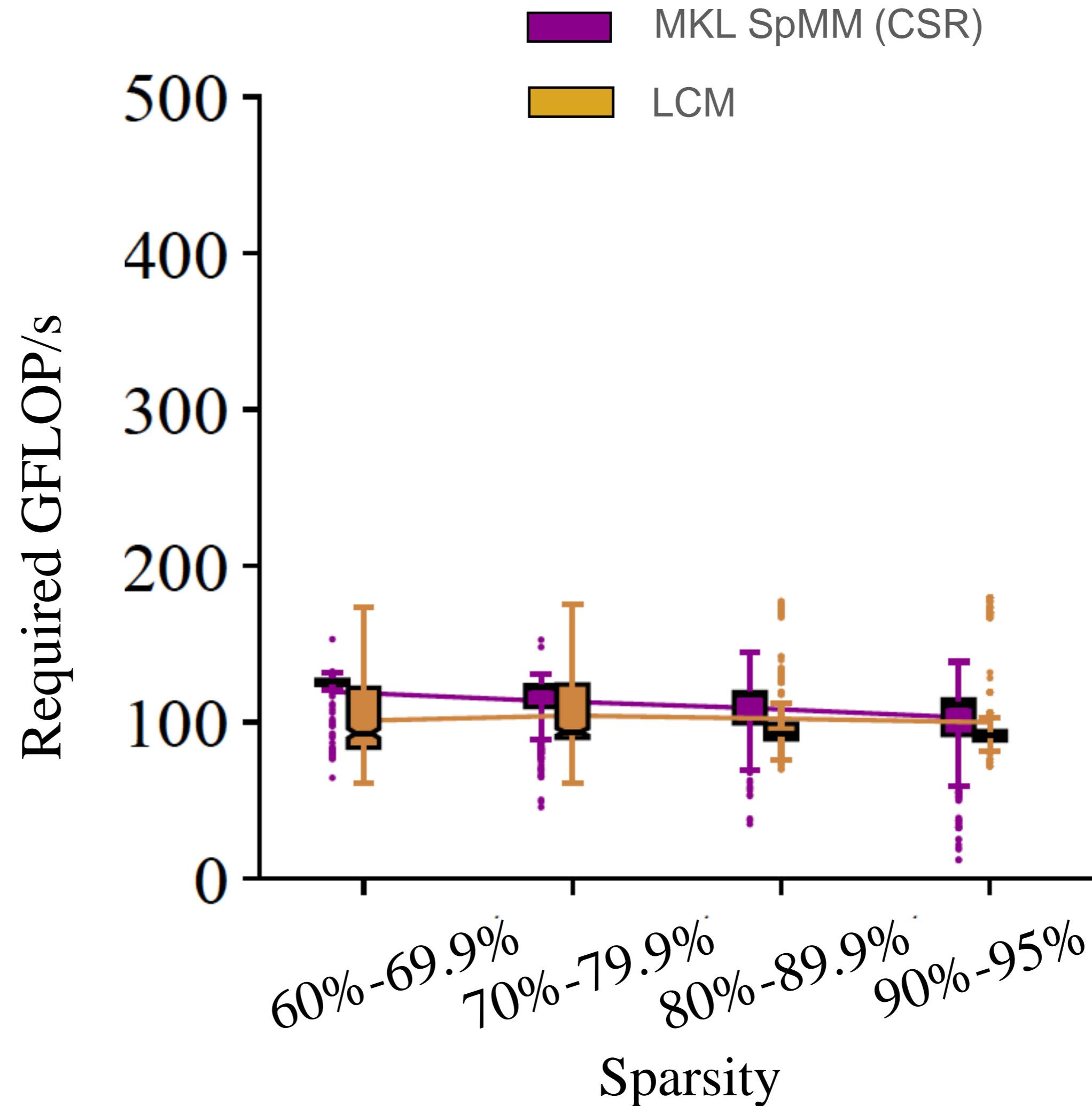
¹T. A. Davis and Y. Hu, "The university of florida sparse matrix collection," *ACM Transactions on Mathematical Software (TOMS)*, vol. 38, no. 1, pp. 1–25, 2011.

²Intel. 2022. Intel Math Kernel Library

LCM on Machine Learning Matrices

Testbed: Over 3000 matrices from Deep Learning Matrix Collection (DLMC¹).

 Machine Learning

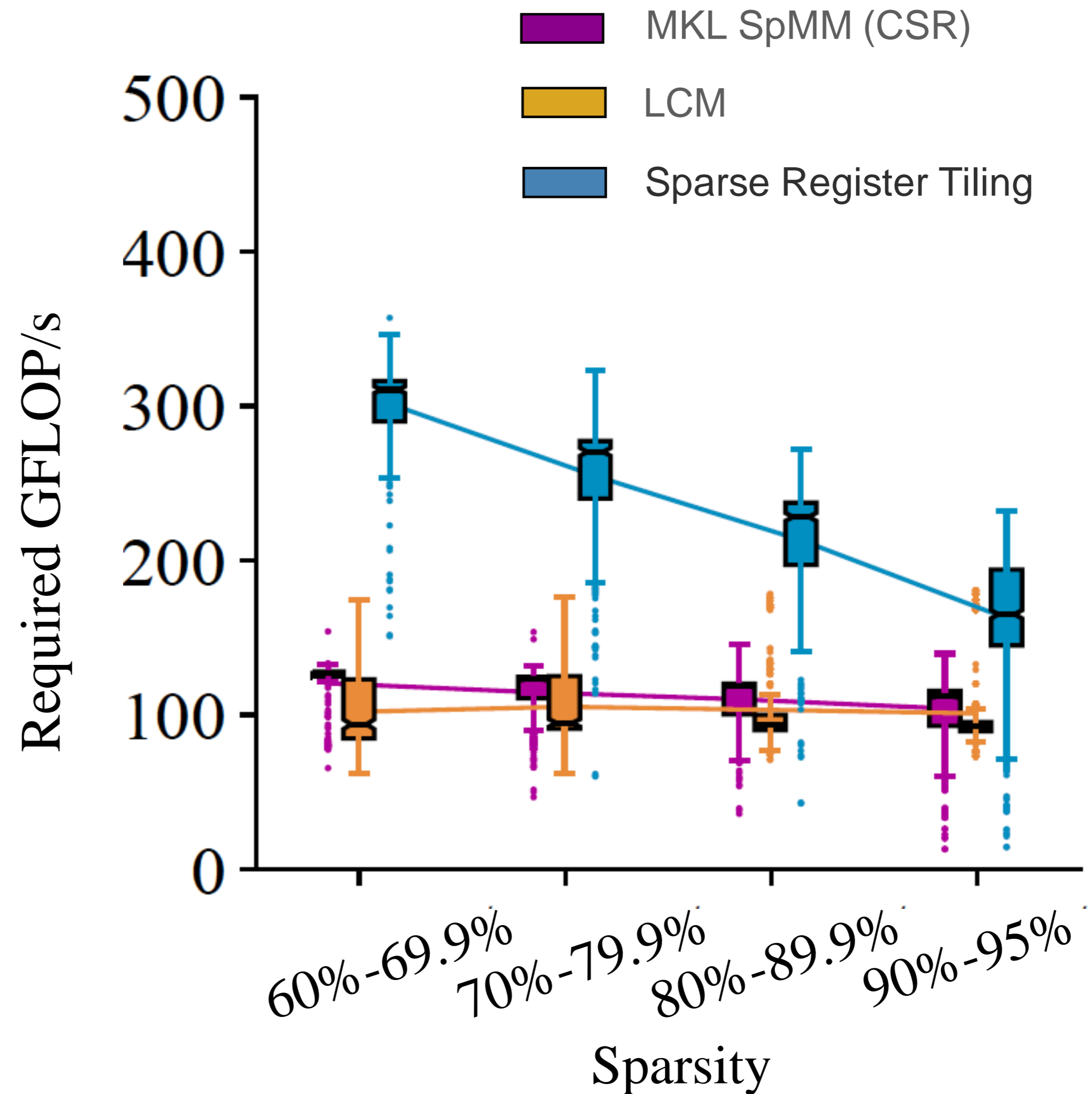


LCM does not improve the performance of *machine learning problems*.

¹Trevor Gale, Matei Zaharia, Cliff Young, and Erich Elsen. 2020. Sparse GPU Kernels for Deep Learning. SC '20.

LCM on Machine Learning Matrices

Testbed: Over 3000 matrices from Deep Learning Matrix Collection (DLMC¹).



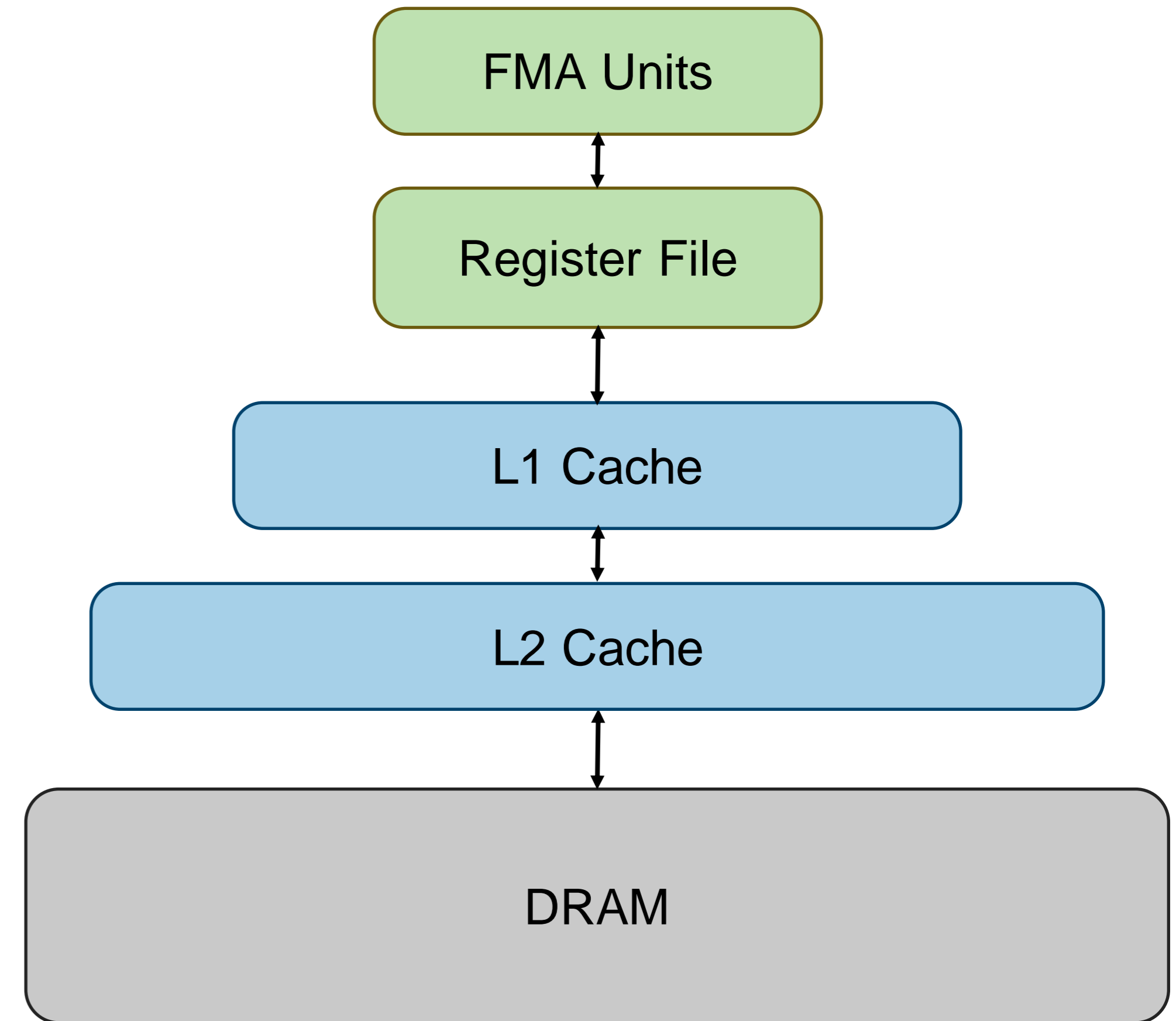
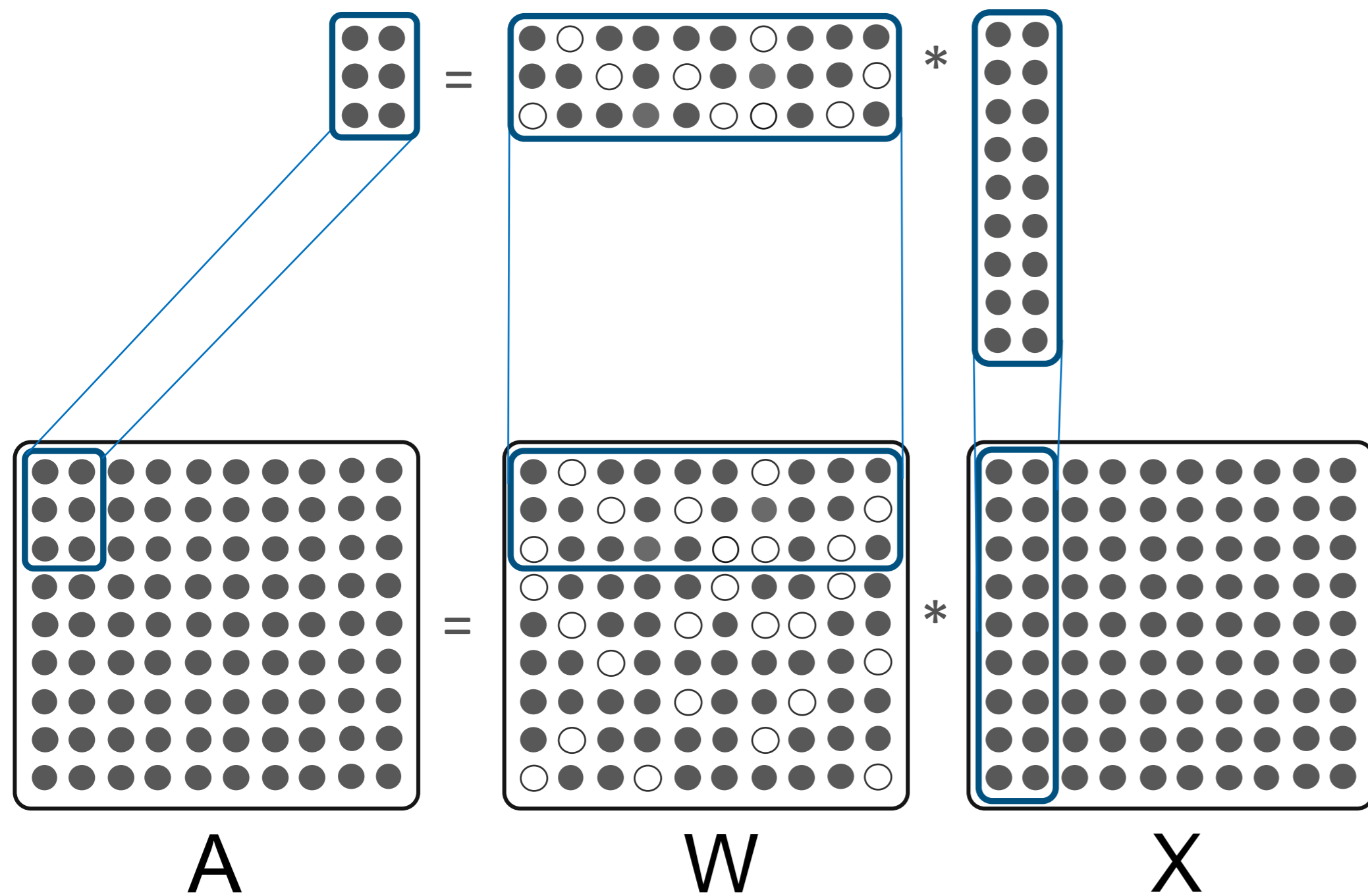
LCM does not improve the performance of *machine learning problems*.

¹Trevor Gale, Matei Zaharia, Cliff Young, and Erich Elsen. 2020. Sparse GPU Kernels for Deep Learning. SC '20.

Cache Tiling for Machine Learning

Because of their random sparsity patterns, in DNNs over 96% of mined patterns are **unstrided** in at least **two of the array accesses!**

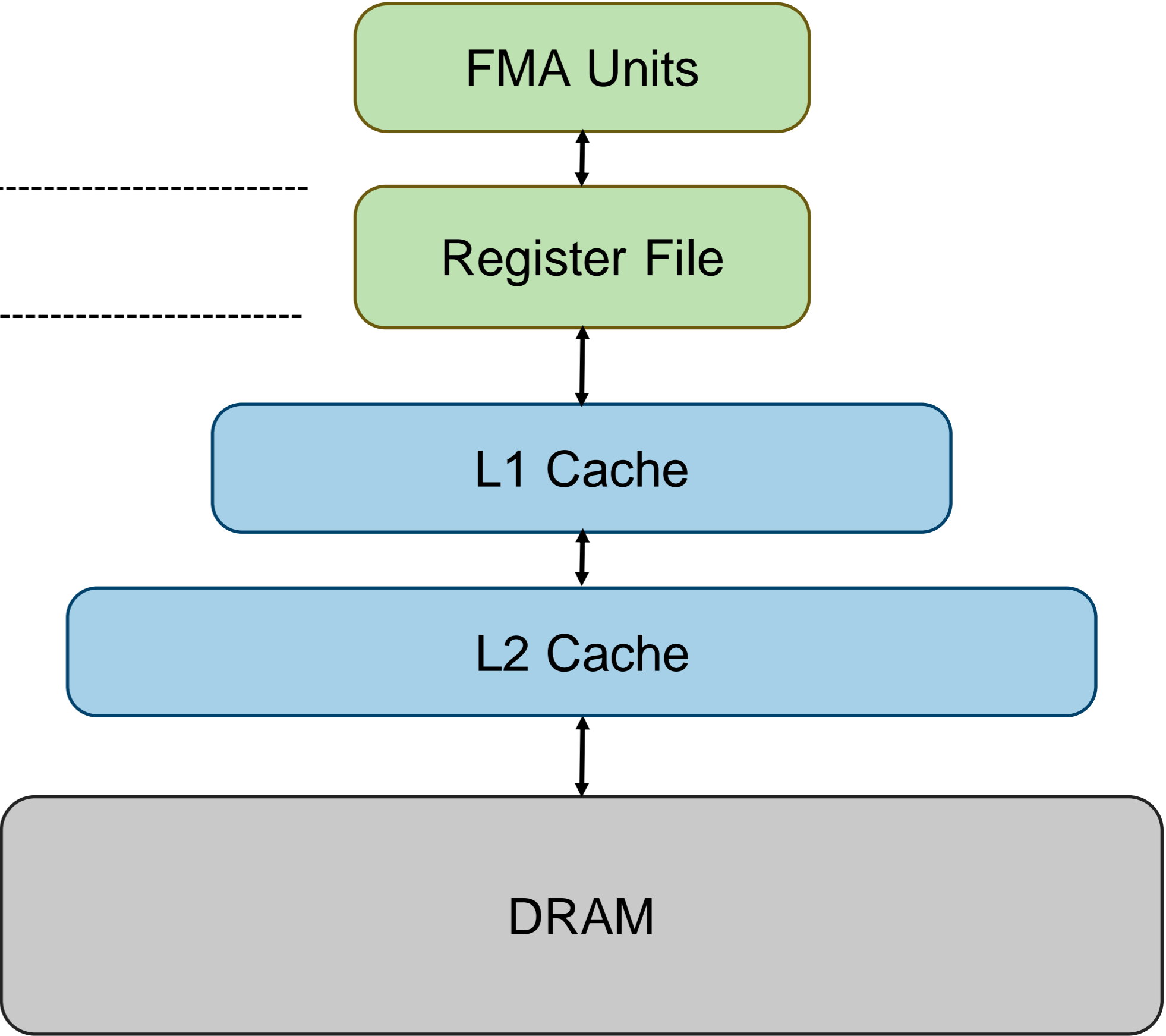
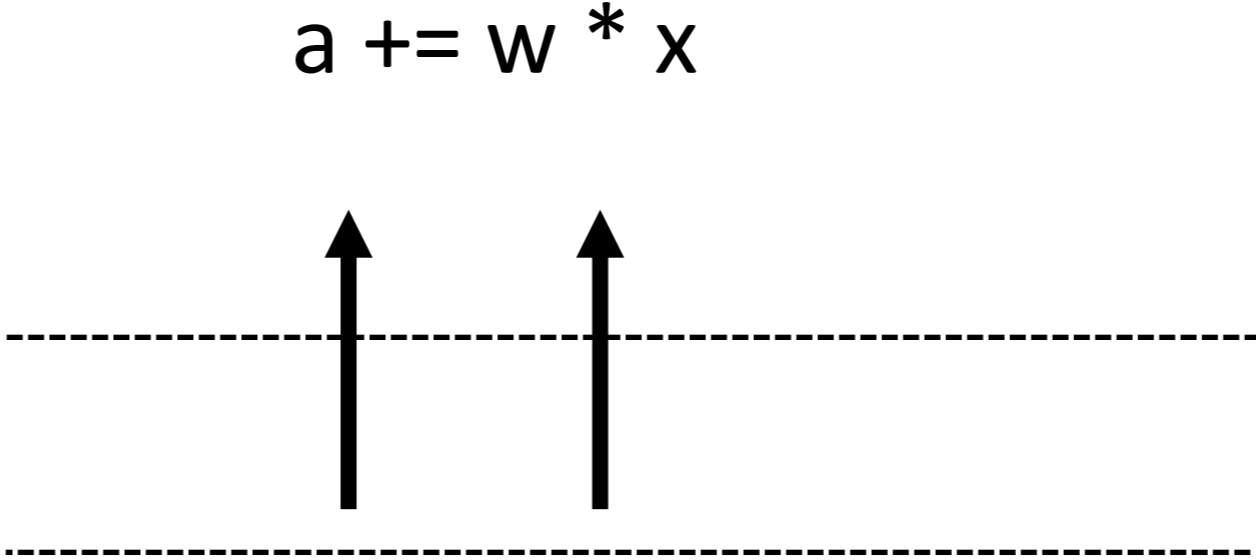
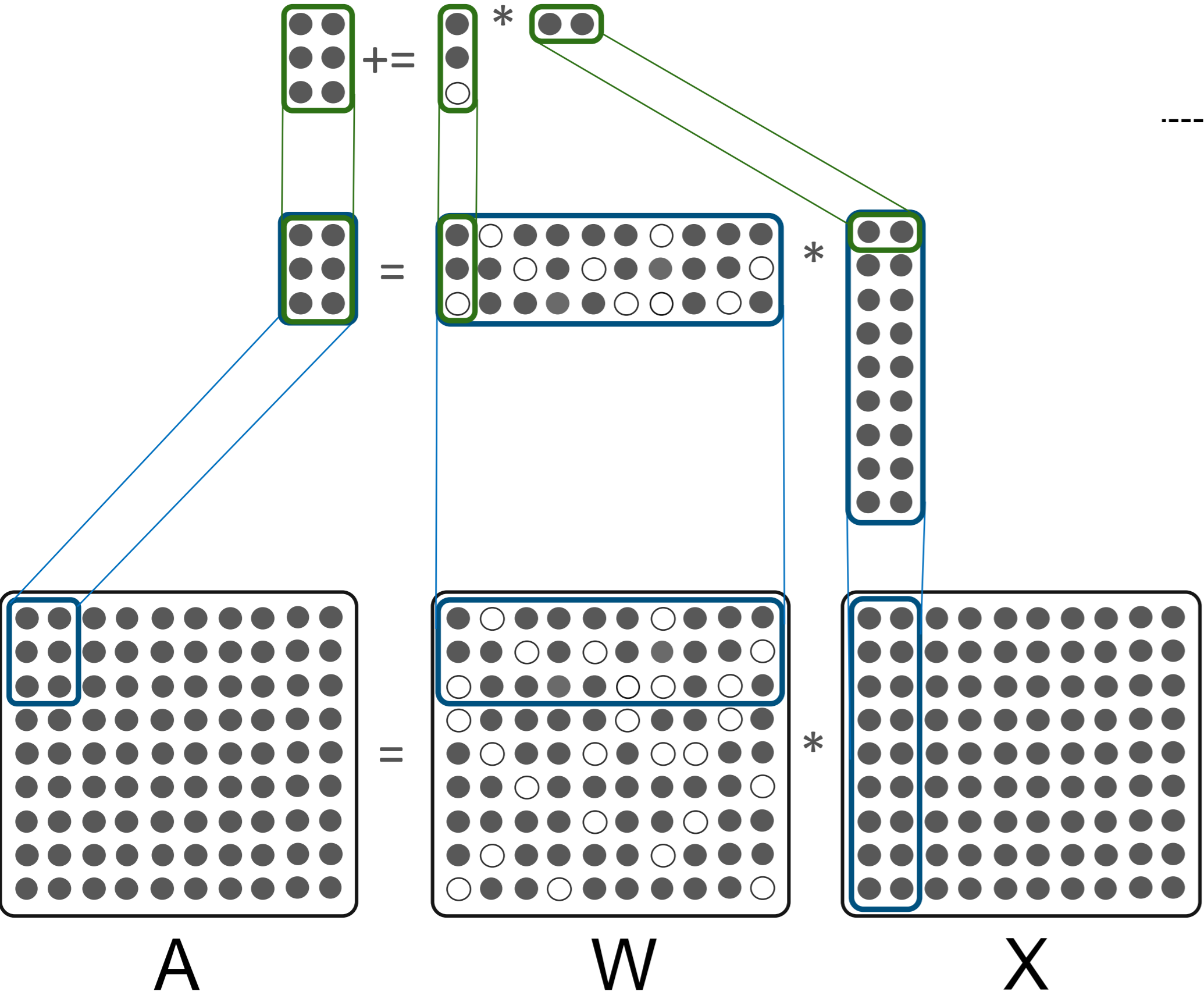
The better reordering method for unstructured matrices for caches is standard tiling (used in dense problems)*.



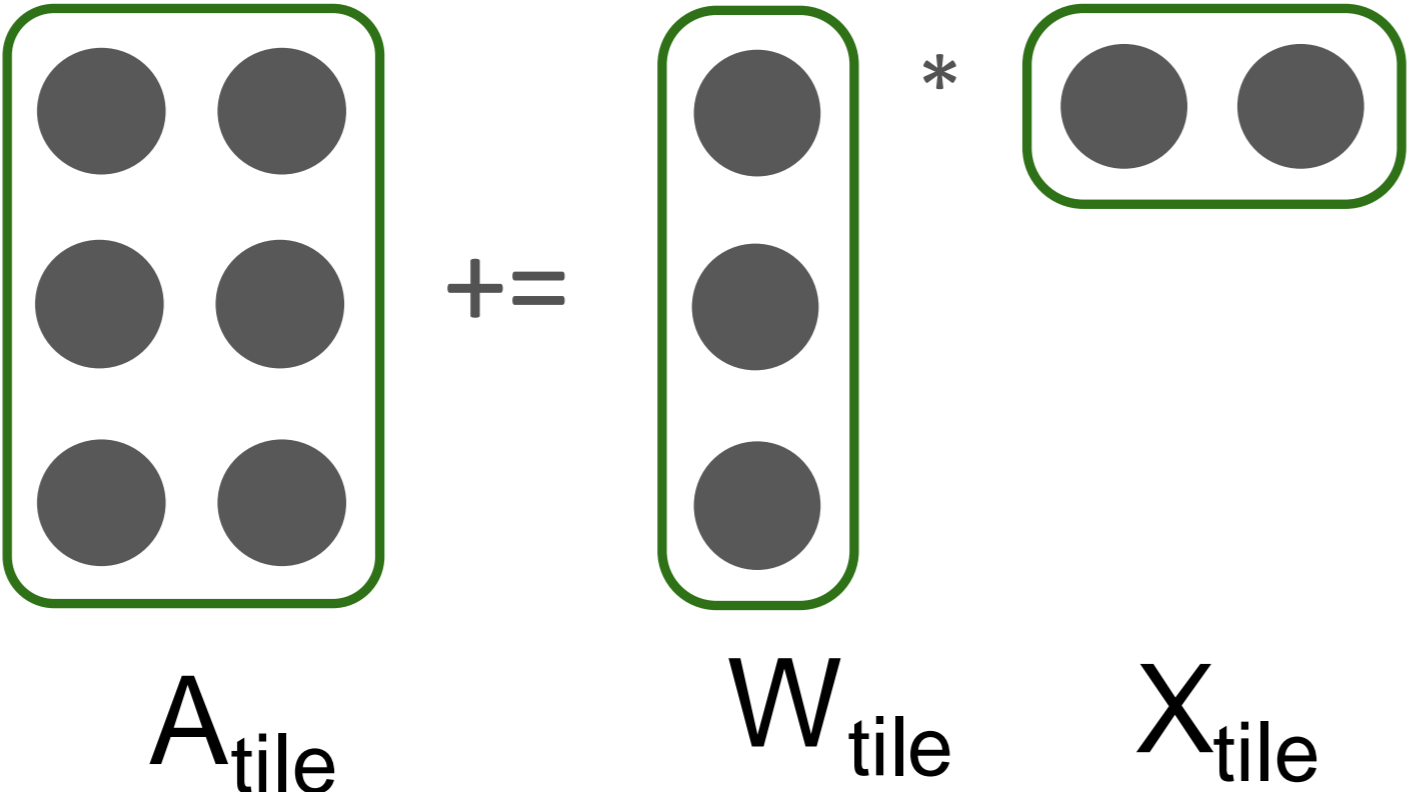
*Working set size discussions in section 3.2 of Register Tiling for Unstructured Sparsity in Neural Network Inference [PLDI'23]

Register Tiling

Improve locality in registers!



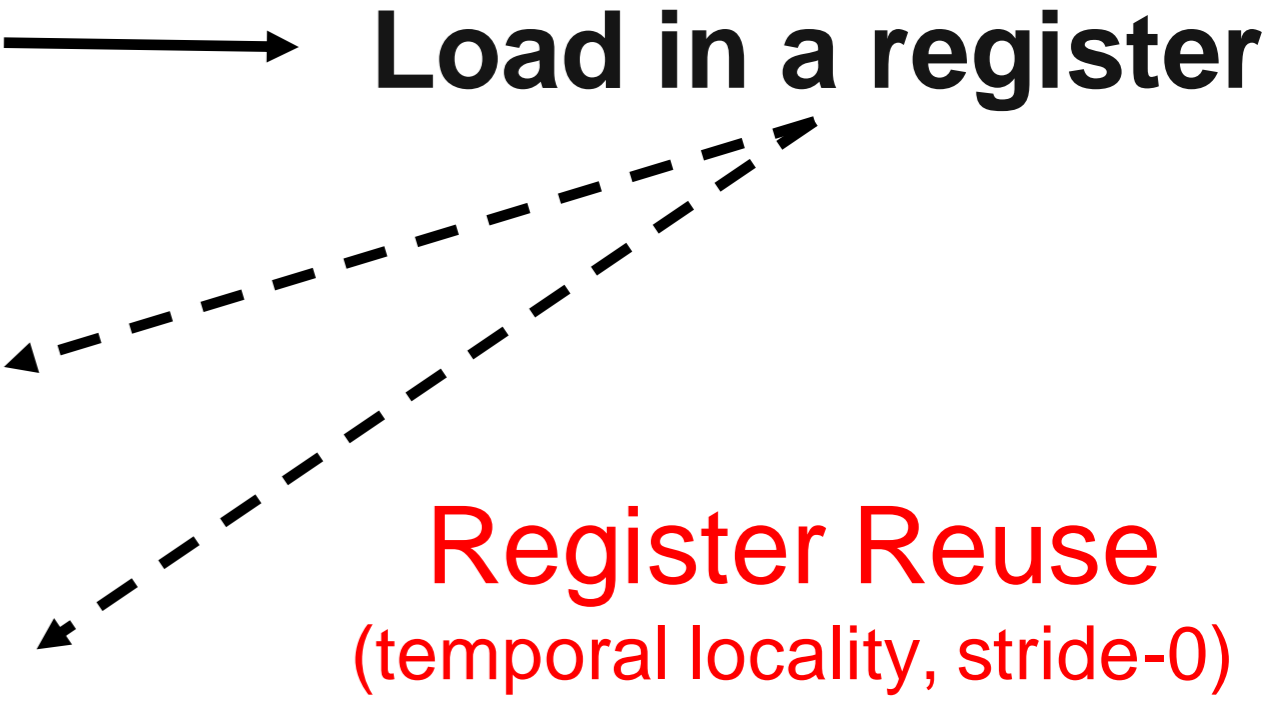
Register Reuse with Dense



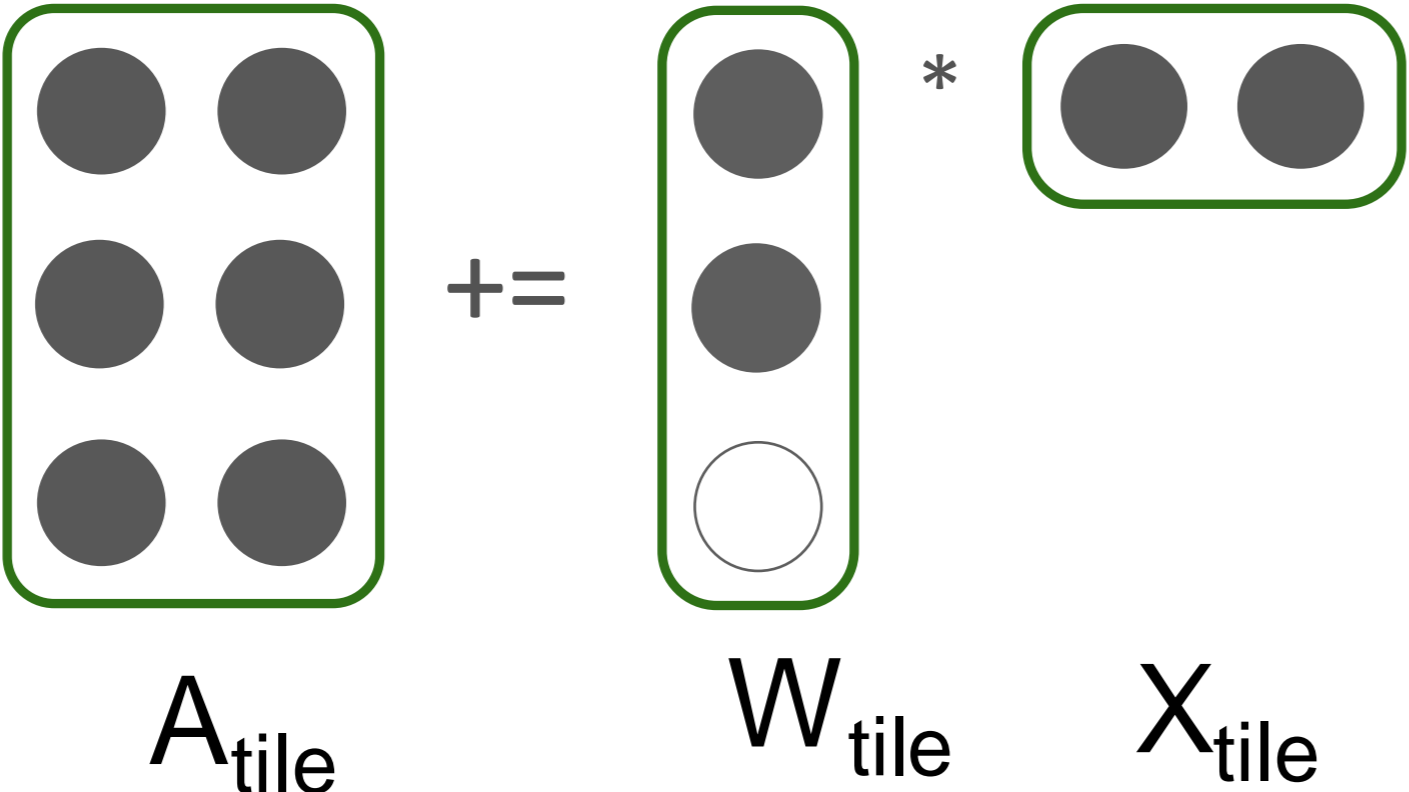
reg1	reg2
reg3	reg4
reg5	reg6
reg7	reg8
reg9	reg10
reg11	reg12

```

A_tile[0, 0] += W_tile[0, k] * X_tile[k, 0]
A_tile[0, 1] += W_tile[0, k] * X_tile[k, 1]
A_tile[1, 0] += W_tile[1, k] * X_tile[k, 0]
A_tile[1, 1] += W_tile[1, k] * X_tile[k, 1]
A_tile[2, 0] += W_tile[2, k] * X_tile[k, 0]
A_tile[2, 1] += W_tile[2, k] * X_tile[k, 1]
    
```



Register Reuse with Sparse



reg1	reg2
reg3	reg4
reg5	reg6
reg7	reg8
reg9	reg10
reg11	reg12

```

if (W_tile[0,k]) {
if (W_tile[1,k]) {
if (W_tile[2,k]) {

```

```

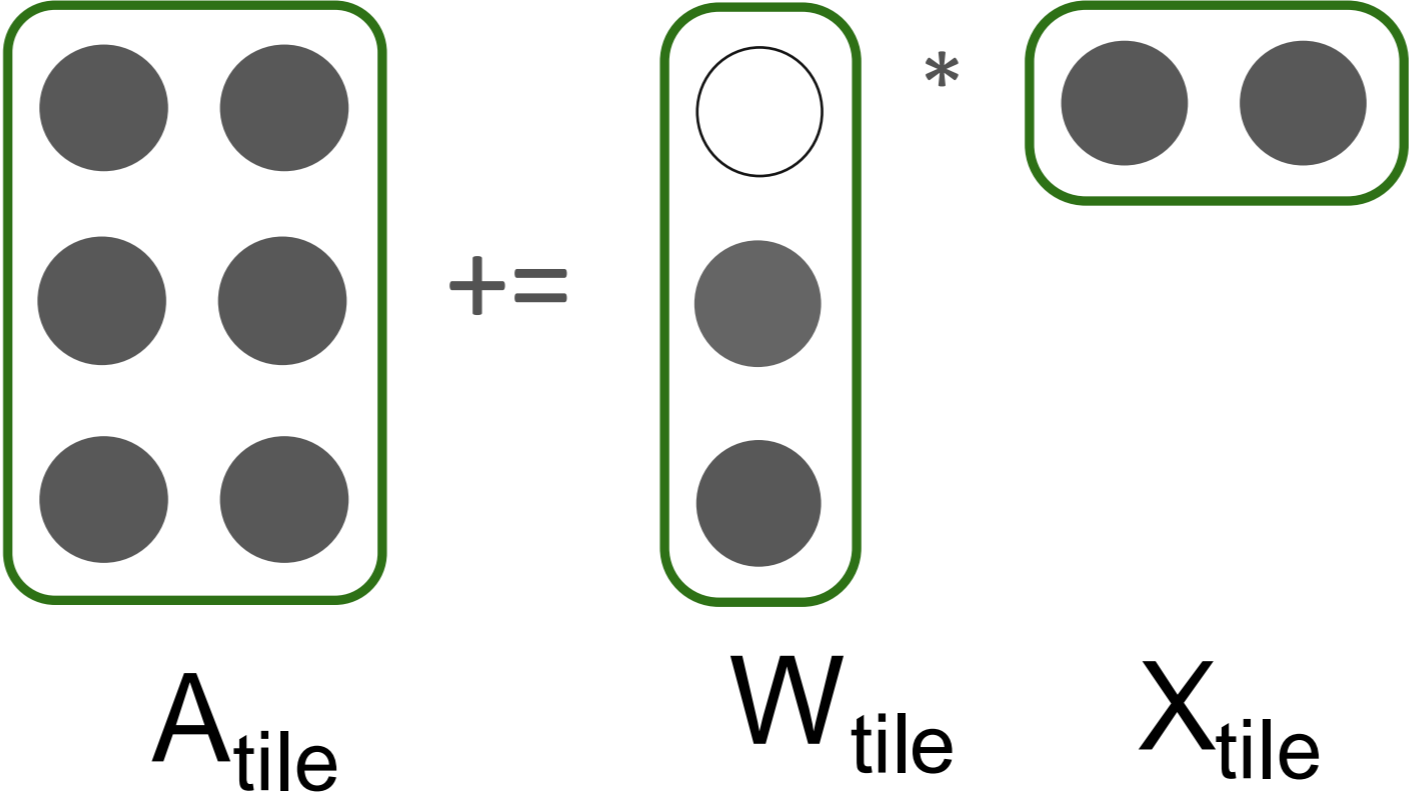
A_tile[0, 0] += W_tile[0, k] * X_tile[k, 0]
A_tile[0, 1] += W_tile[0, k] * X_tile[k, 1]
A_tile[1, 0] += W_tile[1, k] * X_tile[k, 0]
A_tile[1, 1] += W_tile[1, k] * X_tile[k, 1]
A_tile[2, 0] += W_tile[2, k] * X_tile[k, 0]
A_tile[2, 1] += W_tile[2, k] * X_tile[k, 1]

```

Unclear which instructions will execute: **no register reuse**

Overhead of if-conditions

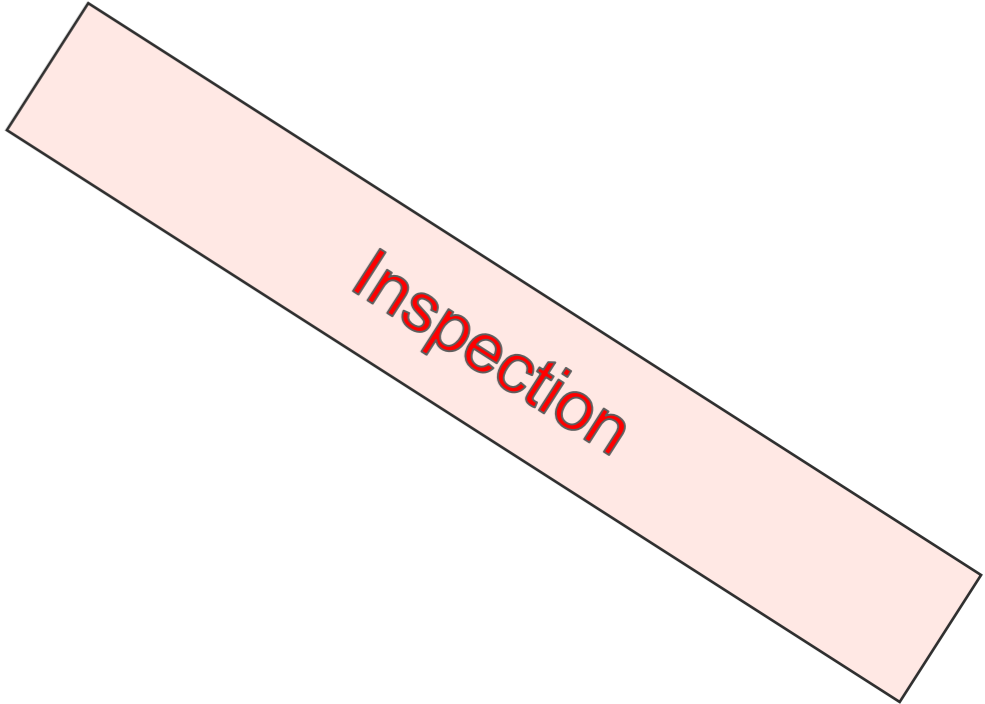
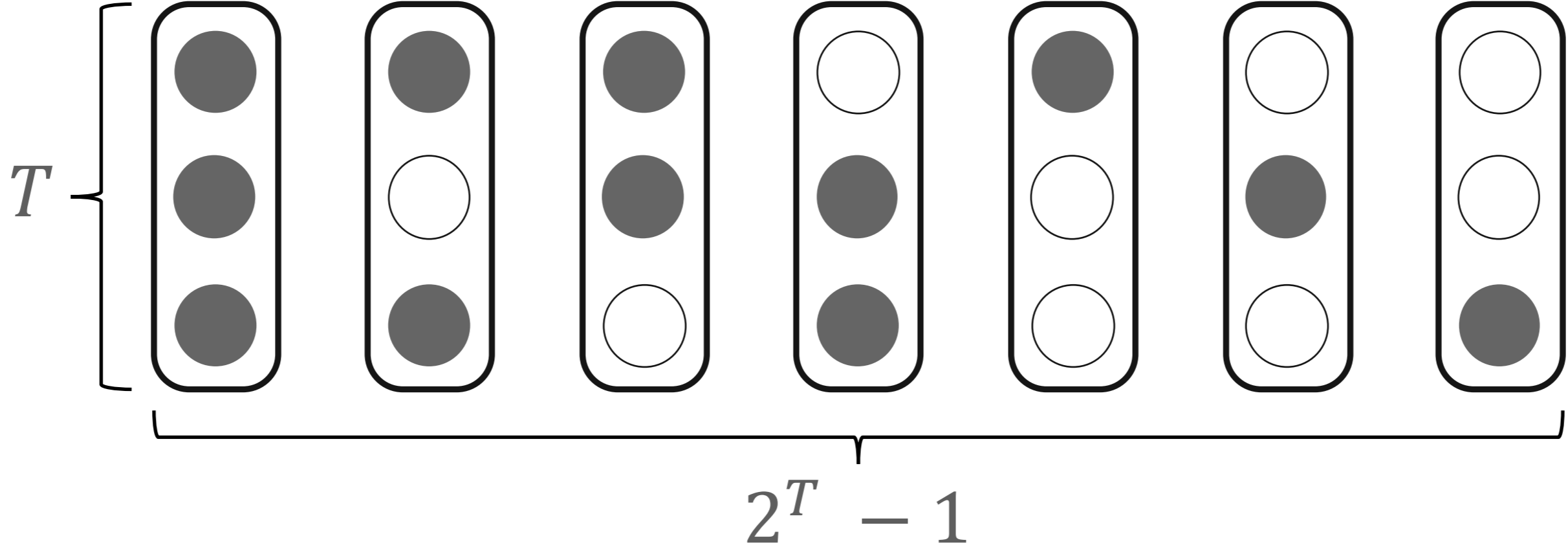
If Tile Sparsity was Known



$A_{\text{tile}}[0, 0] += W_{\text{tile}}[0, k] * X_{\text{tile}}[k, 0]$
 $A_{\text{tile}}[0, 1] += W_{\text{tile}}[0, k] * X_{\text{tile}}[k, 1]$
 $A_{\text{tile}}[1, 0] += W_{\text{tile}}[1, k] * X_{\text{tile}}[k, 0]$
 $A_{\text{tile}}[1, 1] += W_{\text{tile}}[1, k] * X_{\text{tile}}[k, 1]$
 $A_{\text{tile}}[2, 0] += W_{\text{tile}}[2, k] * X_{\text{tile}}[k, 0]$
 $A_{\text{tile}}[2, 1] += W_{\text{tile}}[2, k] * X_{\text{tile}}[k, 1]$

Register reuse will happen because common access patterns are known.

Generate Code for Each Possible Tile Pattern

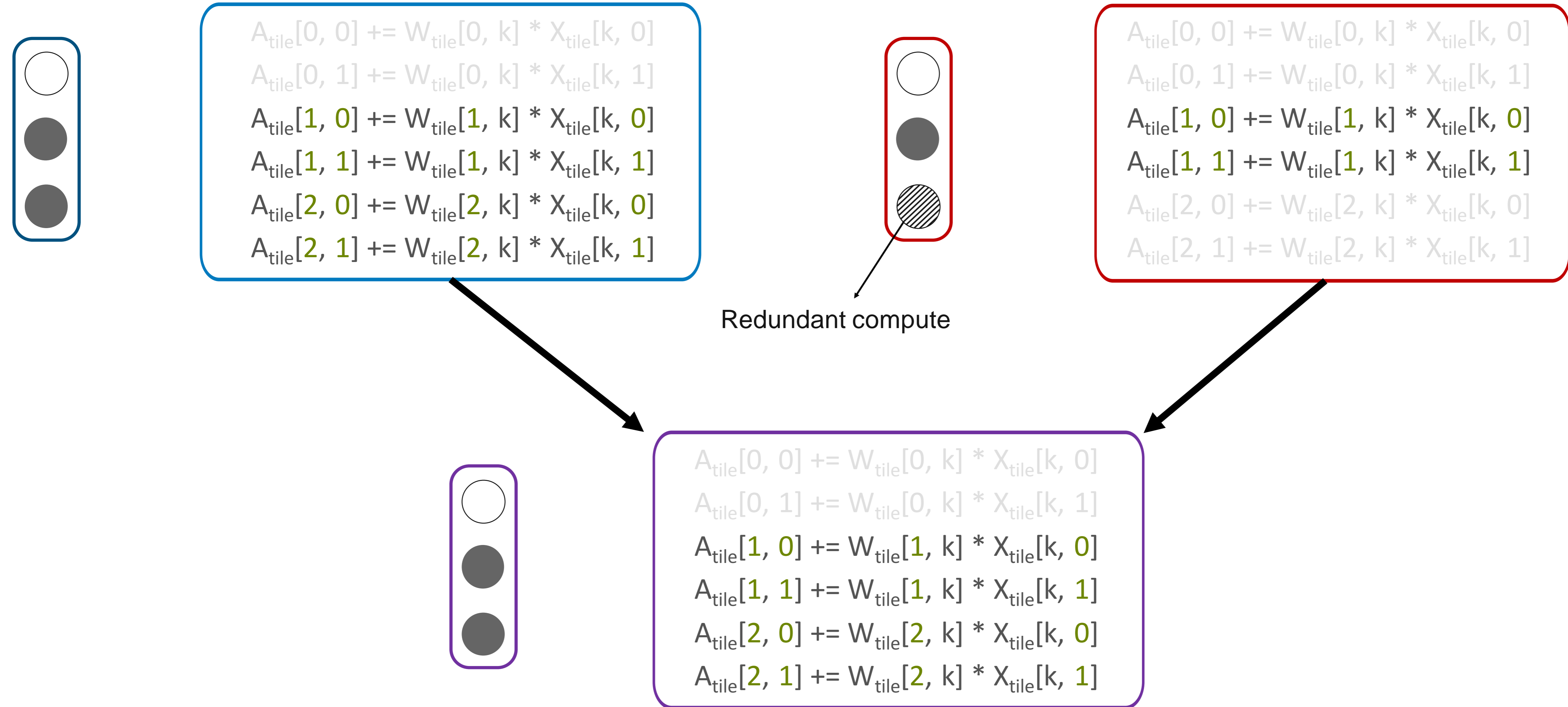


⚠ **Code bloat**

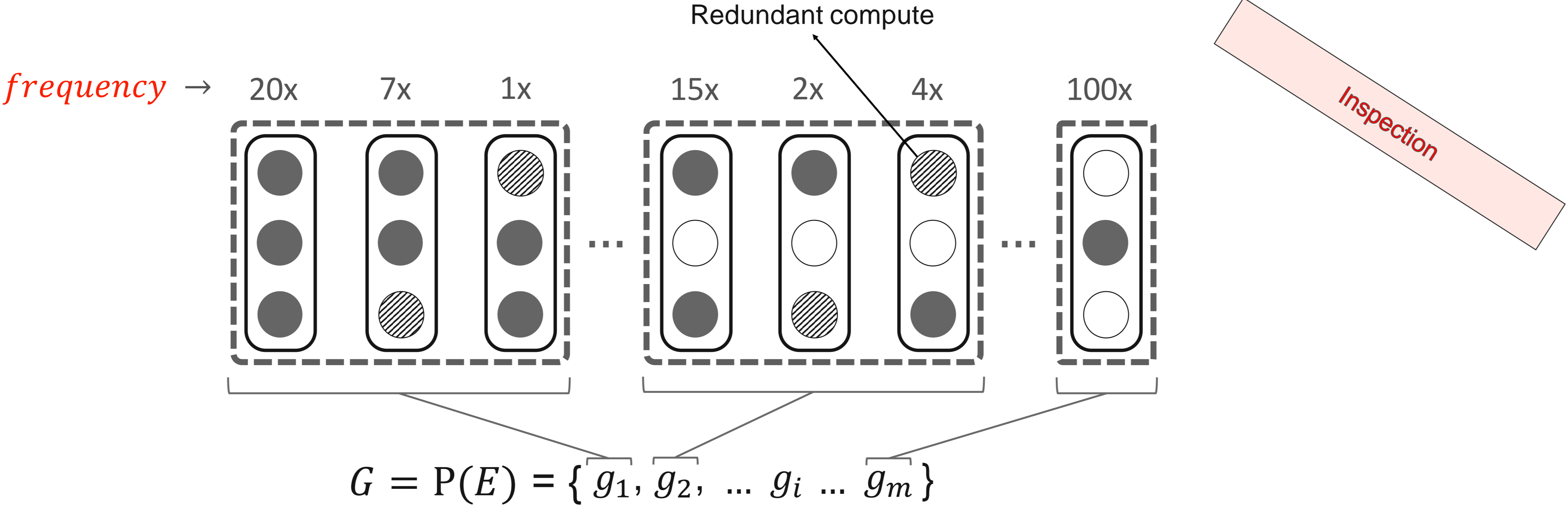
Solution: Sparse Jam Solver



The Jamming Problem



Sparse-Jam Solver



$$\min \sum_{i=1}^{|G|} x_i * (freq(g_i) * exec_cost(g_i) + overhead) \quad x_i \in \{0, 1\}$$

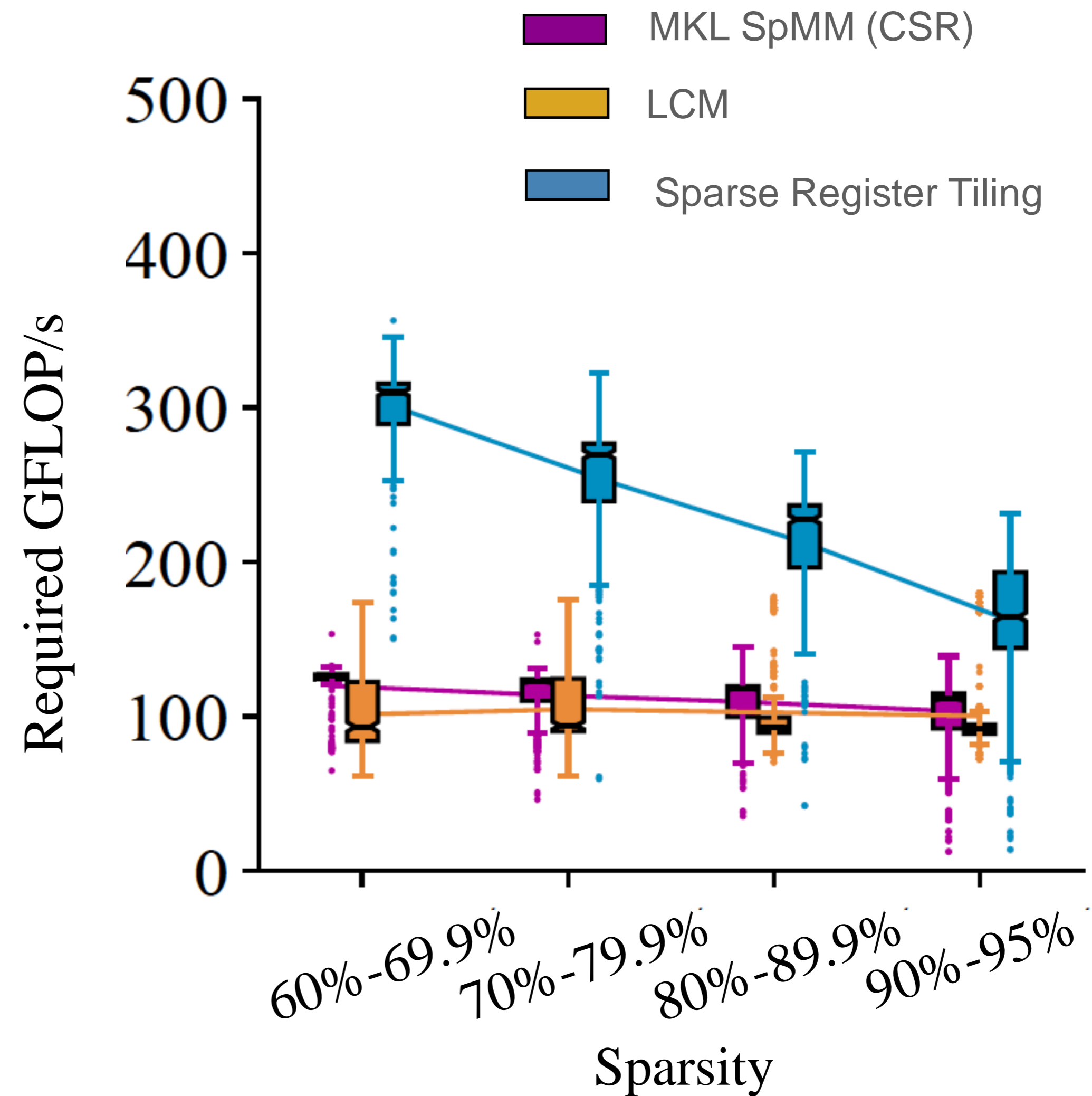
Constraint 1: Ensure all operations are covered

Constraint 2: Size of generated code \leq *CodeSizeLimit*

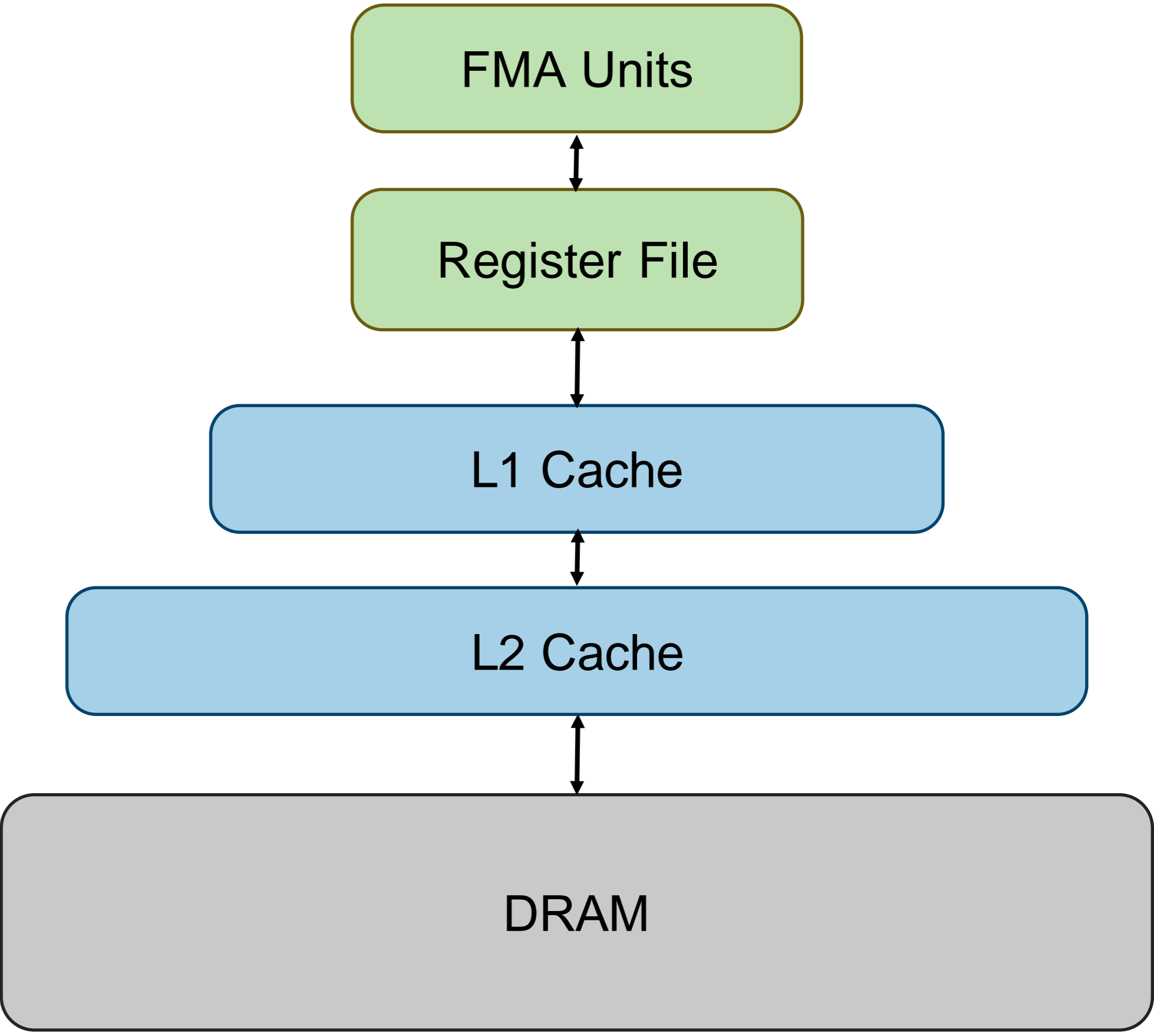
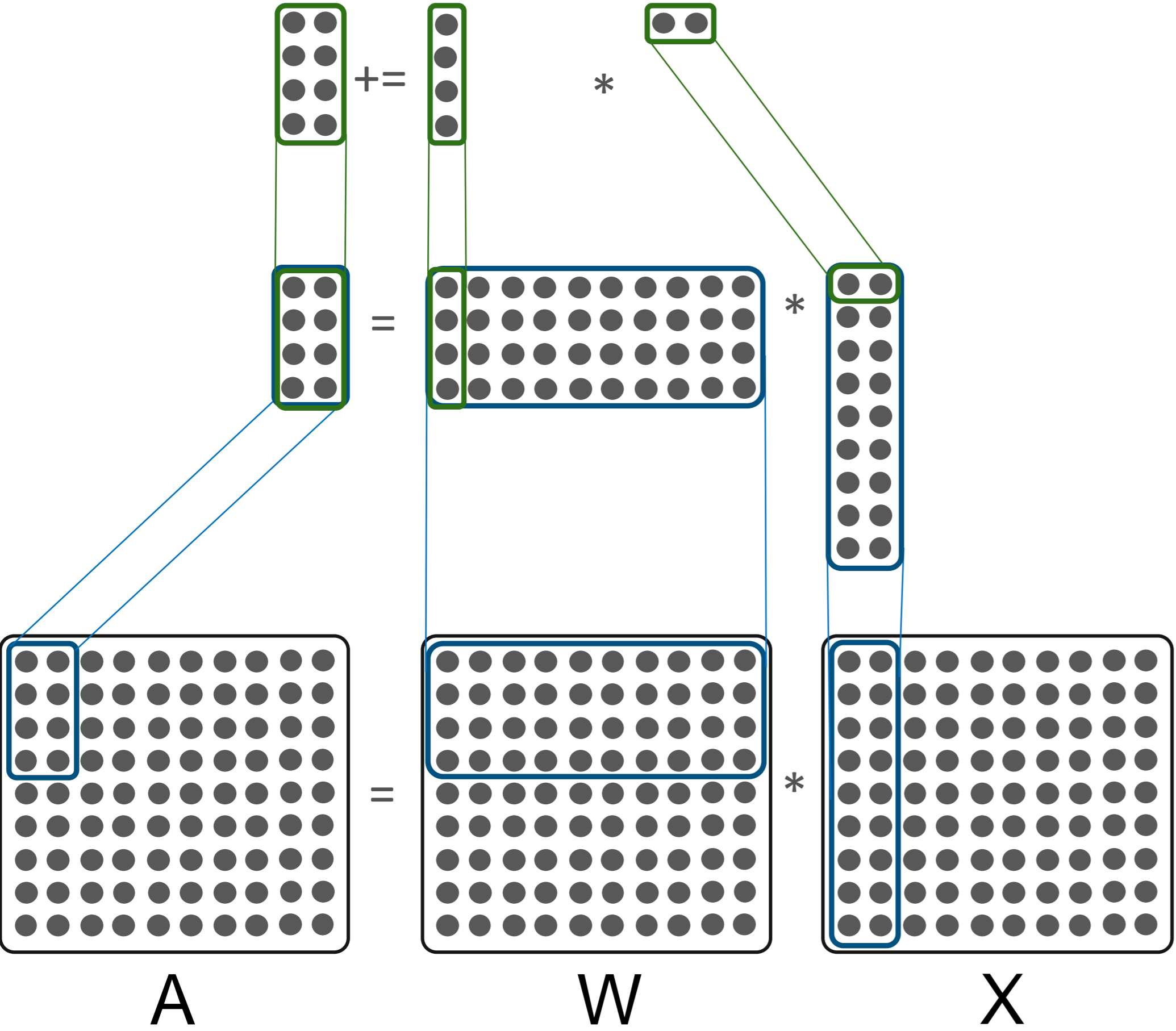
Sparse Register Tiling for Machine Learning

Testbed: **Deep Learning Matrix Collection (DLMC)**¹, a total of **2396 matrices**. A 20 core Cascadelake Xeon(R) Gold 6248 CPU with AVX512.

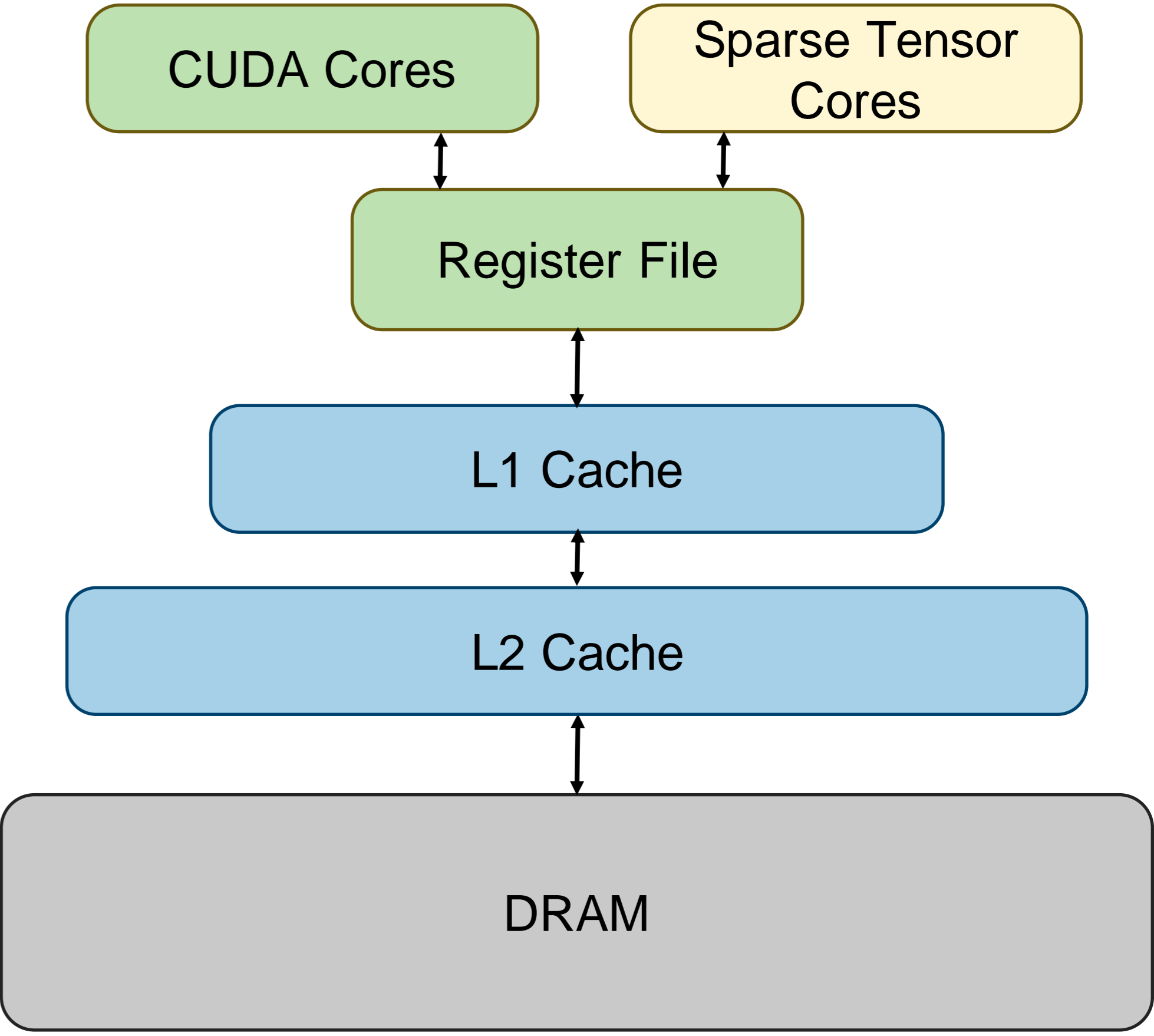
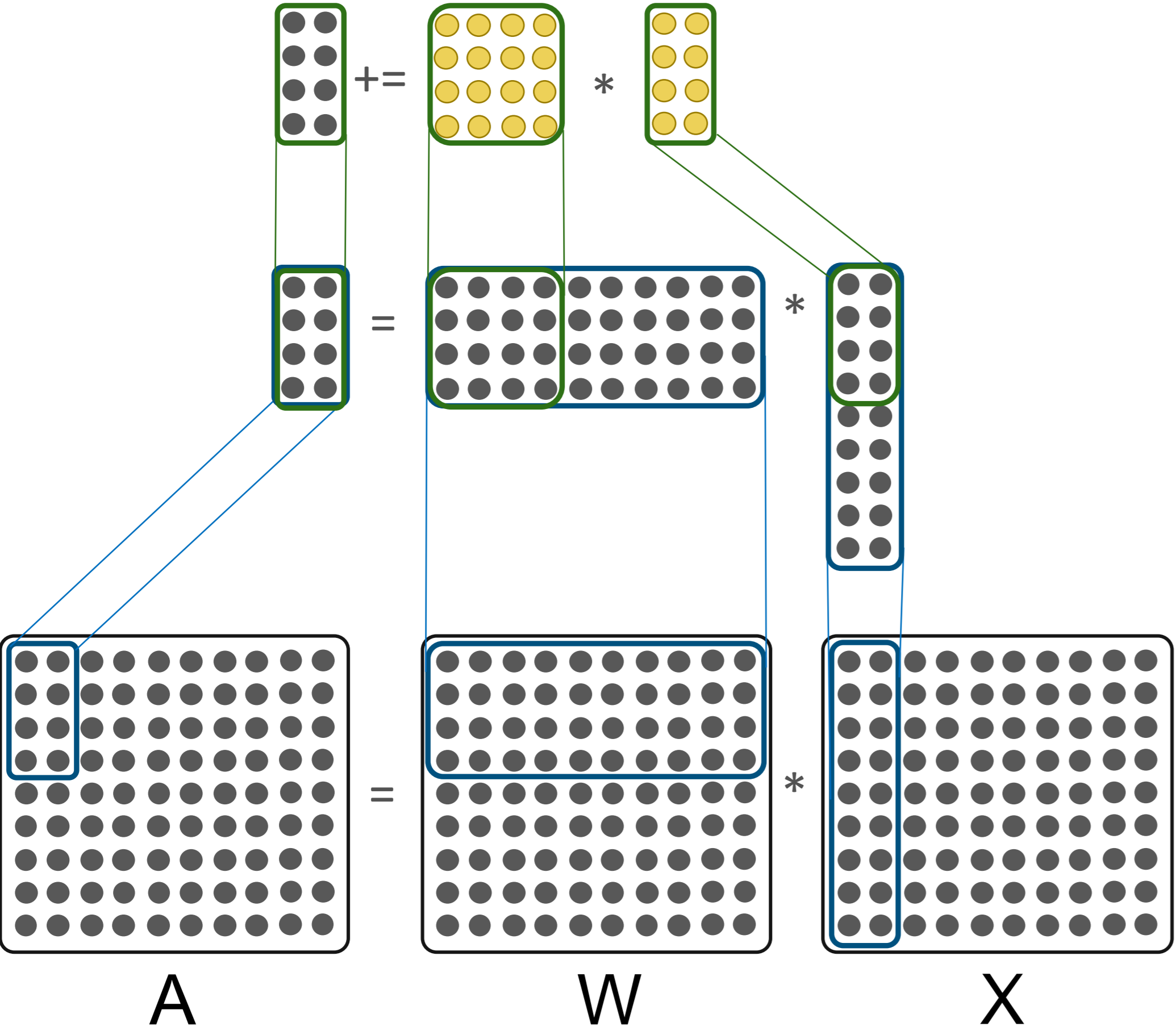
Sparse Register Tiling provides **geomean speedup** of **2.65x**, **1.72x** and **3.23x** over MKL SGEMM, MKL SpMM (CSR) and ASpT² respectively.



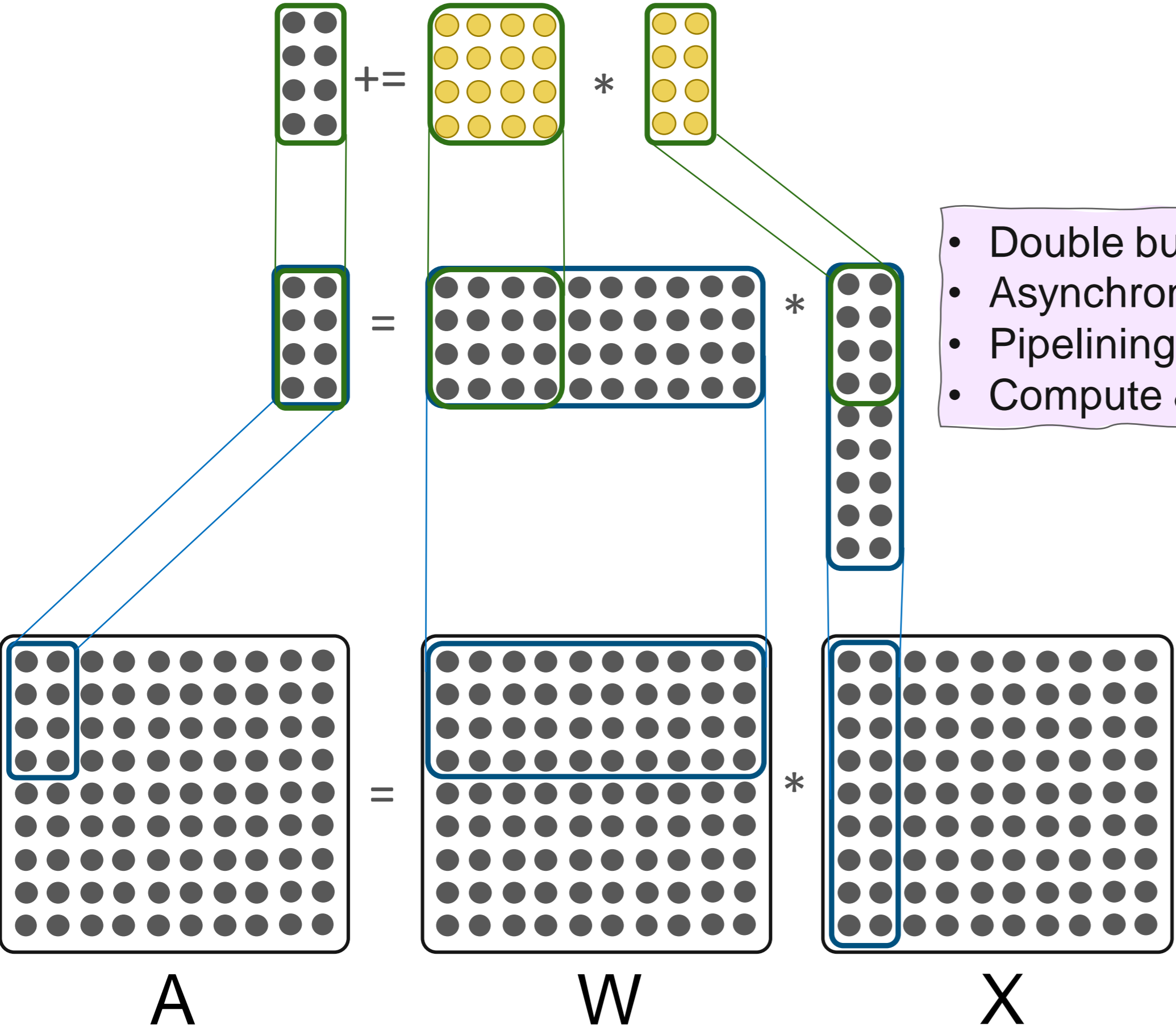
Register Tiling on CPU



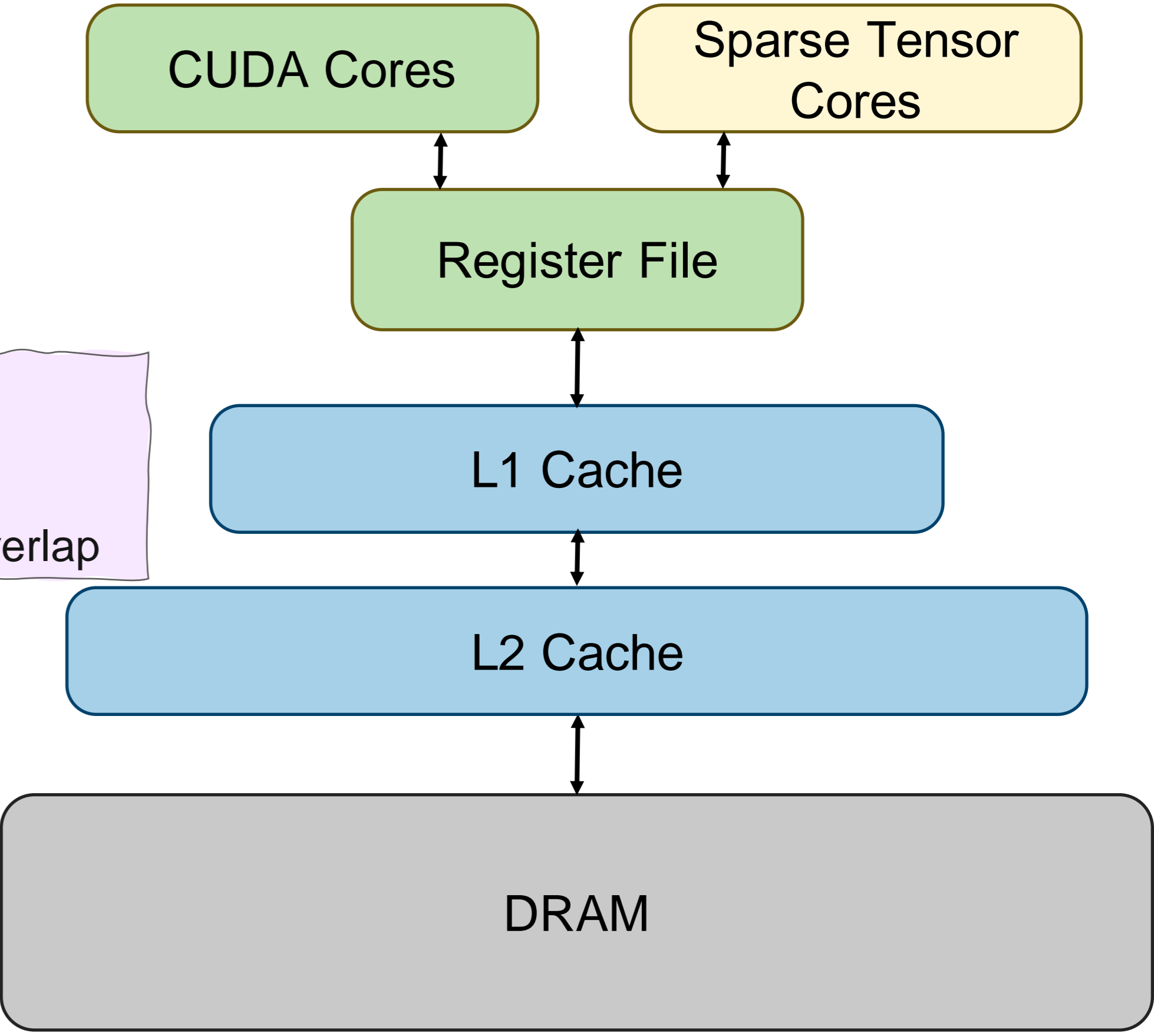
Register Tiling on GPUs



GPU Optimizations Needed to get to Registers

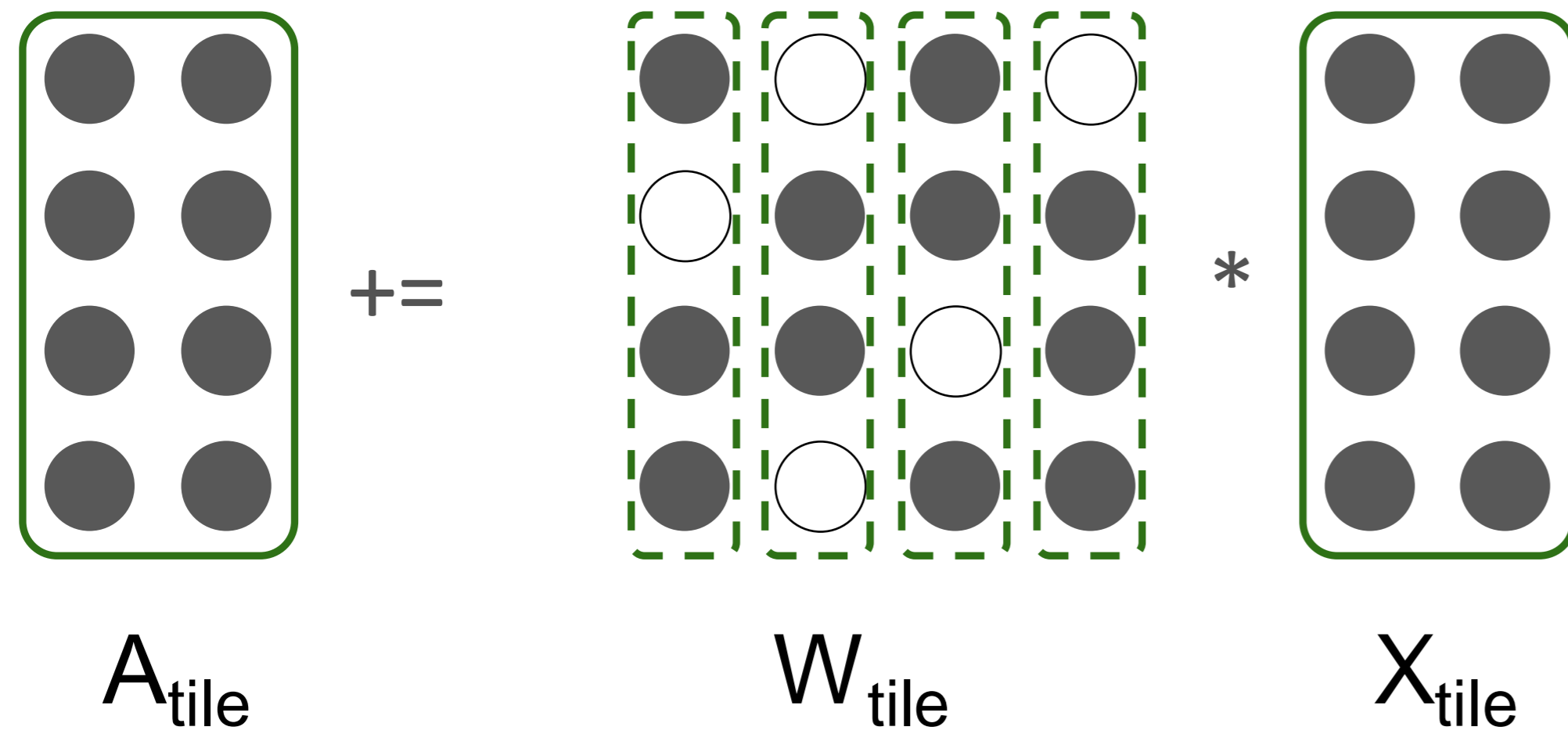


- Double buffering
- Asynchronous mem copies
- Pipelining
- Compute & communication overlap



Sparse Register Tiling on GPUs?

Vertical tiles: register reuse on X



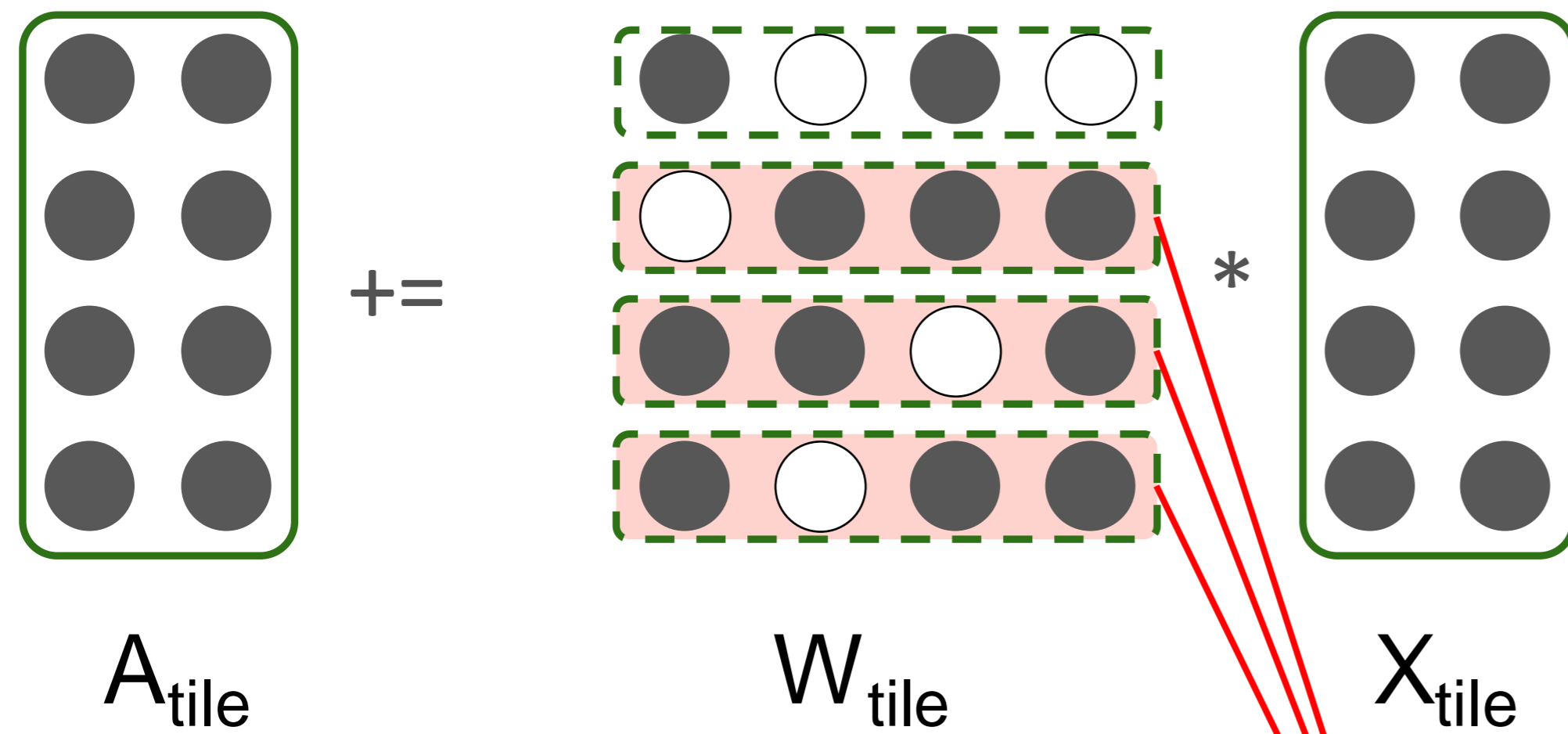
CUDA Cores

~~Sparse Tensor
Cores~~

```
__device__ matmul{  
  A_tile[0, 0] += W_tile2[0, k] * X_tile[k, 0]  
  A_tile[0, 1] += W_tile2[0, k] * X_tile[k, 1]  
  A_tile[2, 0] += W_tile2[2, k] * X_tile[k, 0]  
  A_tile[2, 1] += W_tile2[2, k] * X_tile[k, 1]  
  ...  
}
```

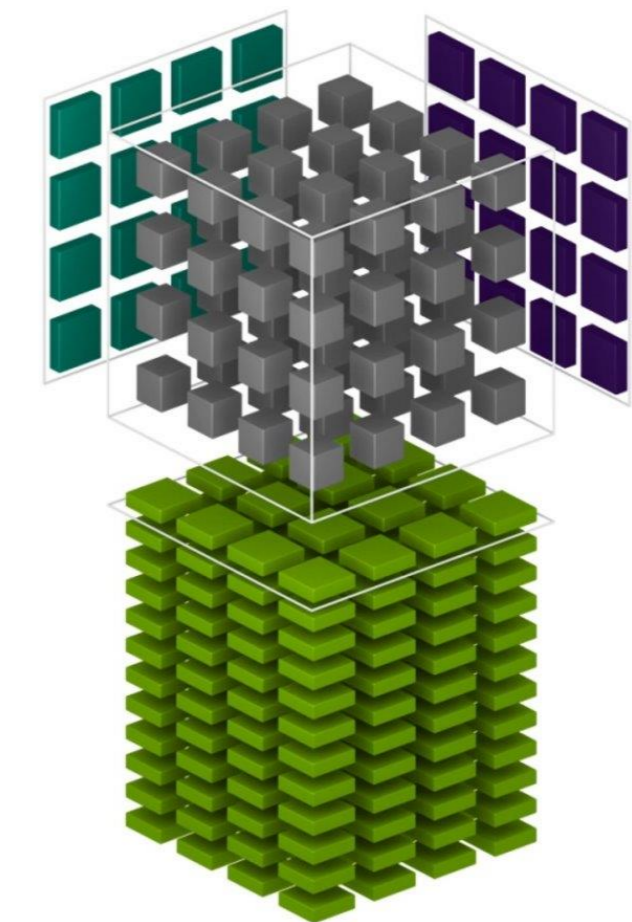
What About the Sparse Tensor Cores?

2:4 Sparsity: At most 2 elements out of every partition of 4 **horizontally** consecutive elements should be non-zeros



❌ Sparse Tensor Cores

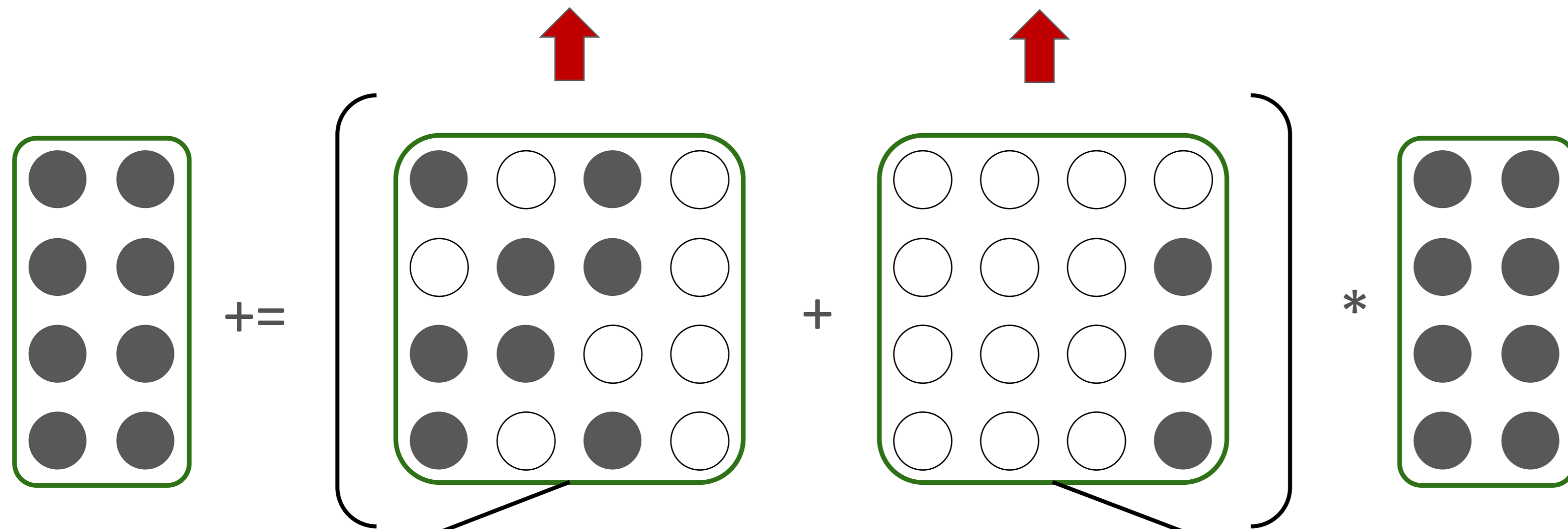
Sparse Tensor Cores



GPU Tile Decomposition

Sparse Tensor Cores

CUDA Cores



```
__device__ void matmul_structure_sparse(*)
{
    asm volatile("mma.sp.sync.aligned.fp16",
        "Wtile1", "Xtile", "Atile")
}
```

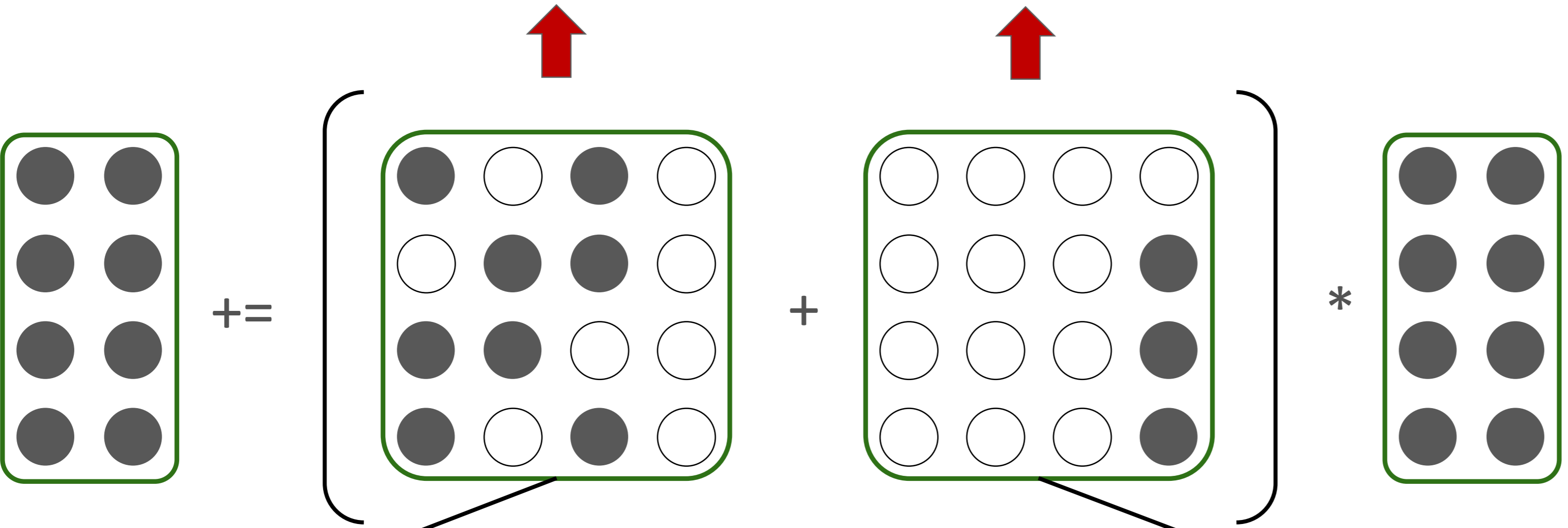
```
__device__ matmul_finegrained(*){
    A_tile[1, 0] += W_tile2[1, k] * X_tile[k, 0]
    A_tile[1, 1] += W_tile2[1, k] * X_tile[k, 1]
    ...
}
```

Tile Decomposition for Unstructured Sparse Model Inference, Liu, Shahsavan, Dehnavi [in-review]

GPU Tile Decomposition & Tile Scheduling

Sparse Tensor Cores

CUDA Cores



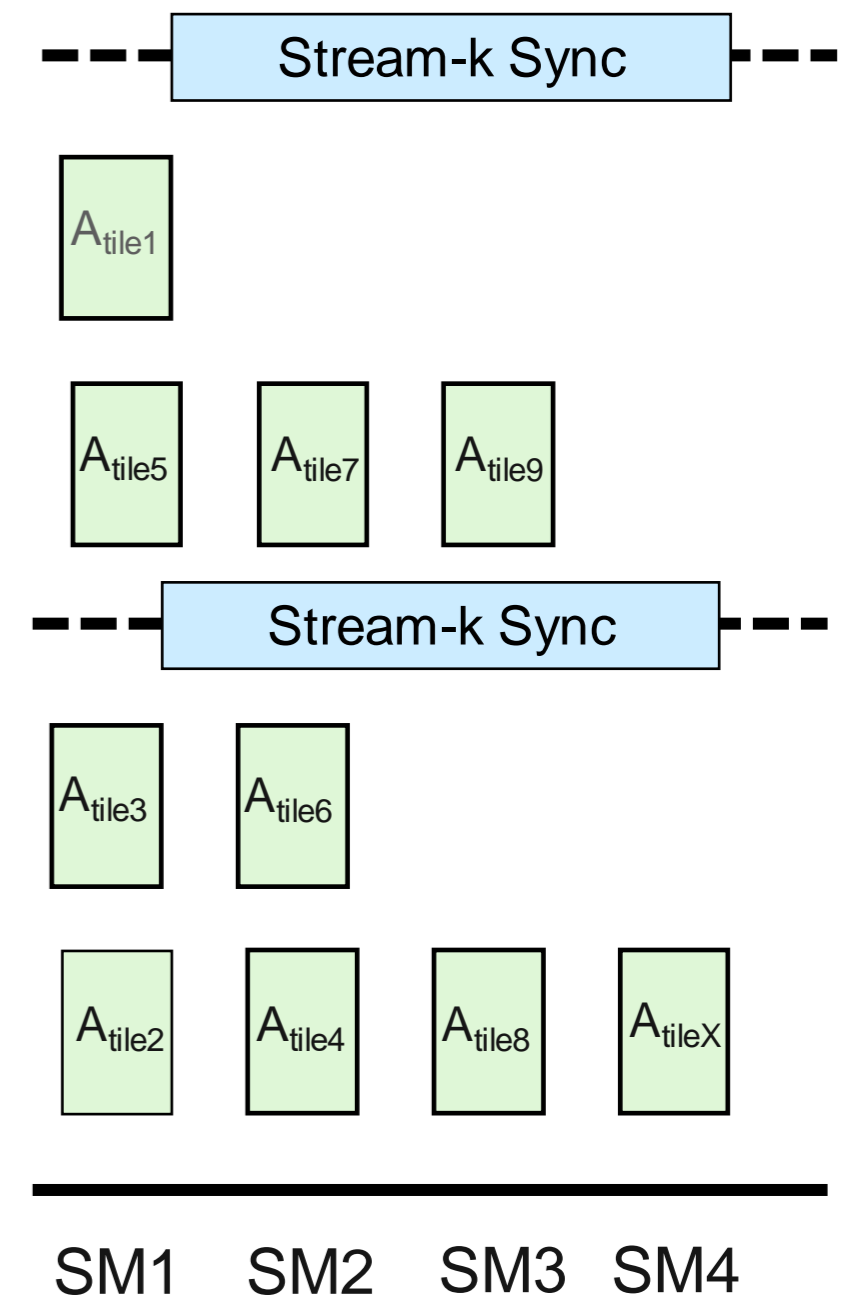
```

__device__ void matmul_structure_sparse(*)
{
    asm volatile("mma.sp.sync.aligned.fp16",
                "Wtile1", "Xtile", "Atile")
}
    
```

```

__device__ matmul_finegrained(*){
    A_tile[1, 0] += W_tile2[1, k] * X_tile[k, 0]
    A_tile[1, 1] += W_tile2[1, k] * X_tile[k, 1]
    ...
}
    
```

Tile Scheduling

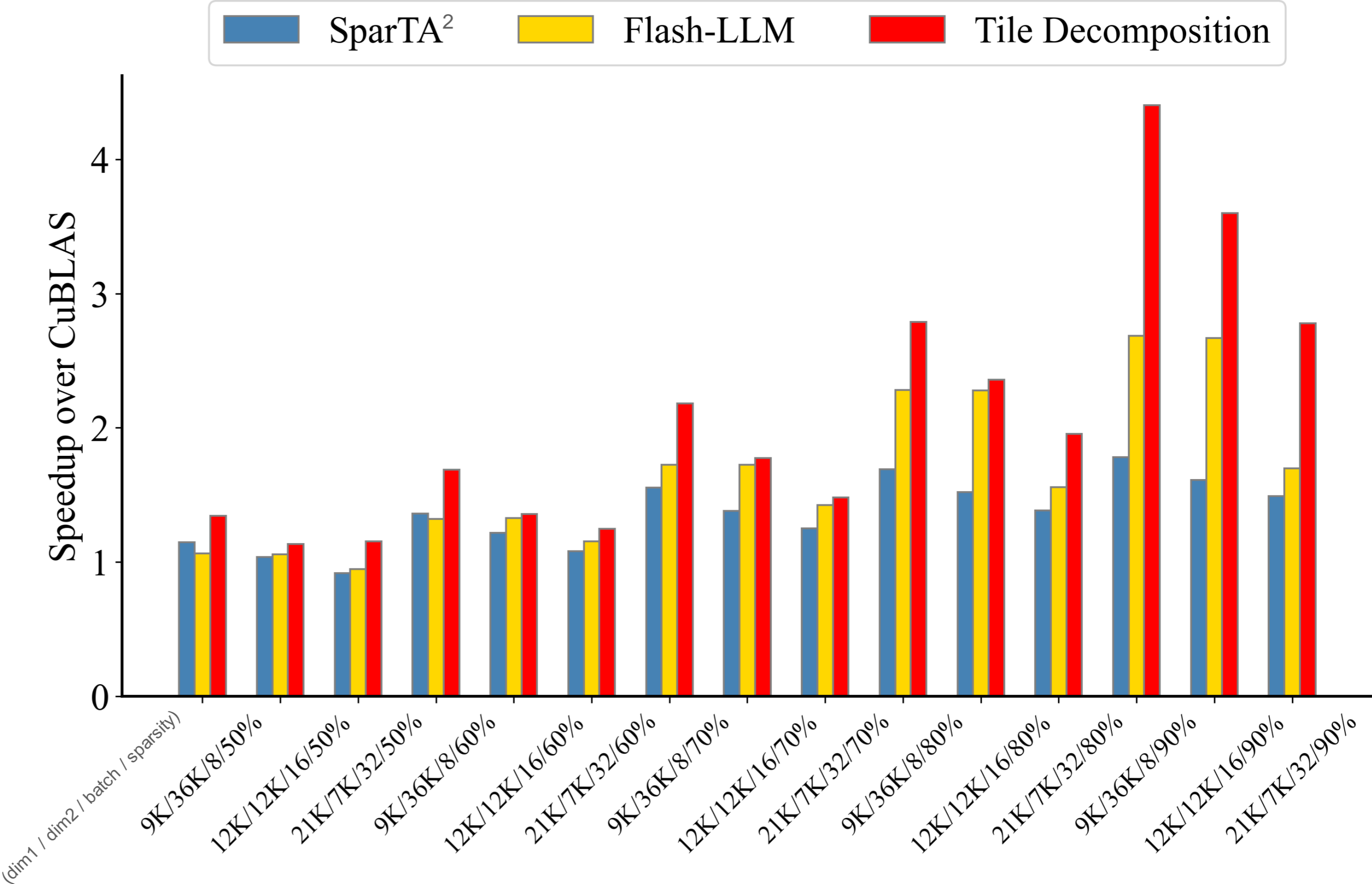


Tile Decomposition for Unstructured Sparse Model Inference, Liu, Shahsavan, Dehnavi [in-review]
 A Framework for Fine-Grained Synchronization of Dependent GPU Kernels, Jangda, et. al. [CGO'24]

Sparse Tile Decomposition on GPUs

Testbed: **GPT Matrices from Flash-LLM**¹,
Nvidia RTX 3080 TI, NVCC 12.2

Tile Decomposition provides **geomean speedup of 2.08x, 1.47x and 1.23x** over CuBLAS (Tensor Cores enabled), SparTA², and Flash-LLM, respectively, for FP16 data type.



¹Flash-llm: Enabling cost-effective and highly-efficient large generative model inference with unstructured sparsity, Xia et. al. VLBD'23

²SparTA: Deep-Learning Model Sparsity via Tensor-with-Sparsity-Attribute, Zheng et. al. OSDI'22

Sparsity Patterns in Machine Learning vs. Engineering

```
L, ... ← factorization(W)  
while (not converged) {
```

...

```
Solve:  $Lx_{k+1} = b$ 
```

```
 $r_{k+1} = r_k - \alpha Wx_{k+1}$ 
```

....

```
}
```

Engineering

PART 2

Solver routines

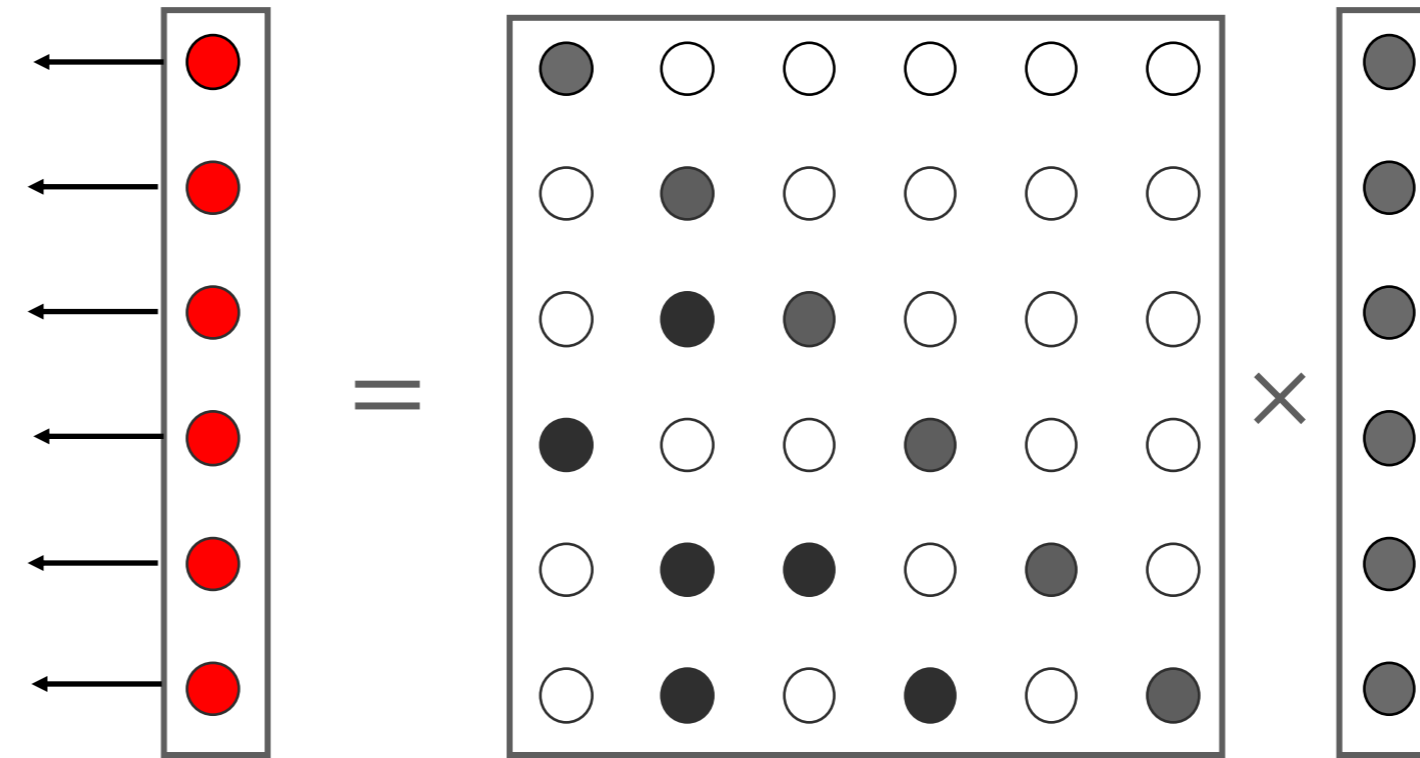
LU factorization, forward-backward solvers, etc.

Matrix Multiplication Routines vs. Solvers

Matrix Multiplication routines:

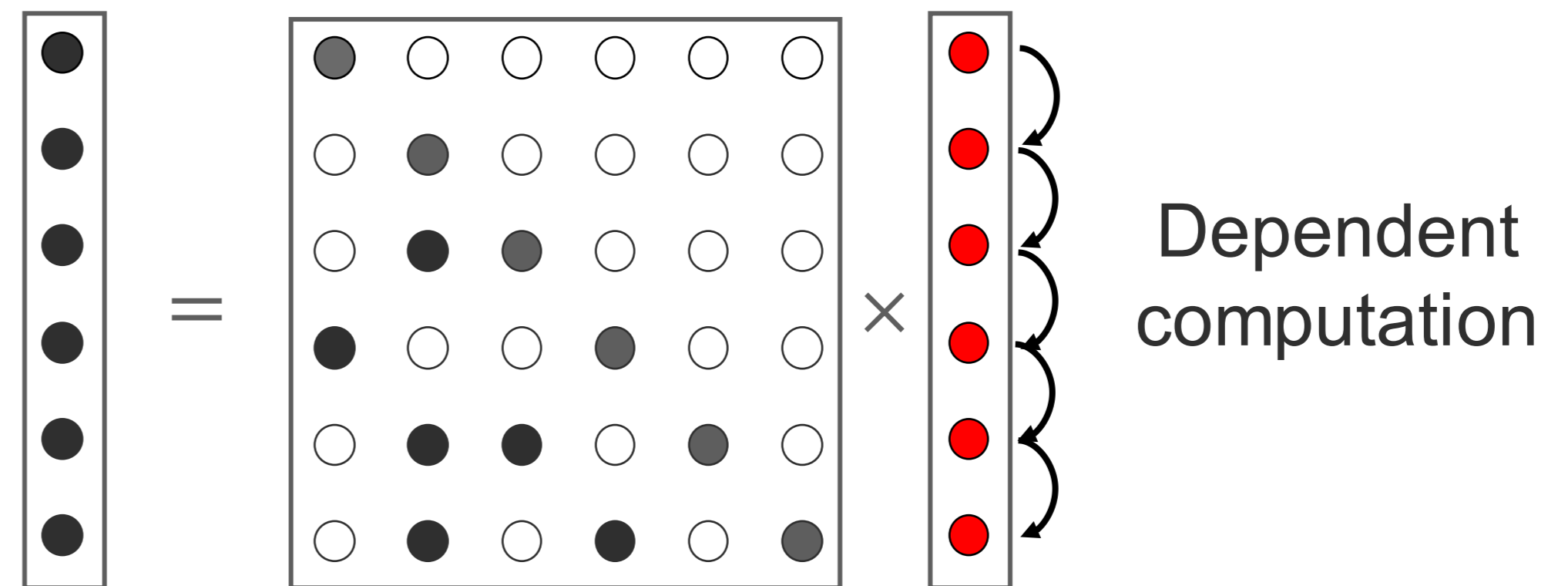
Can be
computed
independently

Unknown



Solvers: e.g. sparse triangular system

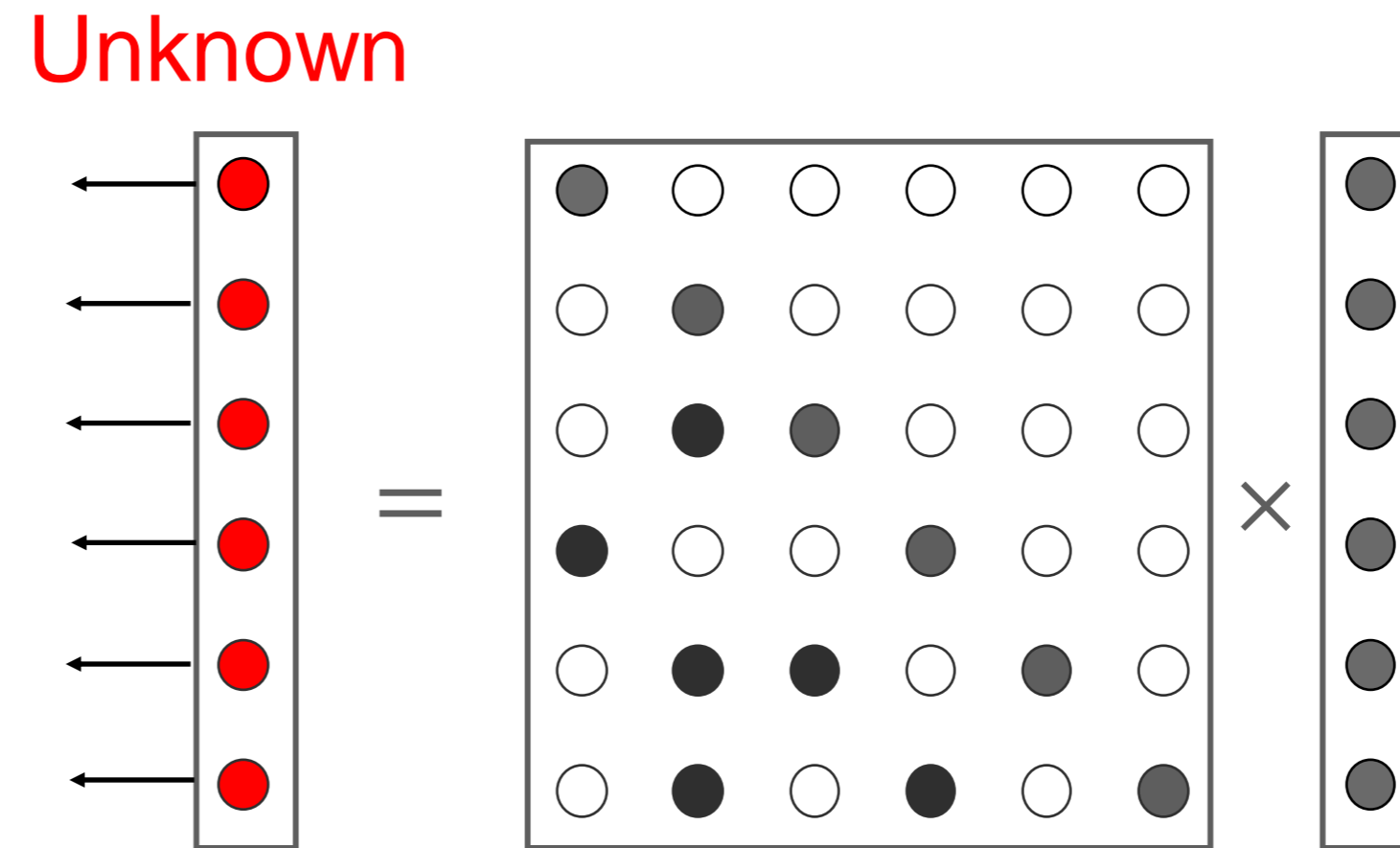
Unknown



Matrix Multiplication Routines vs. Solvers

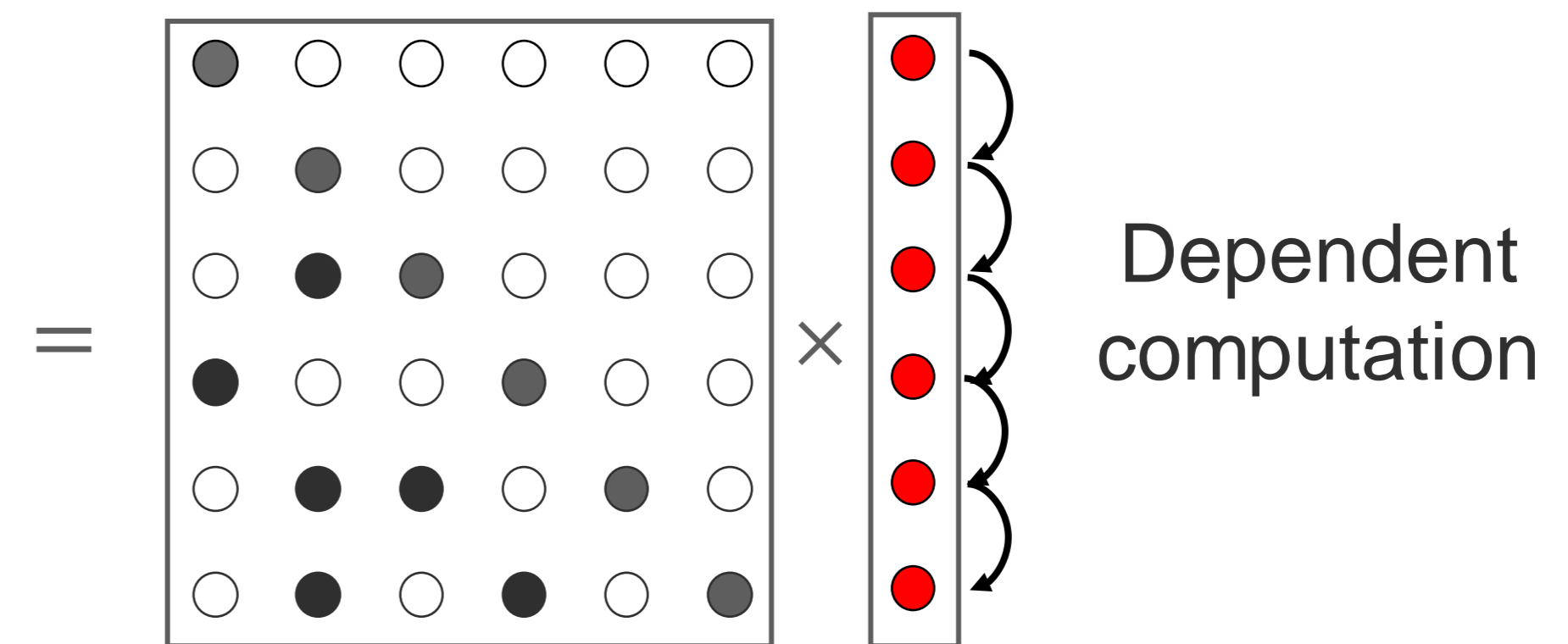
Matrix Multiplication routines:

Can be
computed
independently



Solvers: e.g. sparse triangular system

Unknown



Inspect **data dependence patterns**
for pruning and parallelism, etc.

Solving a Sparse Triangular System

Find the solution to x , $Lx = b$ where L is sparse lower triangular matrix.

$$L = \{n, L_p, L_i, L_x\}$$
$$\begin{bmatrix} 1 & & & & & & & & & \\ & 2 & & & & & & & & \\ & \bullet & 3 & & & & & & & \\ & & & 4 & & & & & & \\ & & & \bullet & 5 & & & & & \\ \bullet & \bullet & \bullet & \bullet & \bullet & 6 & & & & \\ & & & & & & 7 & & & \\ & & & & & \bullet & \bullet & 8 & & \\ & & \bullet & \bullet & \bullet & & & \bullet & 9 & \\ \bullet & & & & & & & \bullet & \bullet & 10 \end{bmatrix} * \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 \\ x_8 \\ x_9 \\ x_{10} \end{bmatrix} = \begin{bmatrix} \bullet \\ \\ \\ \\ \\ \bullet \\ \\ \\ \\ \bullet \end{bmatrix}$$

```
for (j=0; j<n; j++){
    x[j]/=Lx[Lp[j]];
    for (p=Lp[j]+1; p<Lp[j+1]; p++)
        x[Li[p]]-=Lx[p]*x[j];
}
```

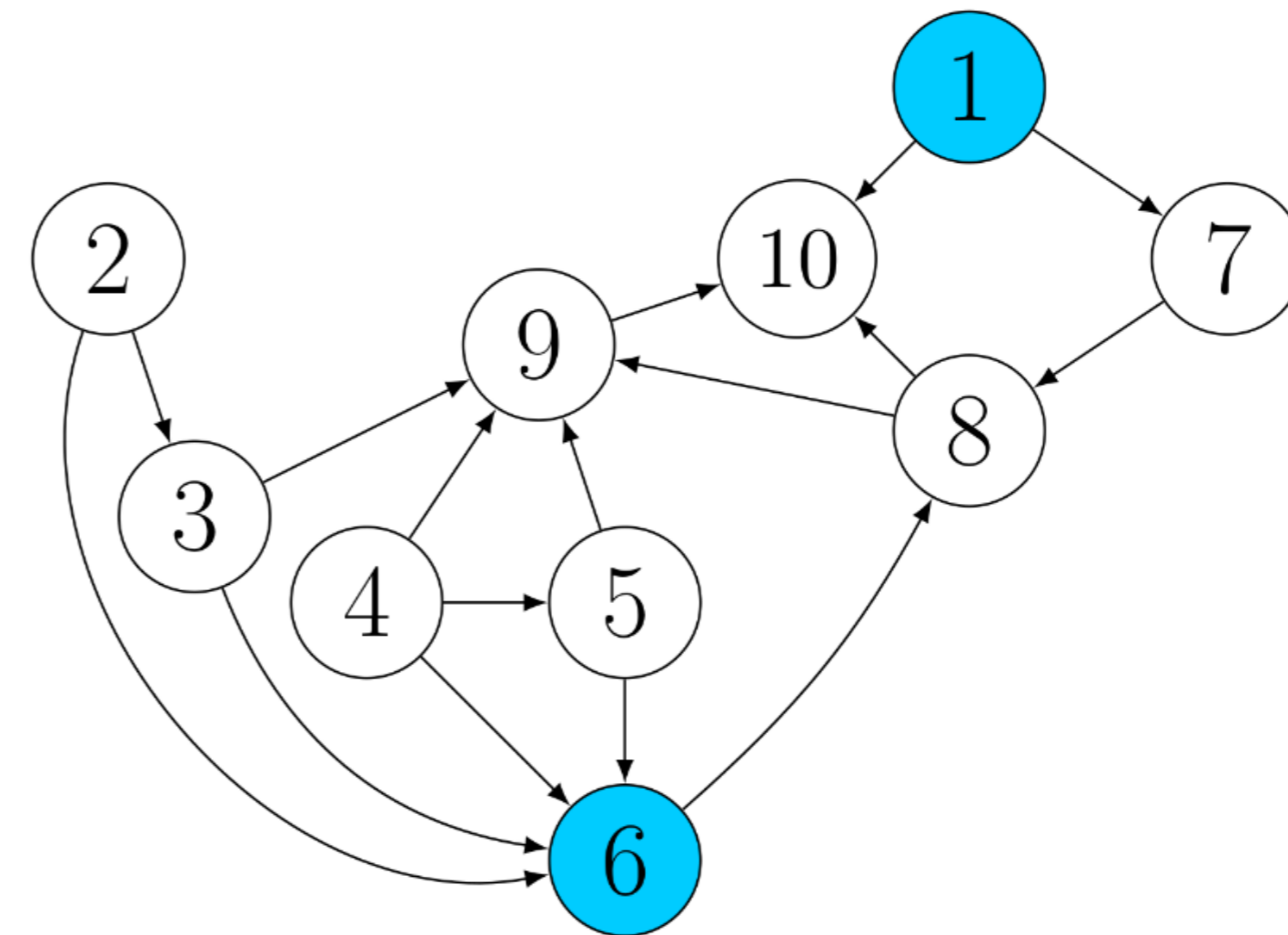
Sparse Analysis in the Sparse Triangular System Solver

$$L = \{n, L_p, L_i, L_x\}$$

$$\begin{bmatrix} 1 & & & & & & & & & \\ & 2 & & & & & & & & \\ & \bullet & 3 & & & & & & & \\ & & & 4 & & & & & & \\ & & & \bullet & 5 & & & & & \\ & \bullet & \bullet & \bullet & \bullet & 6 & & & & \\ & & & & & & 7 & & & \\ \bullet & & & & & & & 8 & & \\ & & \bullet & \bullet & \bullet & & & \bullet & 9 & \\ & & & & & & & \bullet & \bullet & 10 \end{bmatrix} * \begin{bmatrix} x \\ x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 \\ x_8 \\ x_9 \\ x_{10} \end{bmatrix} = \begin{bmatrix} b \\ \bullet \end{bmatrix}$$

Depth First Search (DFS)

Dependence Graph (DG_L)



Inspection

Sparse Analysis in the Sparse Triangular System Solver

$$L = \{n, L_p, L_i, L_x\}$$

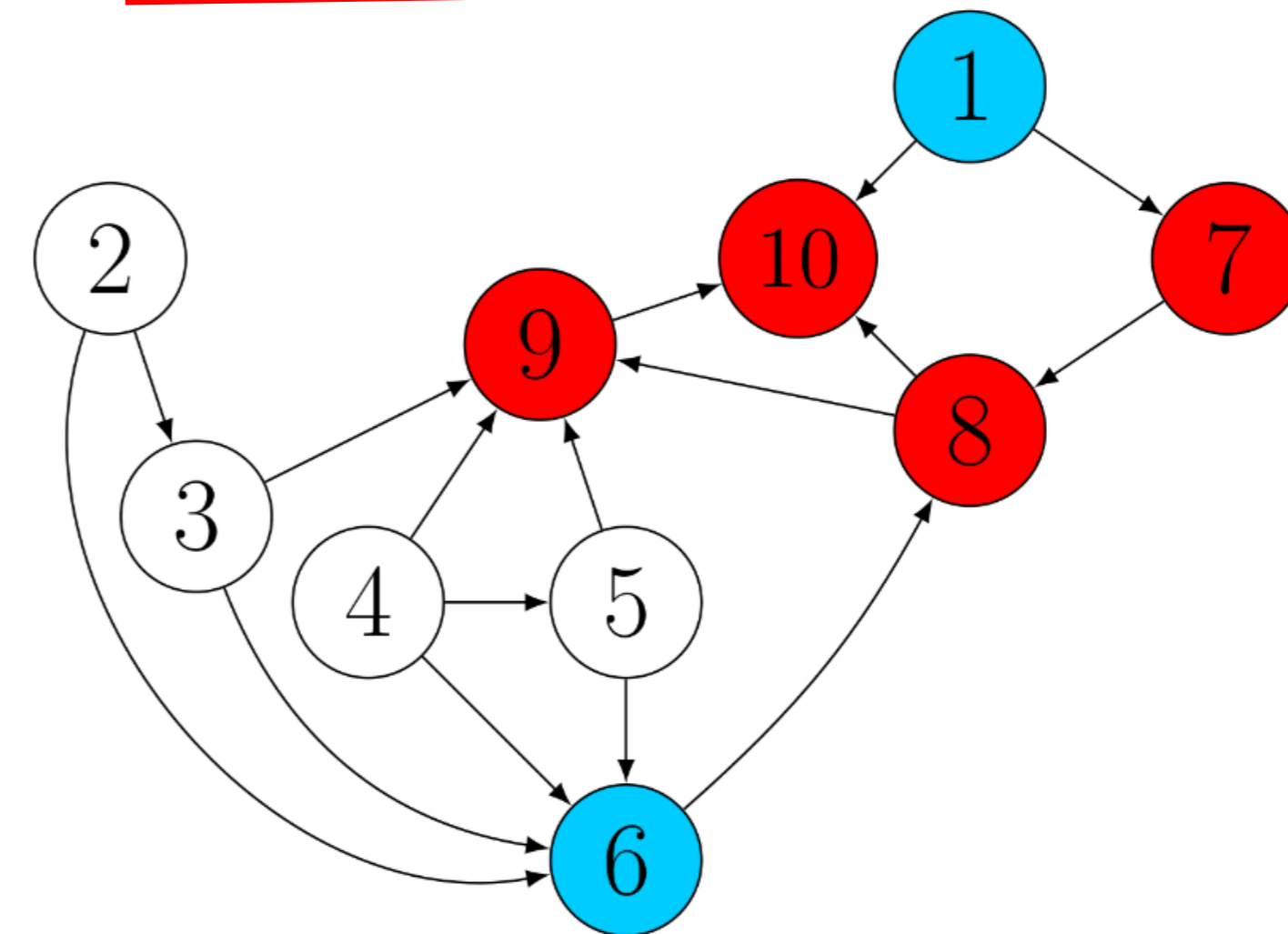
$$\begin{bmatrix} 1 & & & & & & & & & & \\ & 2 & & & & & & & & & \\ & \bullet & 3 & & & & & & & & \\ & & & 4 & & & & & & & \\ & & & \bullet & 5 & & & & & & \\ & \bullet & \bullet & \bullet & \bullet & \bullet & & & & & \\ & \bullet & & & & & \bullet & & & & \\ & & & & & & \bullet & \bullet & & & \\ & & & & & & & \bullet & \bullet & & \\ & & & & & & & & \bullet & \bullet & \\ & & & & & & & & & \bullet & \bullet \end{bmatrix} * \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 \\ x_8 \\ x_9 \\ x_{10} \end{bmatrix} = \begin{bmatrix} \bullet \\ \\ \\ \\ \\ \bullet \\ \\ \\ \\ \bullet \end{bmatrix} b$$

Depth First Search (DFS)

Dependence Graph (DG_L)

$$\beta = \{i | b_i \neq 0\} = \{1, 6\}$$

$$Reach_L(\beta) = \{1, 6, 7, 8, 9, 10\}$$



Inspection

Transforming the Sparse Code

Standard Code

```
for (j=0; j<n; j++){  
  x[j]/=Lx[Lp[j]];  
  
  for (p=Lp[j]+1; p<Lp[j+1]; p++)  
    x[Li[p]]-=Lx[p]*x[j];  
}
```

Prune the iteration space:
Reachset = {6,1,7,8,9,10}

Symbolically-Guided Code

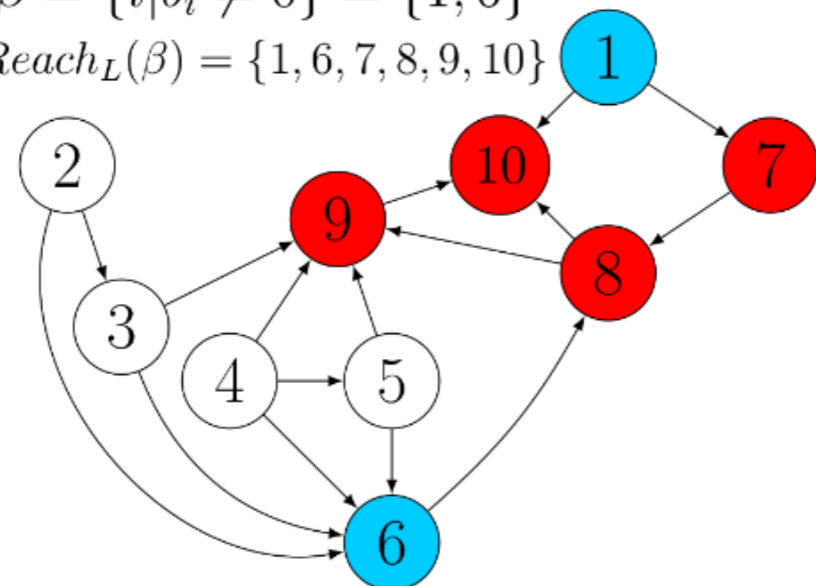
```
for (px=0; px<RSsize; px++) {  
  j=reachset[px];  
  x[j]/=Lx[Lp[j]]  
  p=Lp[j]+1;  
  for (; p<Lp[j+1]; p++)  
    x[Li[p]]-=Lx[p]*x[j];  
}
```

Code Gen

Dependence Graph (DG_L)

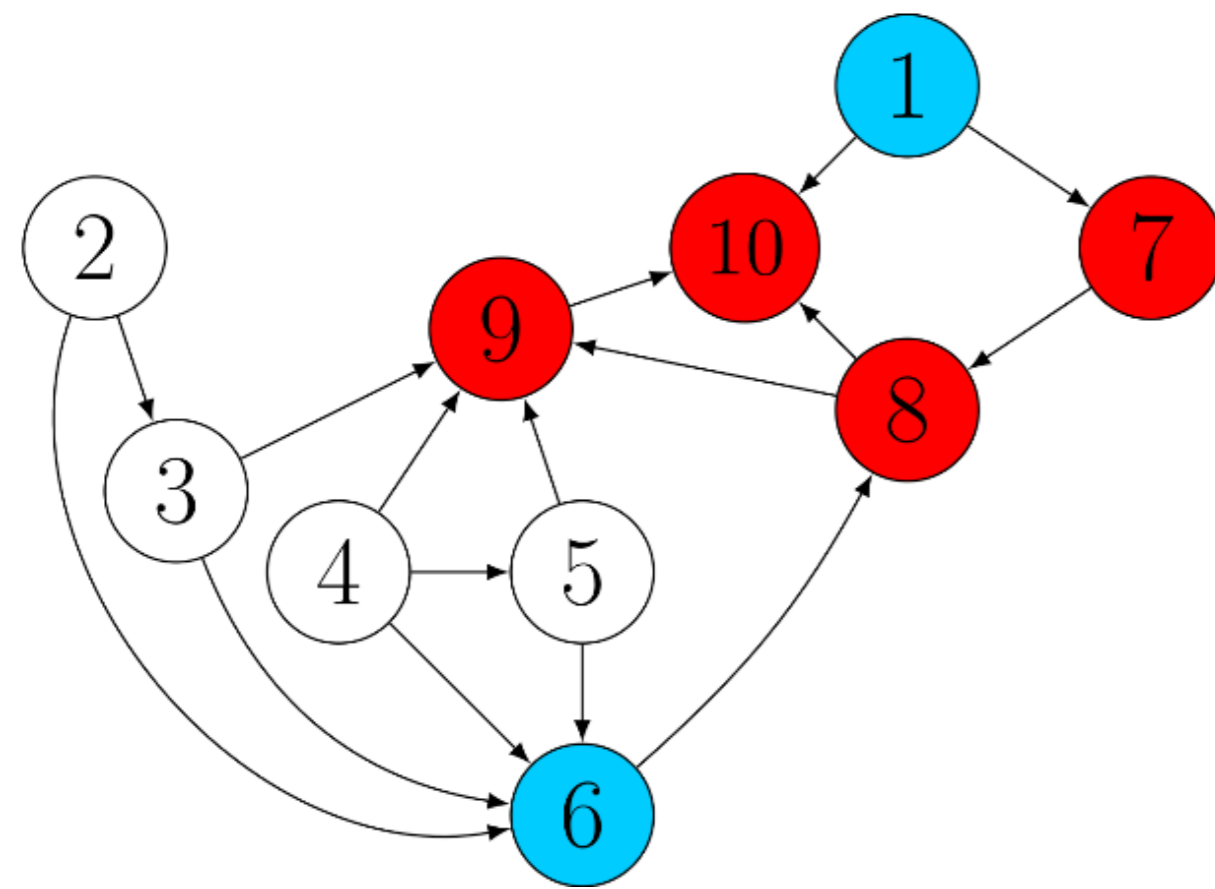
$\beta = \{i | b_i \neq 0\} = \{1, 6\}$

$Reach_L(\beta) = \{1, 6, 7, 8, 9, 10\}$

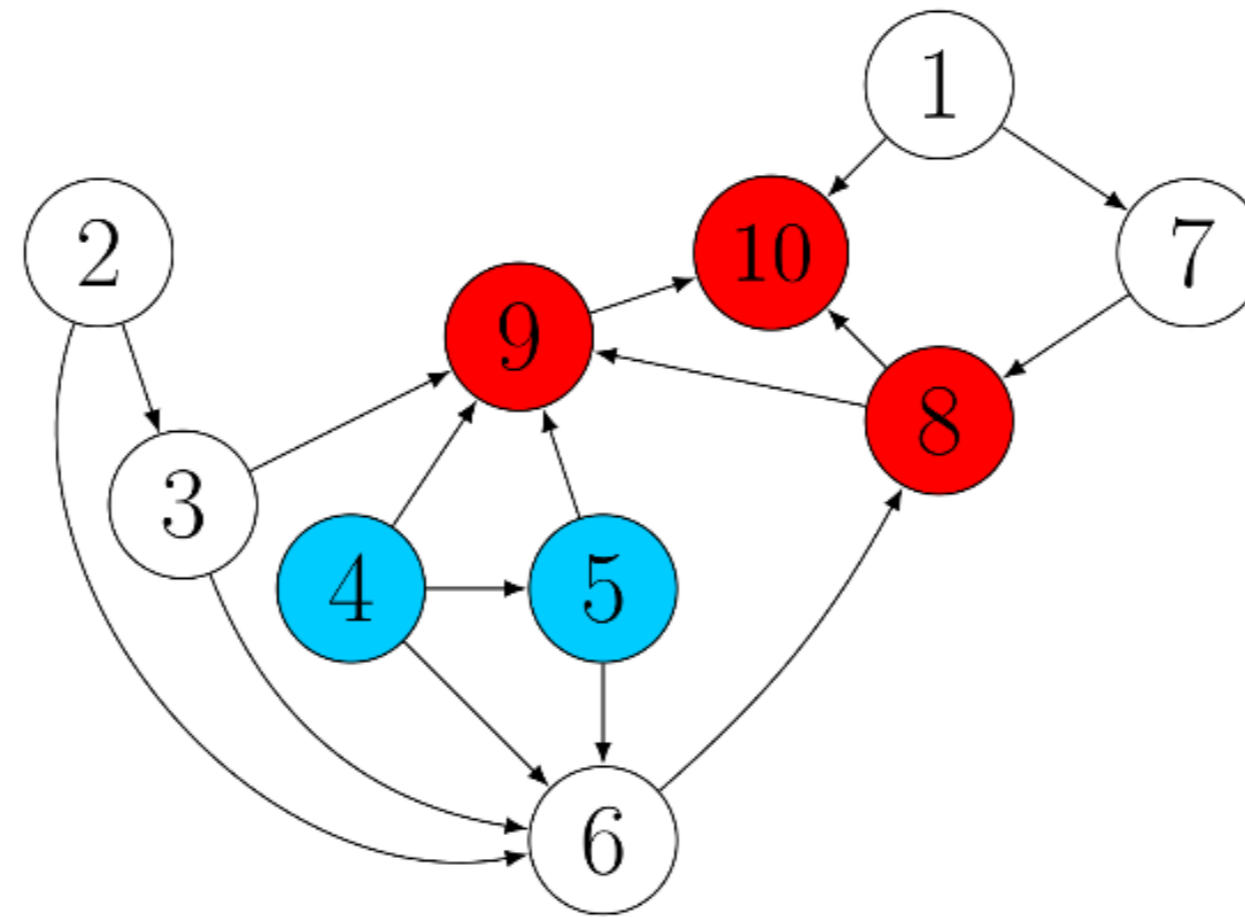


Sparse Inspection and Code Transformations for Sparse Solvers

Sympiler [SC'17]: Inspects for single core optimizations, i.e. pruning, blocking, etc.



Depth First Search
For **Pruning**



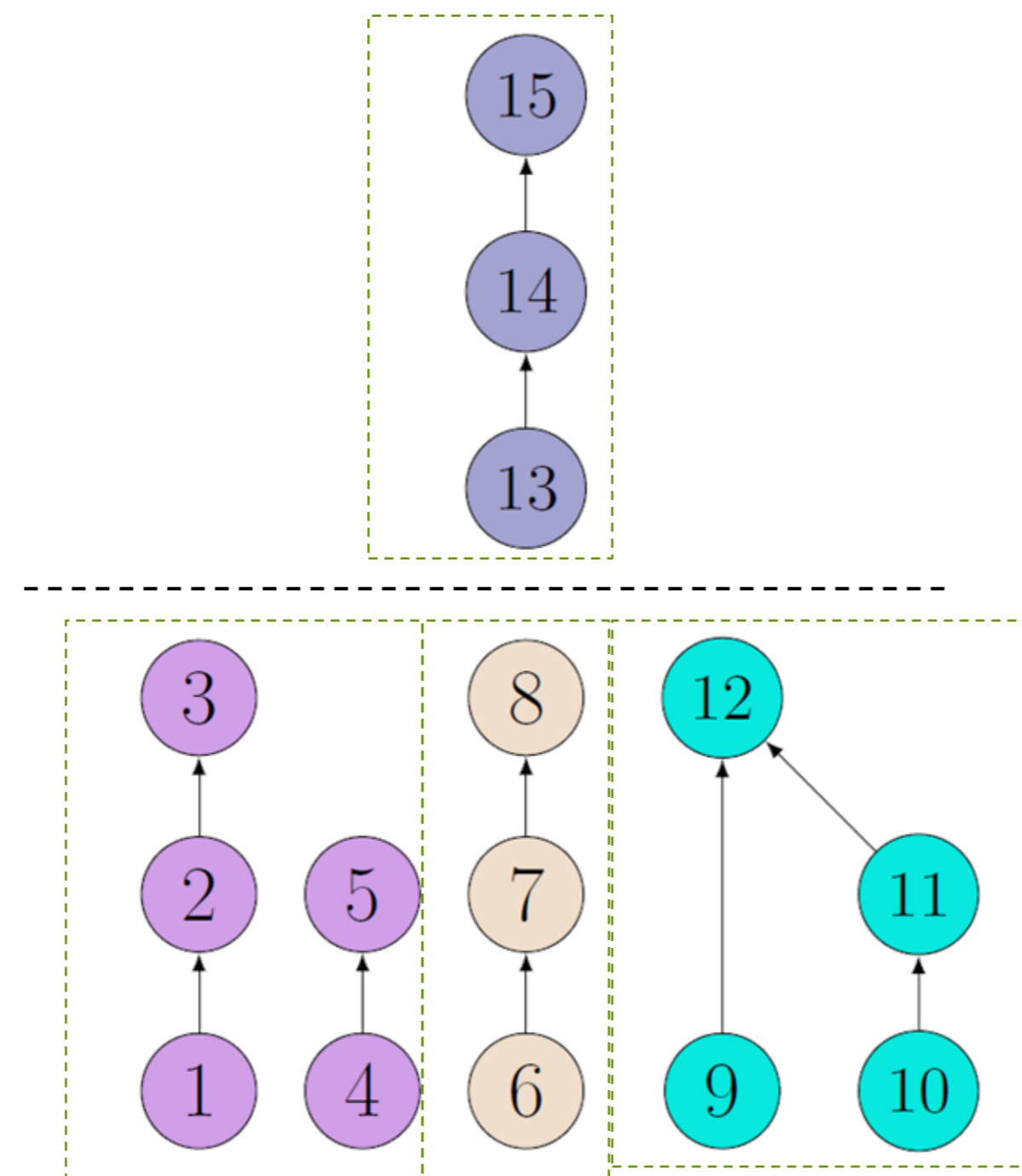
Node Equivalence for
Blocking

Sparse Inspection and Code Transformations for Sparse Solvers

Sympiler [SC'17]: Inspects for single core optimizations, i.e. pruning, blocking, etc.

ParSy [SC'18] and HDAGG [IPDPS'22]: Create well balanced workloads for multicore execution.

+



Load-Balanced Level
Coarsening for **Parallelism**

Sparse Inspection and Code Transformations for Sparse Solvers

Sympiler [SC'17]: Inspects for single core optimizations, i.e. pruning, blocking, etc.

ParSy [SC'18] and HDAGG [IPDPS'22]: Create well balanced workloads for multicore execution

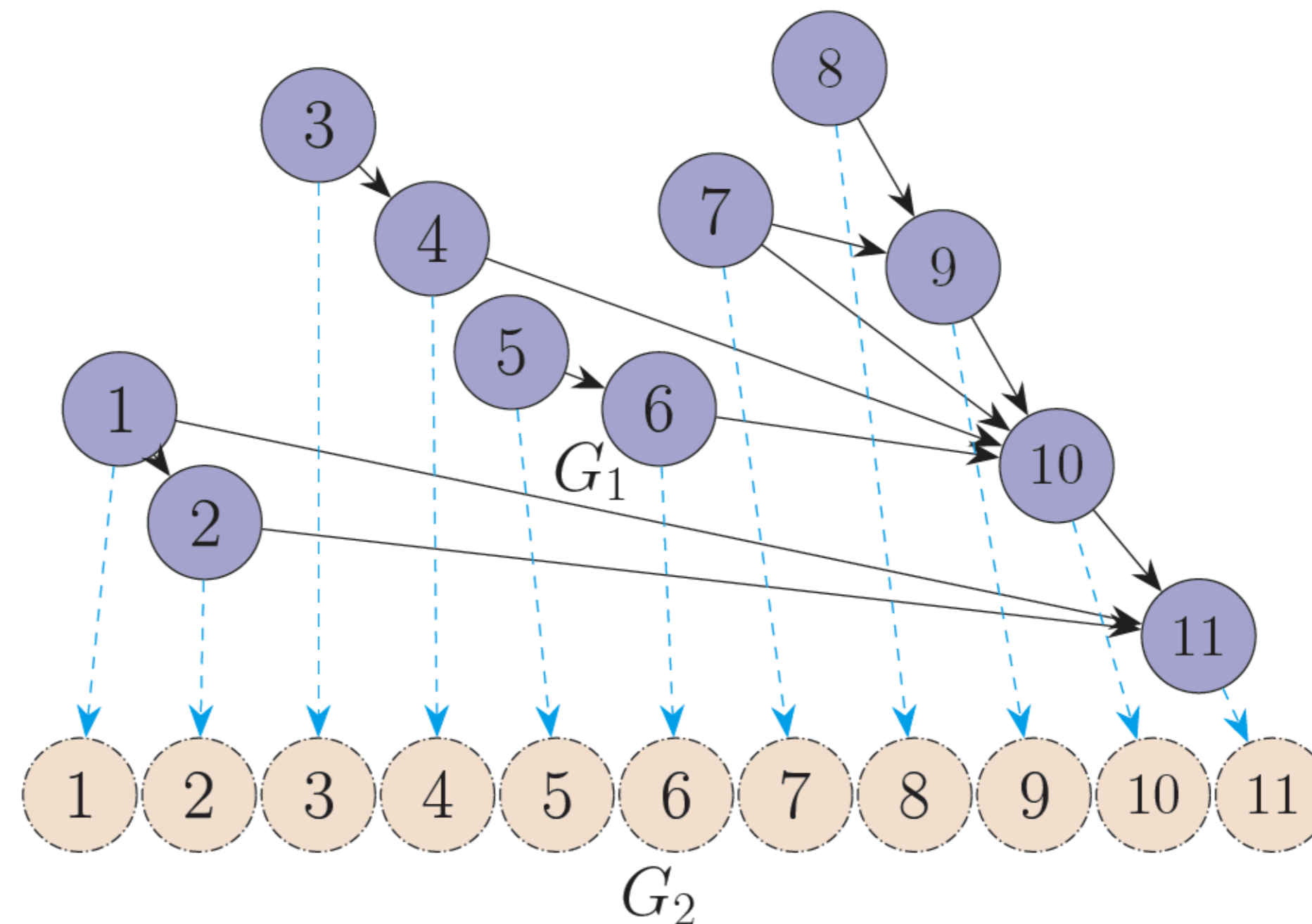
Sparse Fusion [SC'23]: Inspects the graphs of multiple operations to generate parallel fused code.

+

+

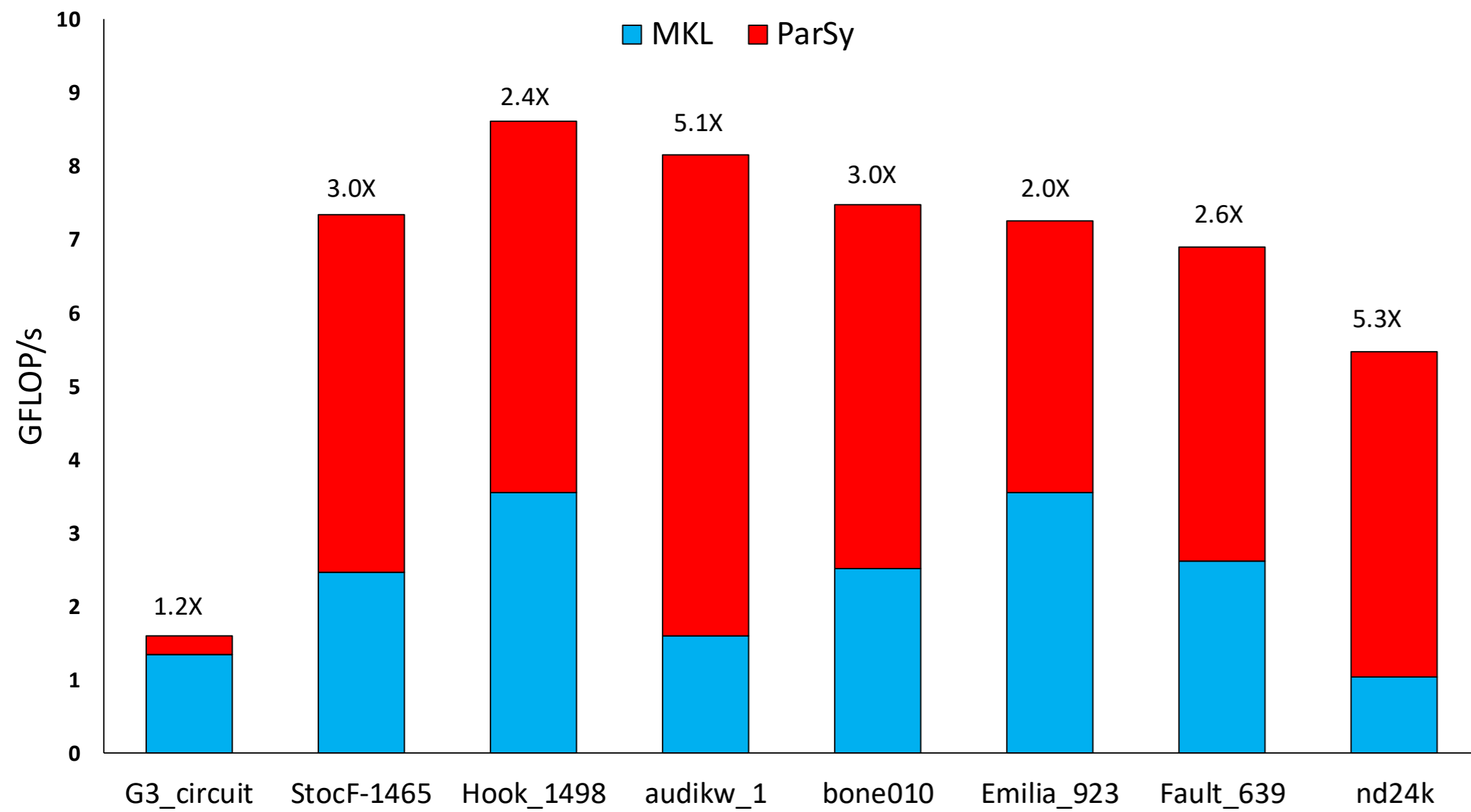
$$Lx = b$$

$$y = Ax$$



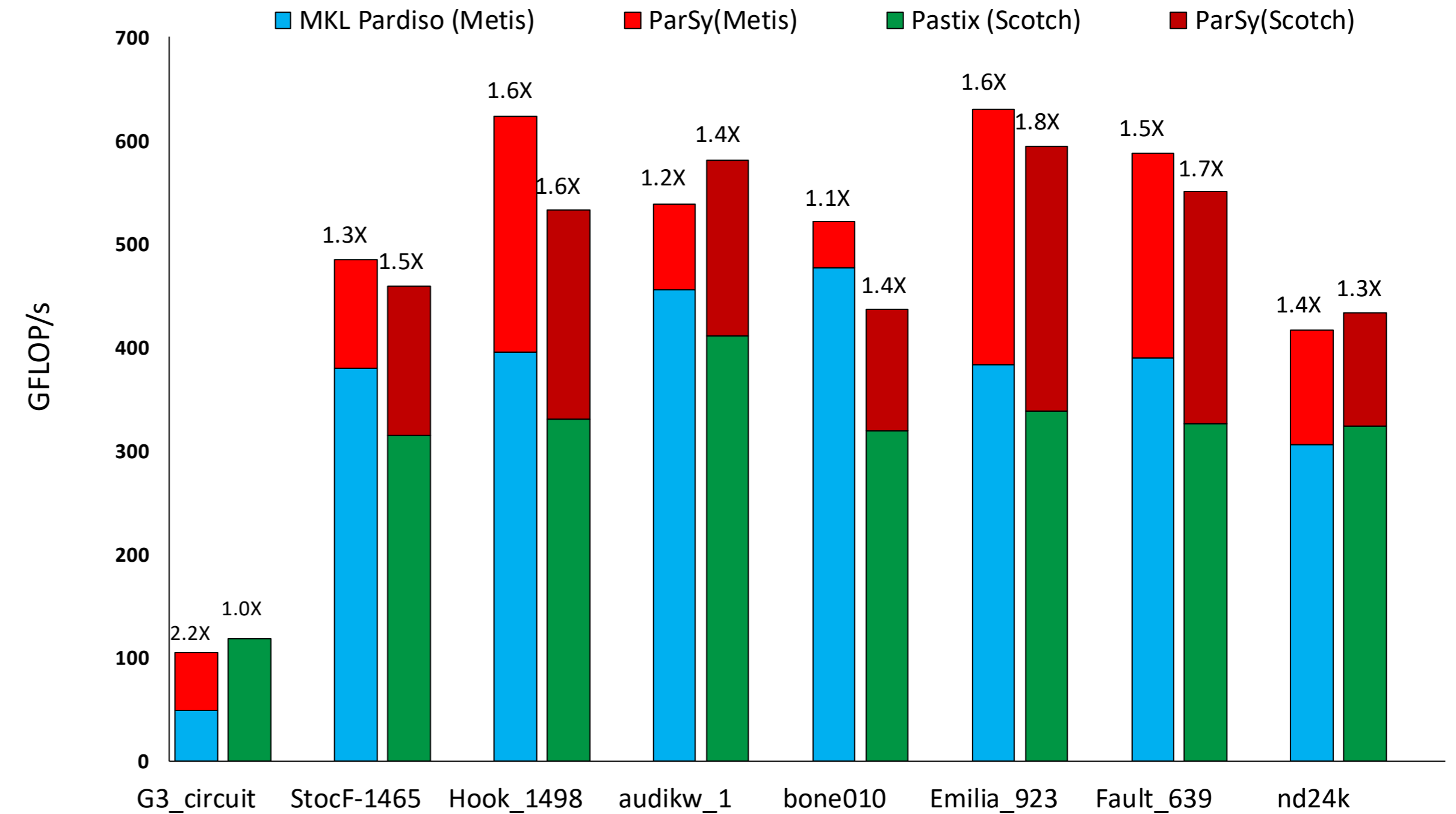
ParSy vs Competitors

Triangular Solve



ParSy is faster than Intel **MKL** with an average speedup of **2.5x**.

Cholesky Factorization



ParSy is faster than **Pardiso** with an average speedup of **1.5x**.

There is no One-Size-Fits-All Approach!

Thesis

Sparsity is often **static** or changes moderately in most simulations.

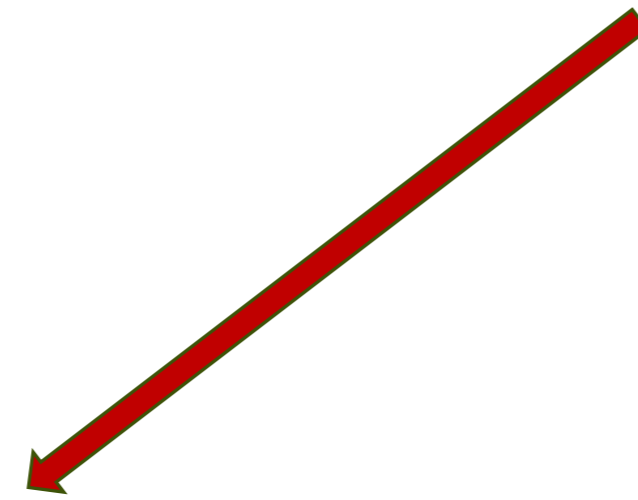
Key contribution

Inspect sparsity patterns to automatically generate highly-optimized code for sparsity.

There is no One-Size-Fits-All Approach!

Thesis

Sparsity is often static or **changes moderately** in most simulations.



Tweak the systems:

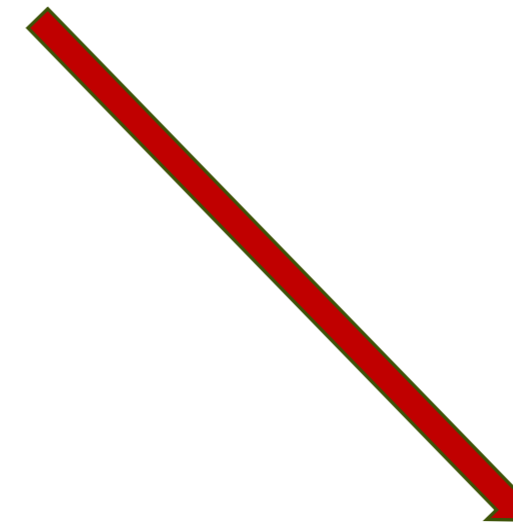
- ❖ Relax sparsity-specific specialization
- ❖ Runtime checks
- ❖ Just-in-time compilation

Typically reduces our performance gains

There is no One-Size-Fits-All Approach!

Thesis

Sparsity is often static or **changes moderately** in most simulations.



Change the algorithms!

- ❖ Constrained/unconstrained optimization in graphics
- ❖ Machine Learning Training

A Constrained Optimization Problem: QP Solvers

NASOQ: A Numerically Accurate Sparsity-Oriented Quadratic Program Solver [Siggraph'20]

$$\min_x \frac{1}{2} x^T H x + q^T x$$

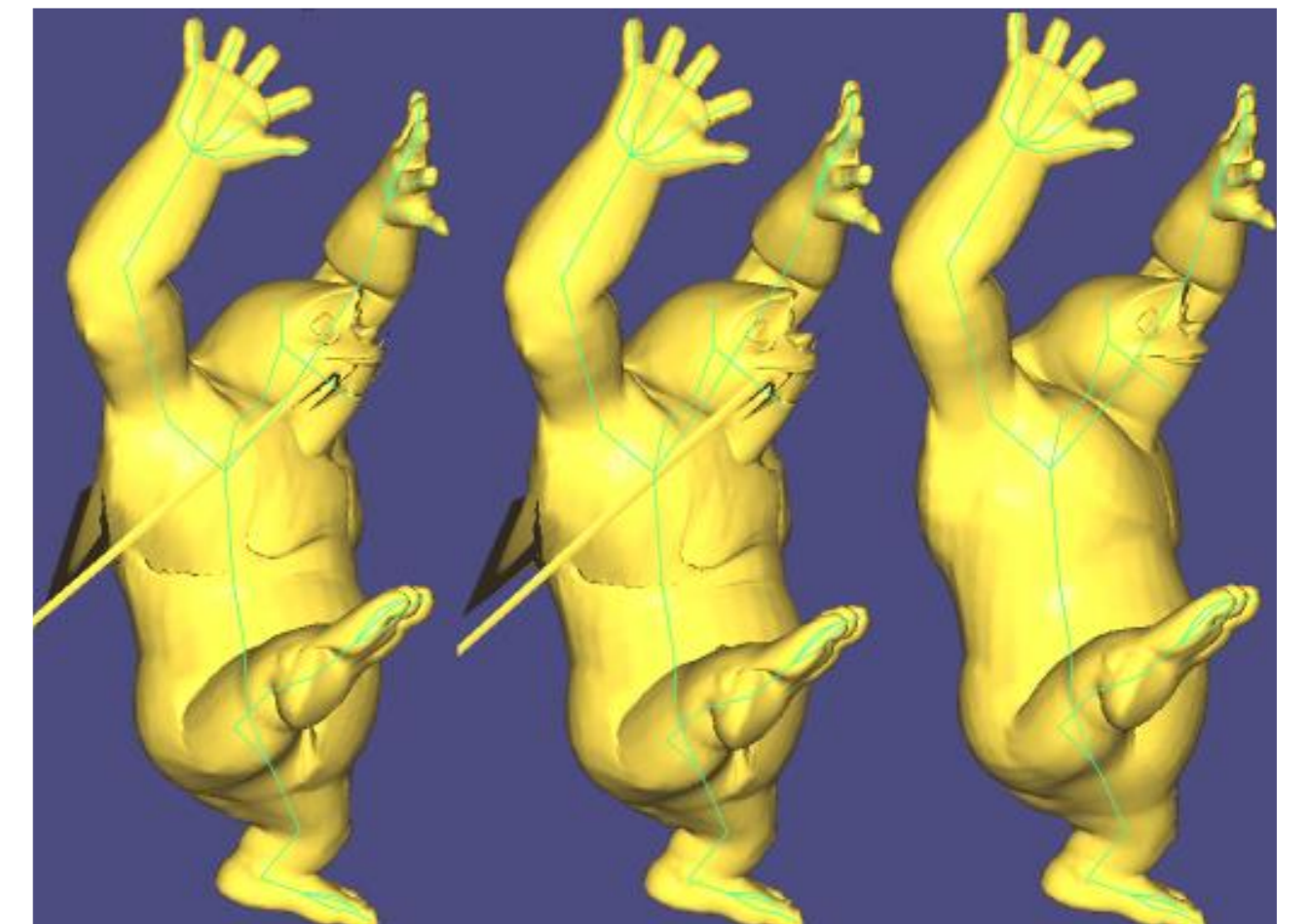
$$Ax = b$$

$$Cx \leq d$$

Objective

Constraints

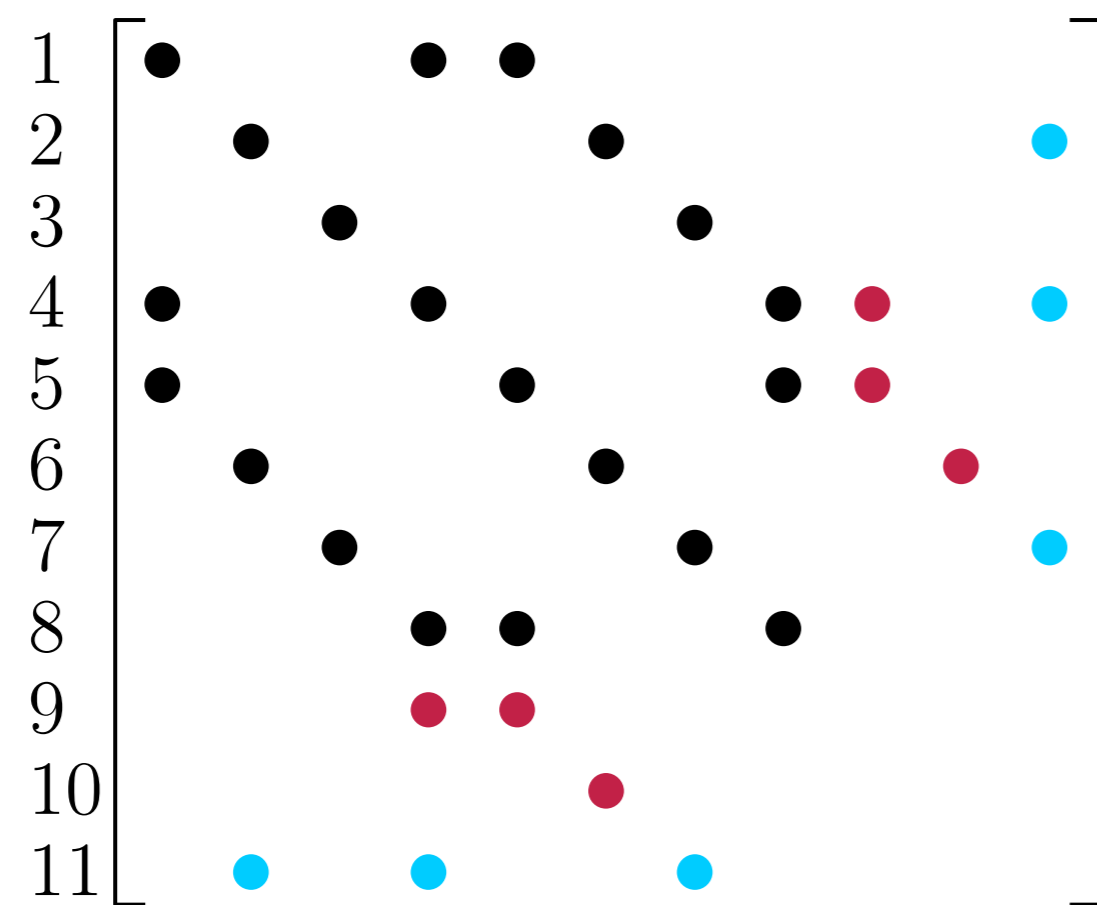
A case for dynamic sparsity



<https://nasoq.github.io/>

The Goldfarb and Idnani Algorithm

In the optimality phase, for each added or removed constraint a new KKT system has to be solved.



KKT matrix

Optimality Phase

while(x_k is not optimal) {

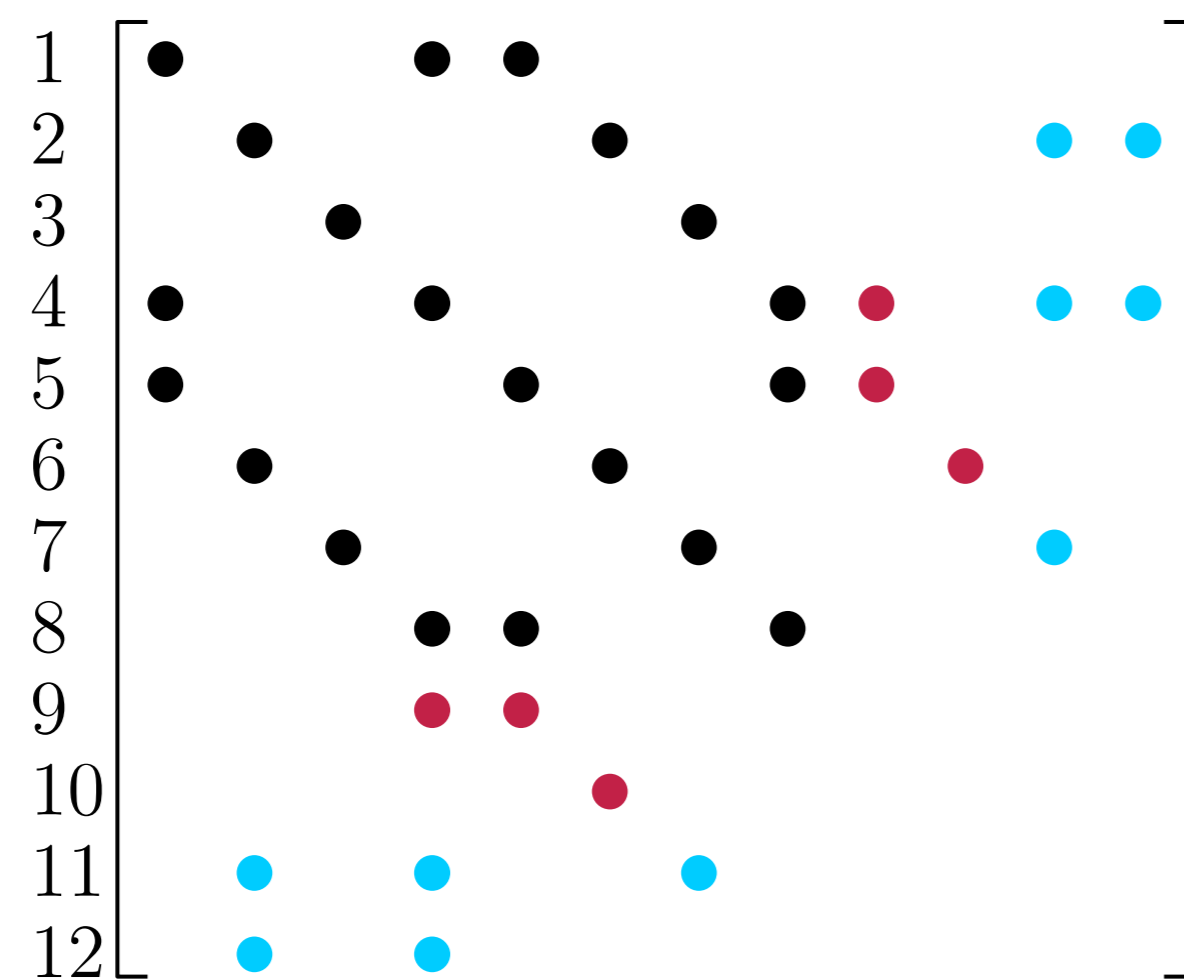
$$[\alpha, \pi, \mu] \leftarrow \begin{bmatrix} H & A^T & C_w^T \\ A & 0 & 0 \\ C_w & 0 & 0 \end{bmatrix}^{-1} \begin{bmatrix} C_{act} \\ 0 \\ 0 \end{bmatrix}$$

...
}

Factorize the Updated KKT Matrix

The Goldfarb and Idnani Algorithm

Solving incrementally increasing KKT systems, has large overheads and thus leads to poor scalability.



KKT matrix

Optimality Phase

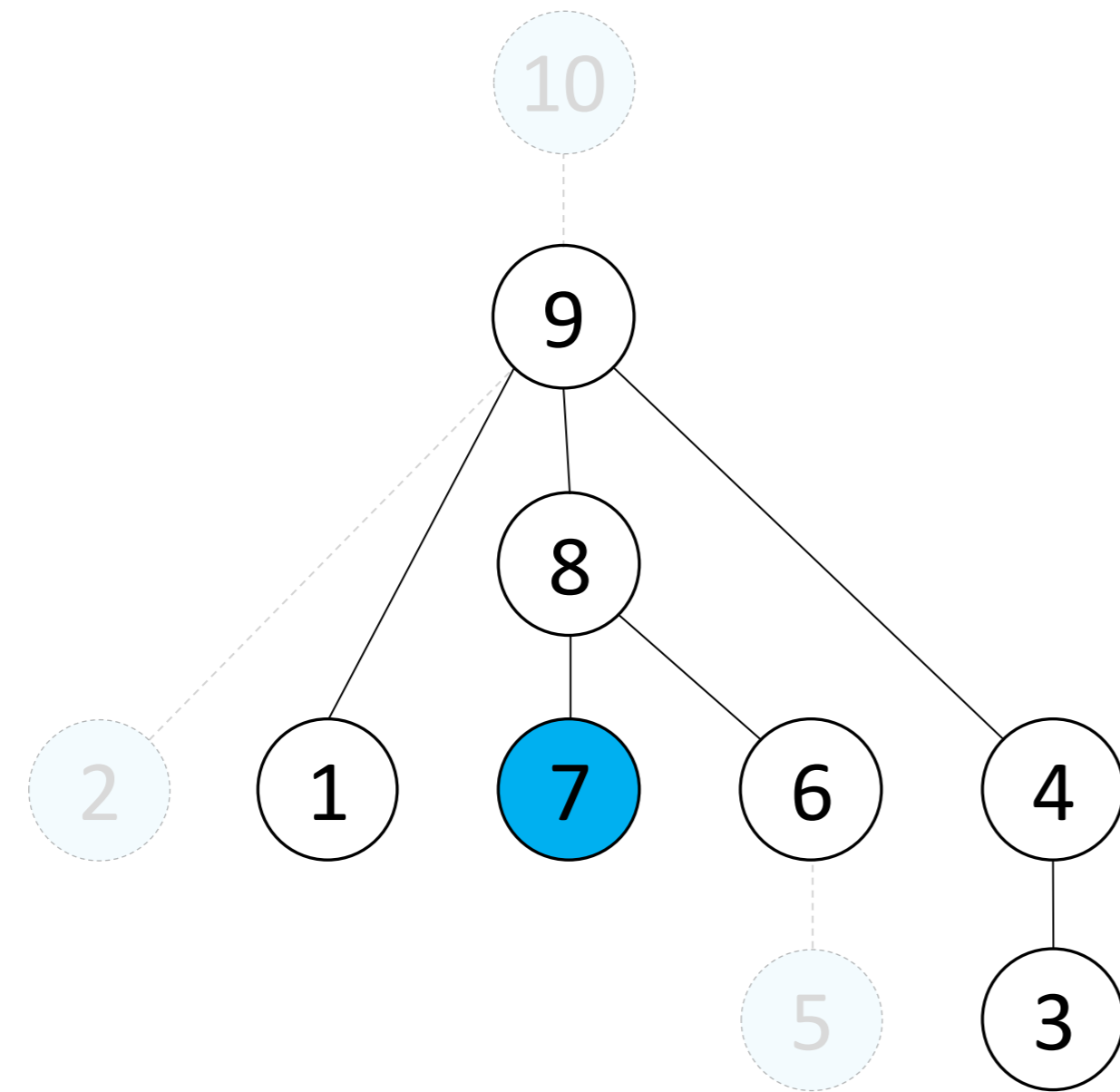
while(x_k is not optimal) {

$$[\alpha, \pi, \mu] \leftarrow \begin{bmatrix} H & A^T & C_w^T \\ A & 0 & 0 \\ C_w & 0 & 0 \end{bmatrix}^{-1} \begin{bmatrix} C_{act} \\ 0 \\ 0 \end{bmatrix}$$

...
}

Factorize the Updated KKT Matrix

Sparsity updates: the Optimality Phase

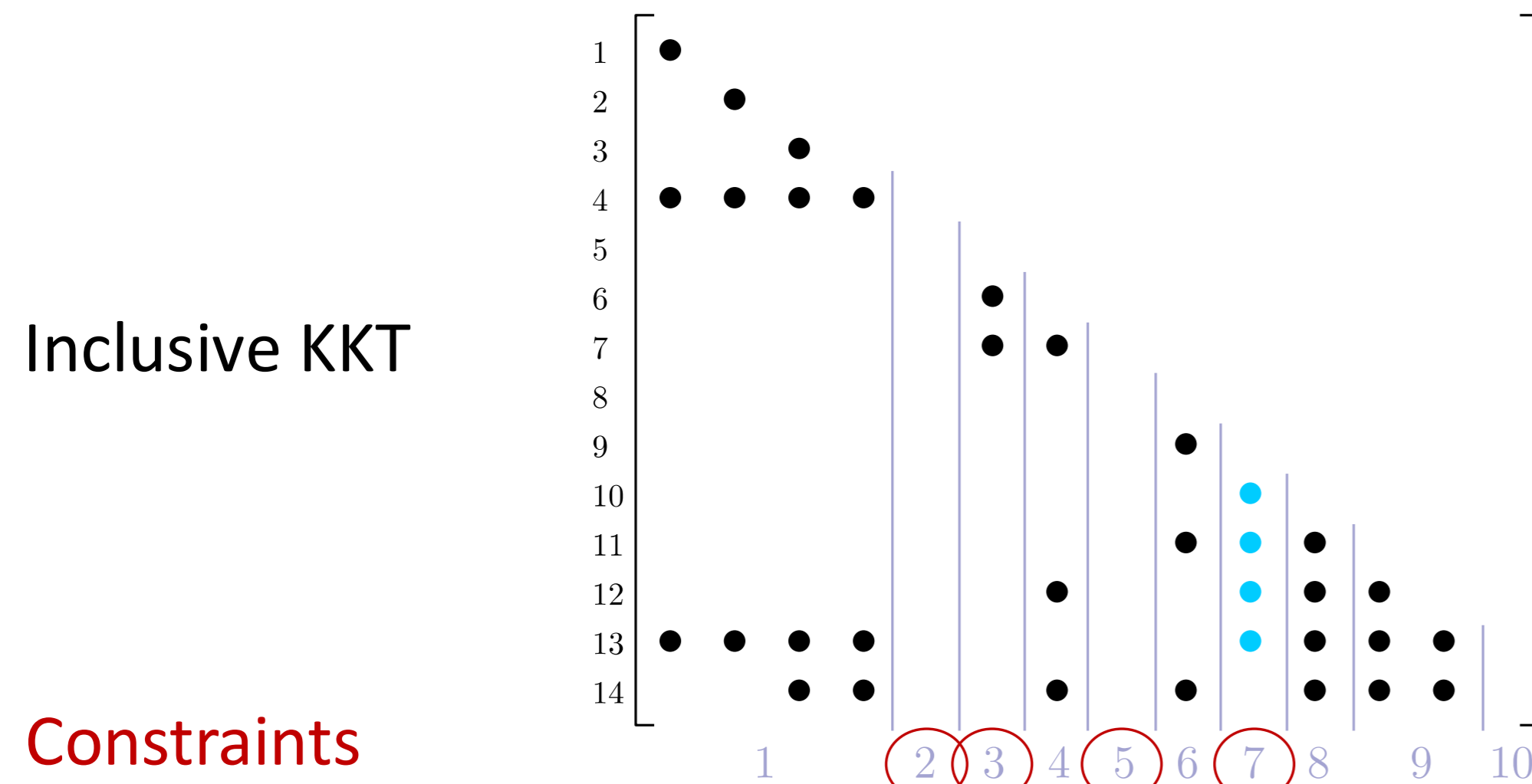


Optimality Phase

```

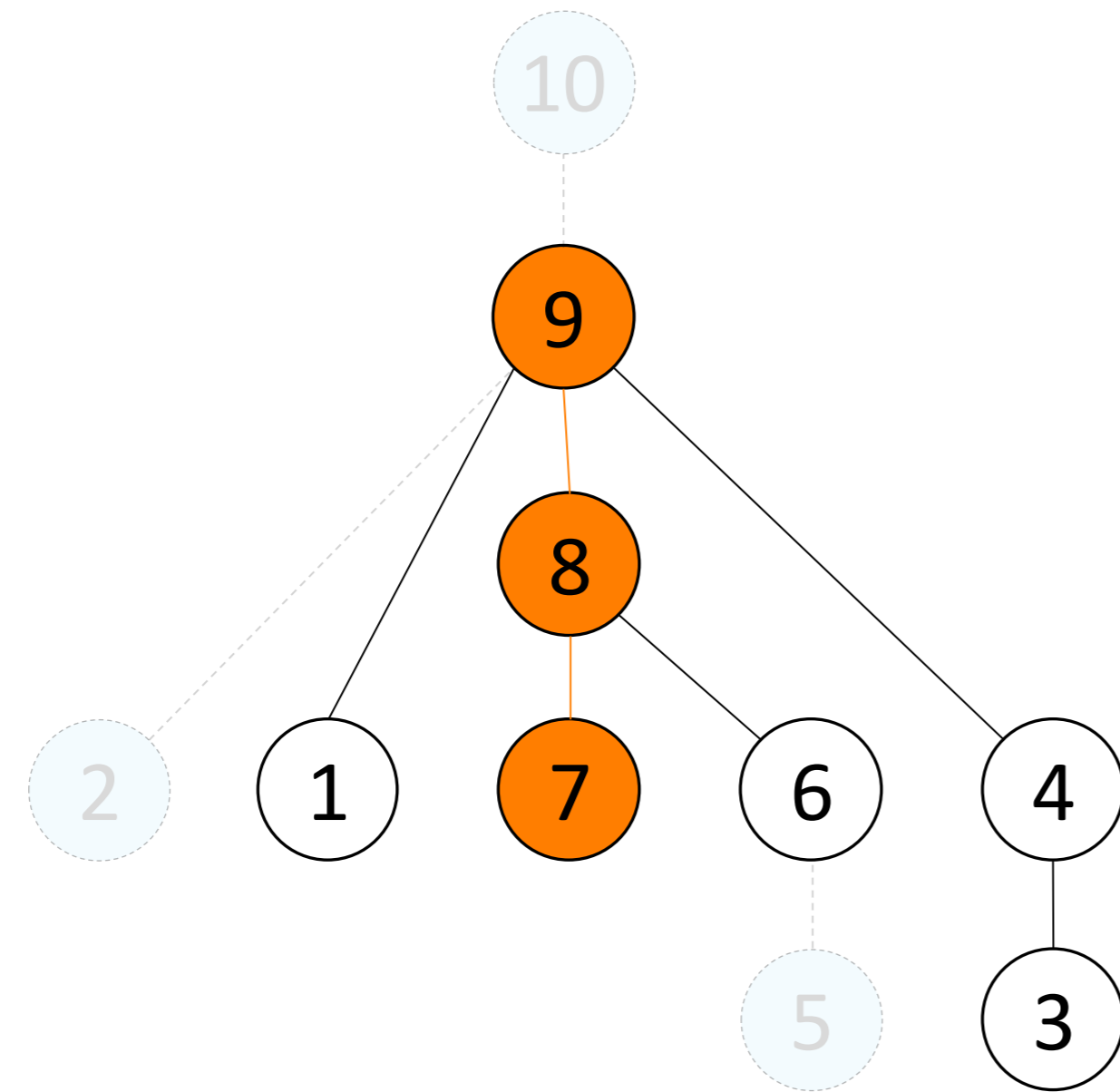
while( $x_k$  is not optimal) {
     $[\alpha, \pi, \mu] \leftarrow \begin{bmatrix} H & A^T & C_w^T \\ A & 0 & 0 \\ C_w & 0 & 0 \end{bmatrix}^{-1} \begin{bmatrix} C_{act} \\ 0 \\ 0 \end{bmatrix}$ 
    ...
}

```



For each added inequality constraint, the corresponding column in the inclusive KKT matrix and the dependence graph (elimination tree in this case) are activated.

Sparsity updates: the Optimality Phase



Optimality Phase

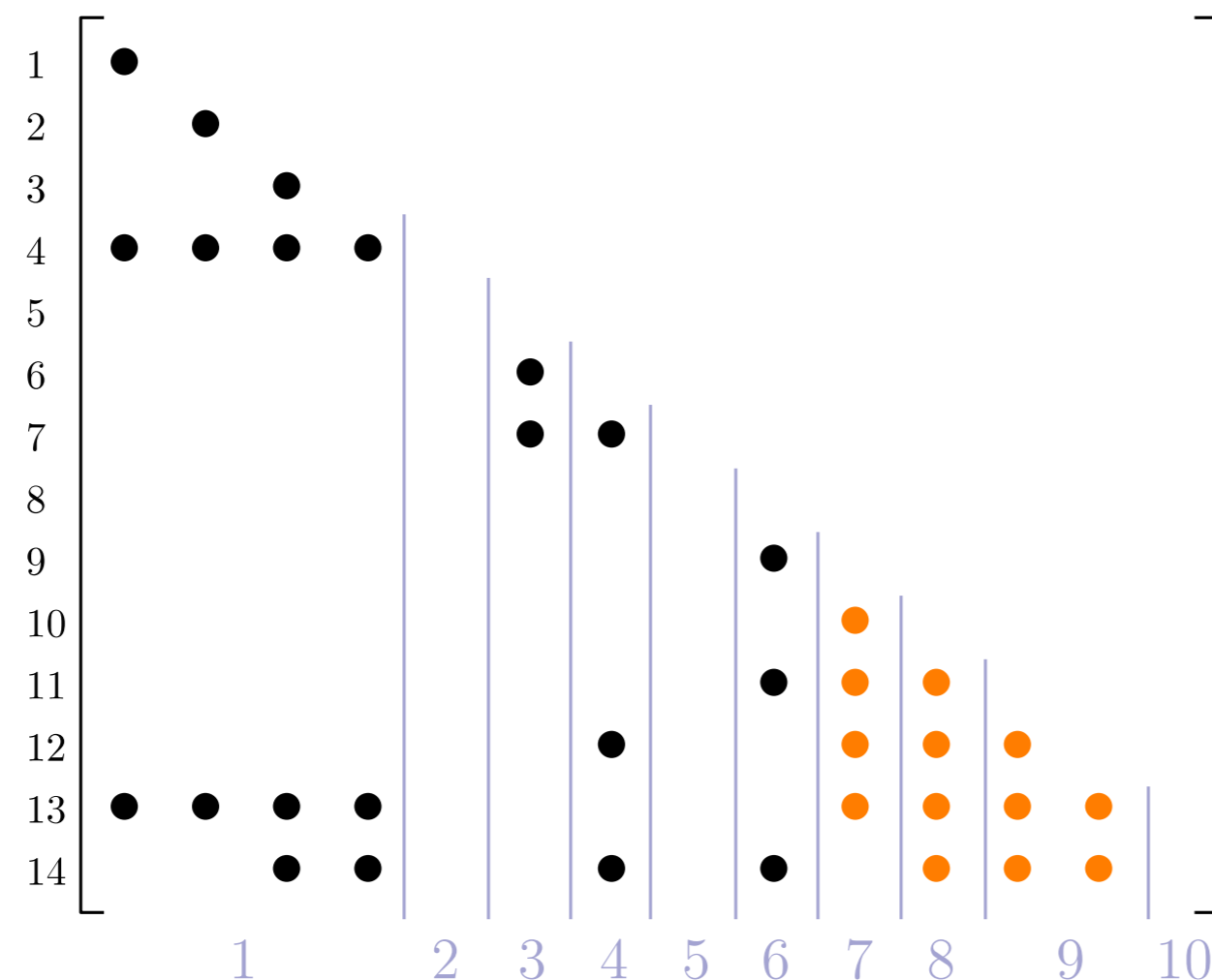
```

while( $x_k$  is not optimal) {
     $[\alpha, \pi, \mu] \leftarrow \begin{bmatrix} H & A^T & C_w^T \\ A & 0 & 0 \\ C_w & 0 & 0 \end{bmatrix}^{-1} \begin{bmatrix} C_{act} \\ 0 \\ 0 \end{bmatrix}$ 
    ...
}

```

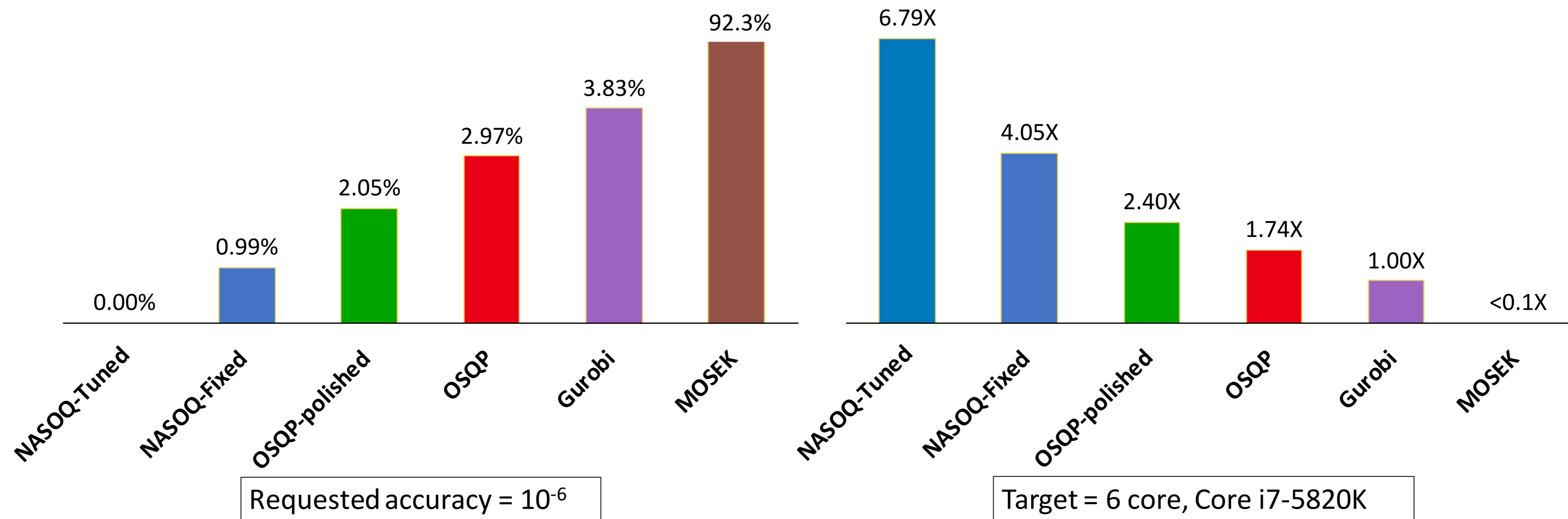
Only the columns affected by the newly added constraint need to be numerically updated, the rest of the values are reused.

Inclusive KKT



Results: NASOQ vs Others

NASOQ has the lowest **failure rate** AND the fastest **solve time**



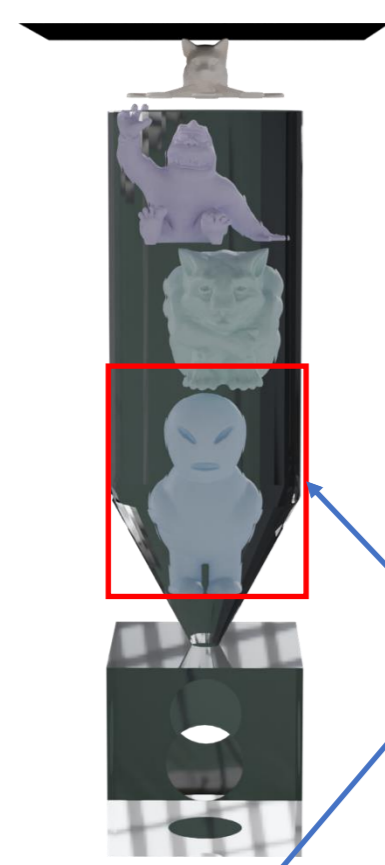
Gurobi is the baseline for speedups



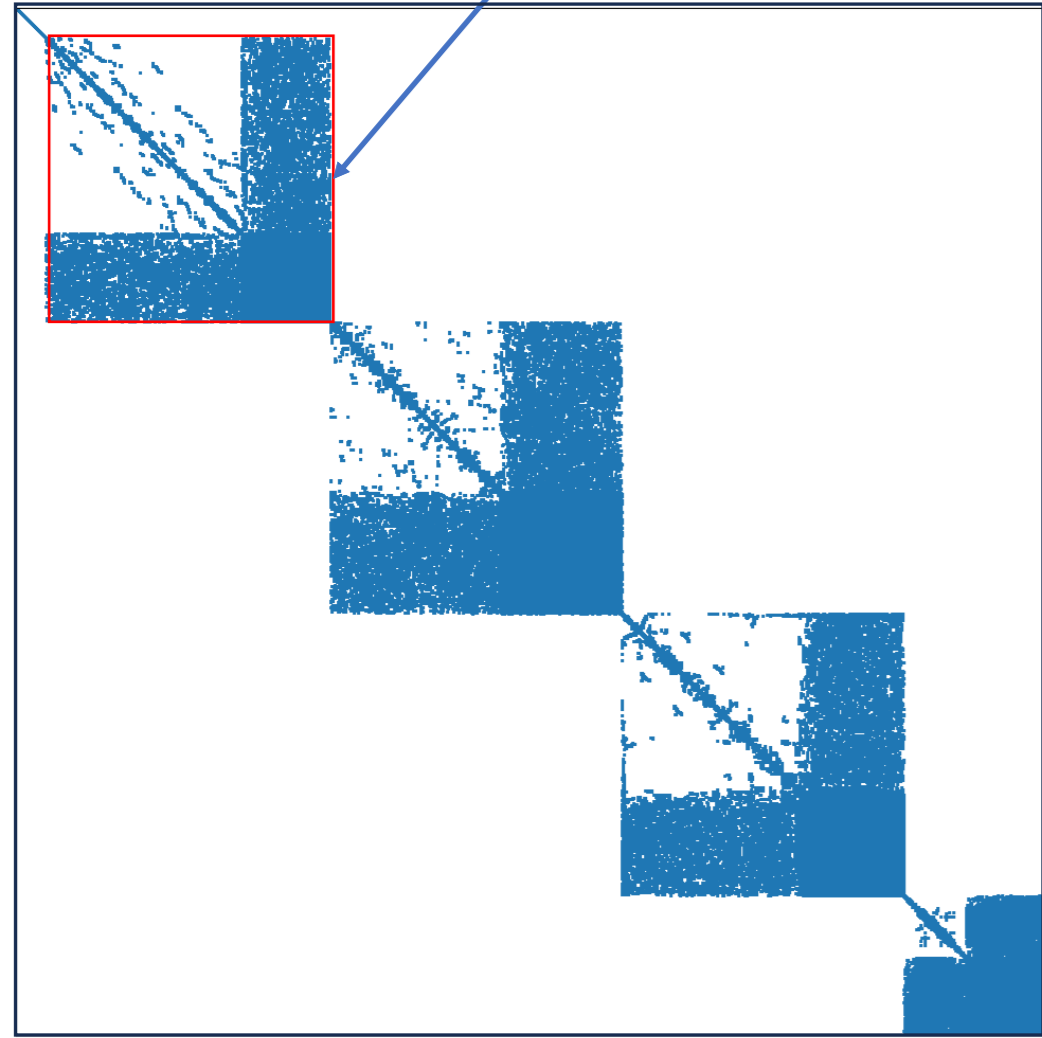
<https://nasoq.github.io/>

Unconstrained Optimization Problem: Newton Solver for Contact

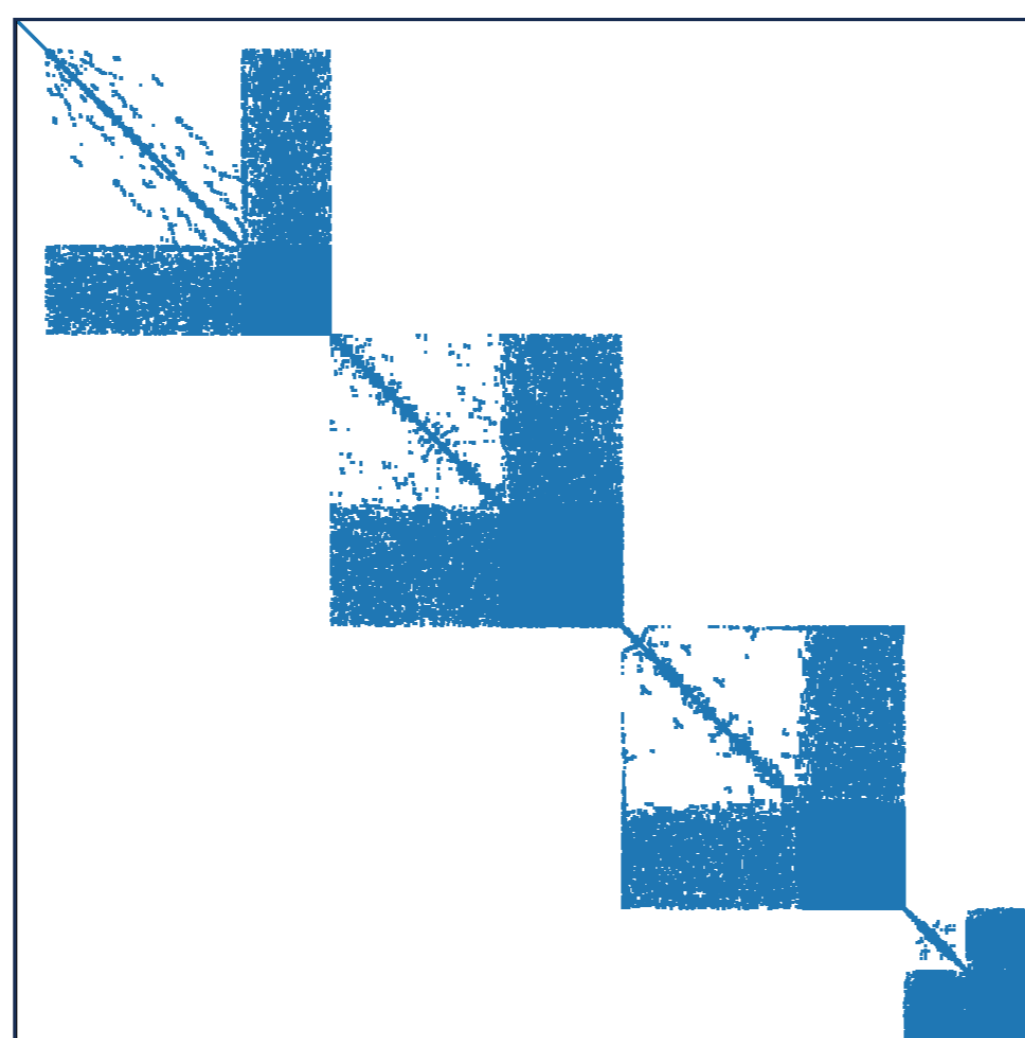
A case for dynamic sparsity



Hessian sparsity pattern comes from the geometry



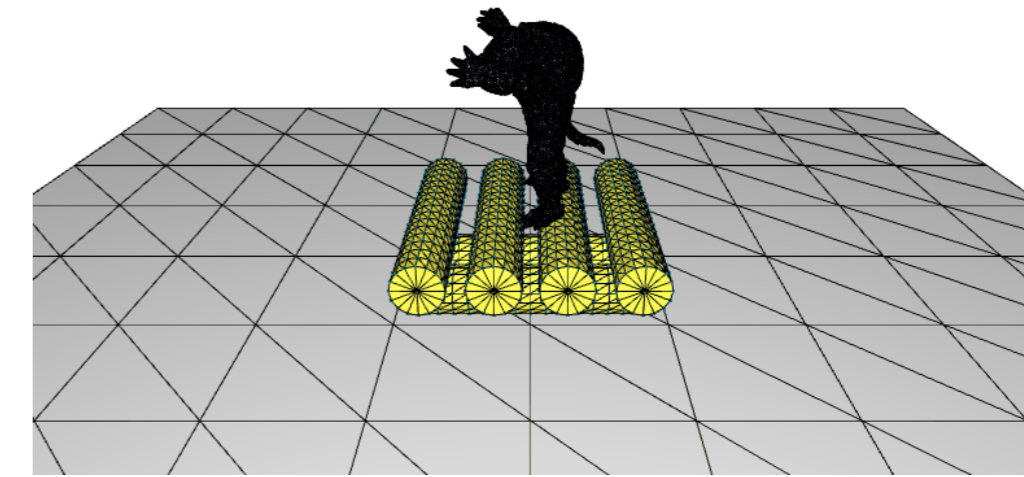
Frame 13: H_0



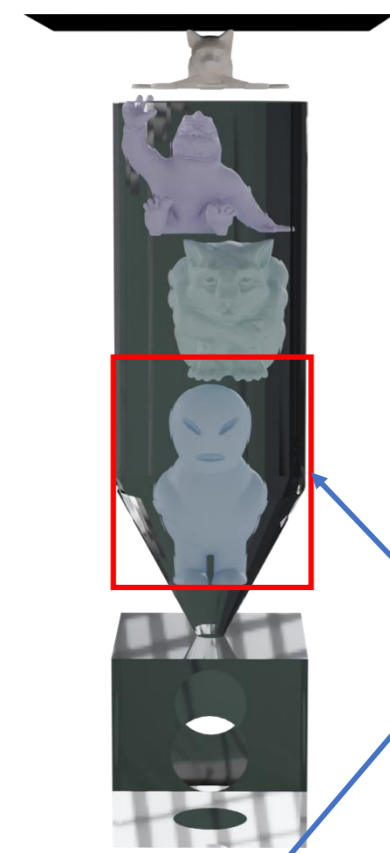
Frame 27: H_{10}

Unconstrained Newton Solver

```
For each frame:  
  While( $q_k$  is not optimal){  
     $d \leftarrow -H_k^{-1} g$   
    Compute step length  $t$   
     $q_{k+1} = q_k + t * d$   
     $H_{k+1} = \text{computeHessian}(H_k)$   
  }
```



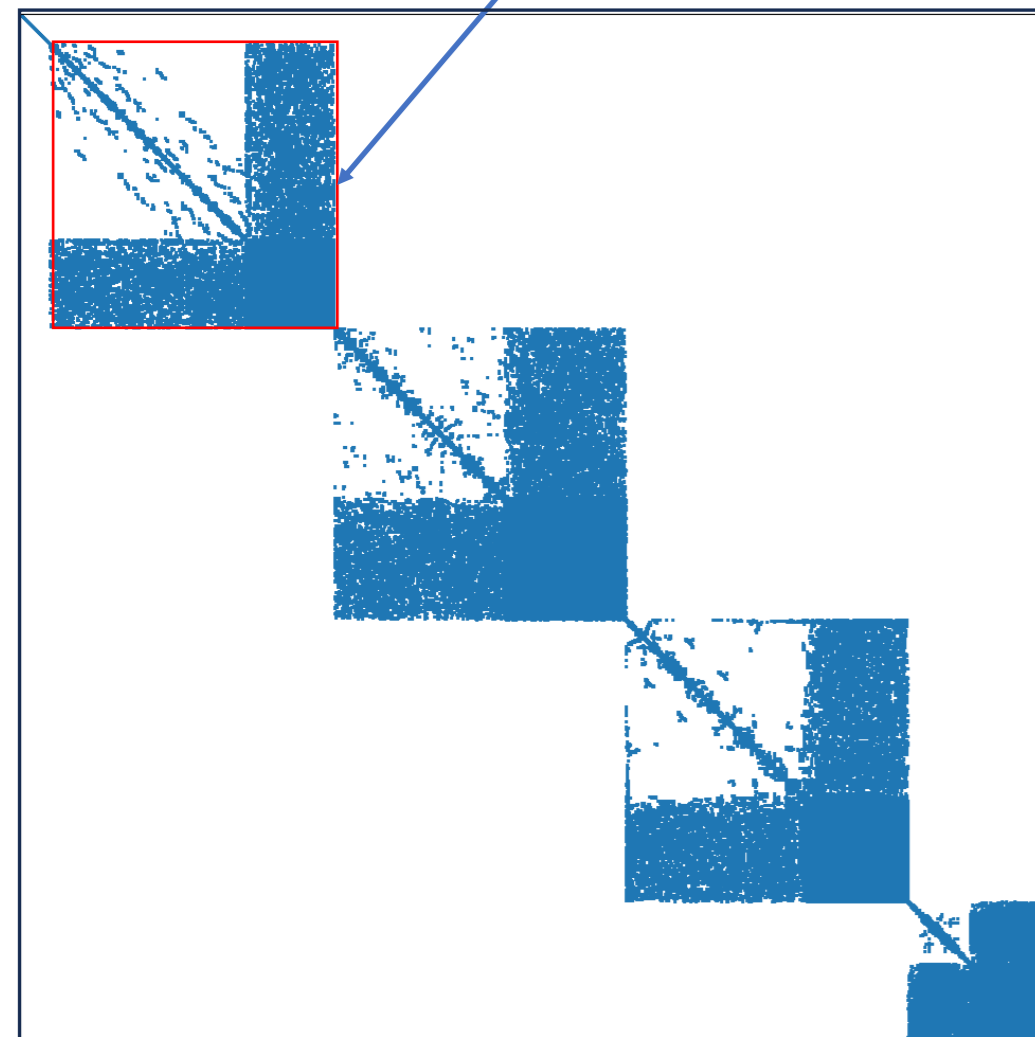
Moderate Changes to the Hessian Sparsity



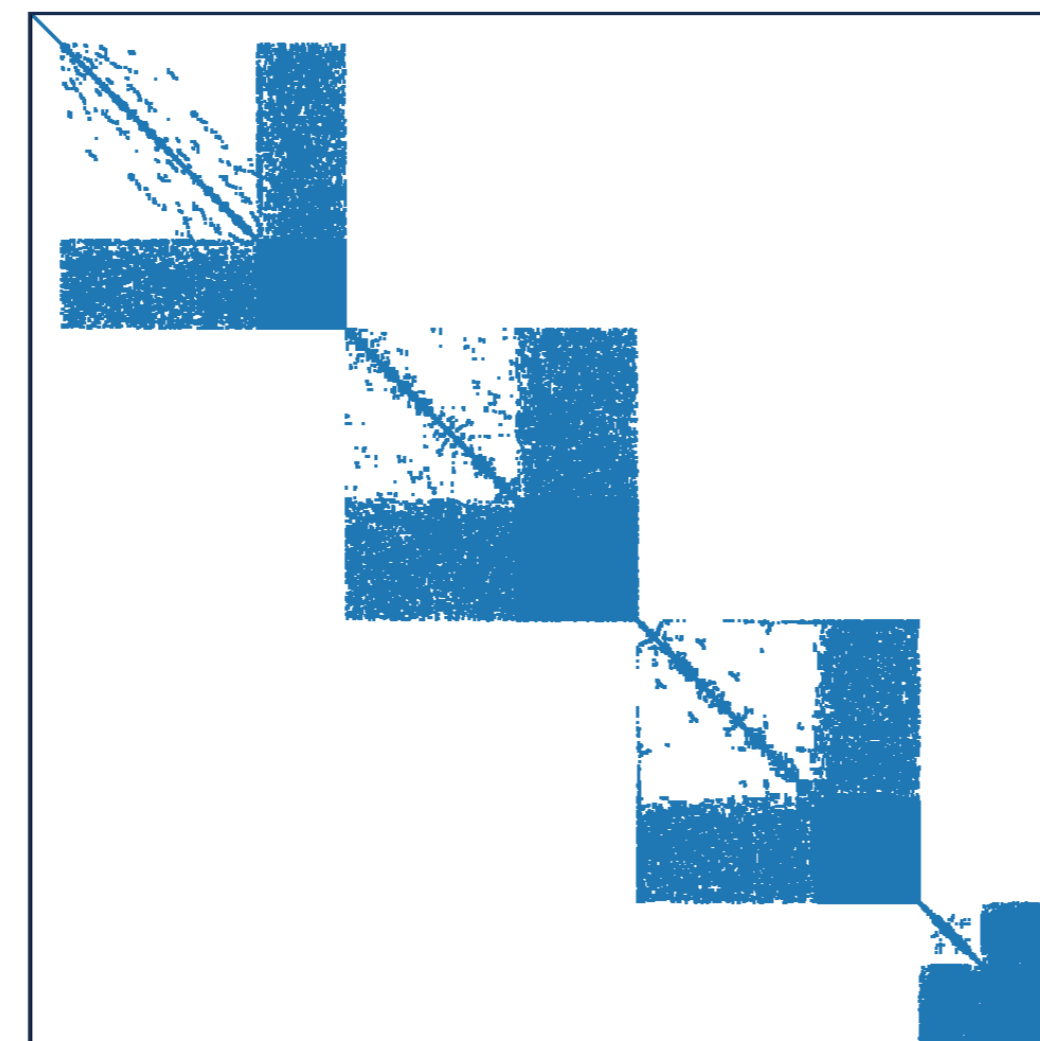
Hessian sparsity pattern comes from the geometry



Contact
→

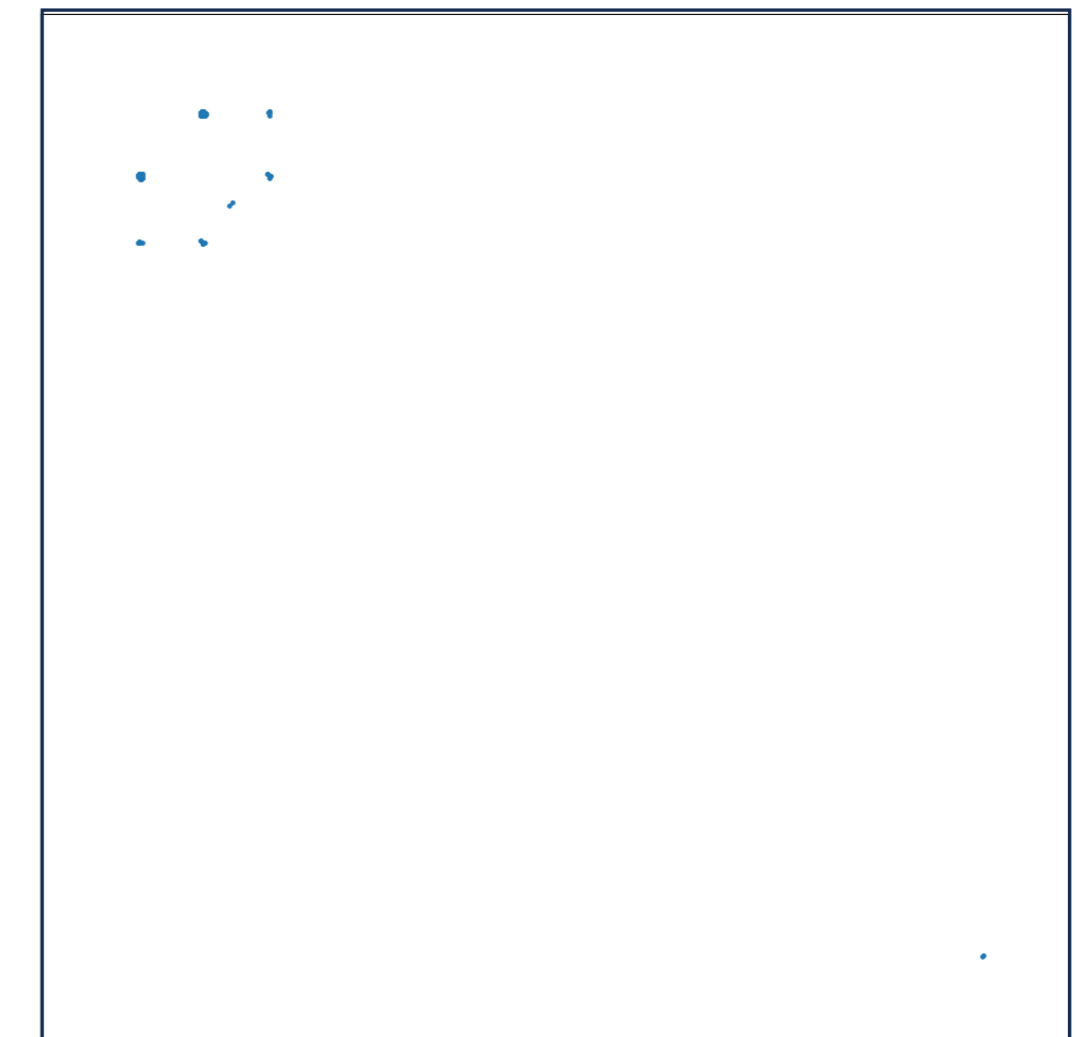


Frame 13: H_0



Frame 27: H_{10}

Difference in sparsity
→

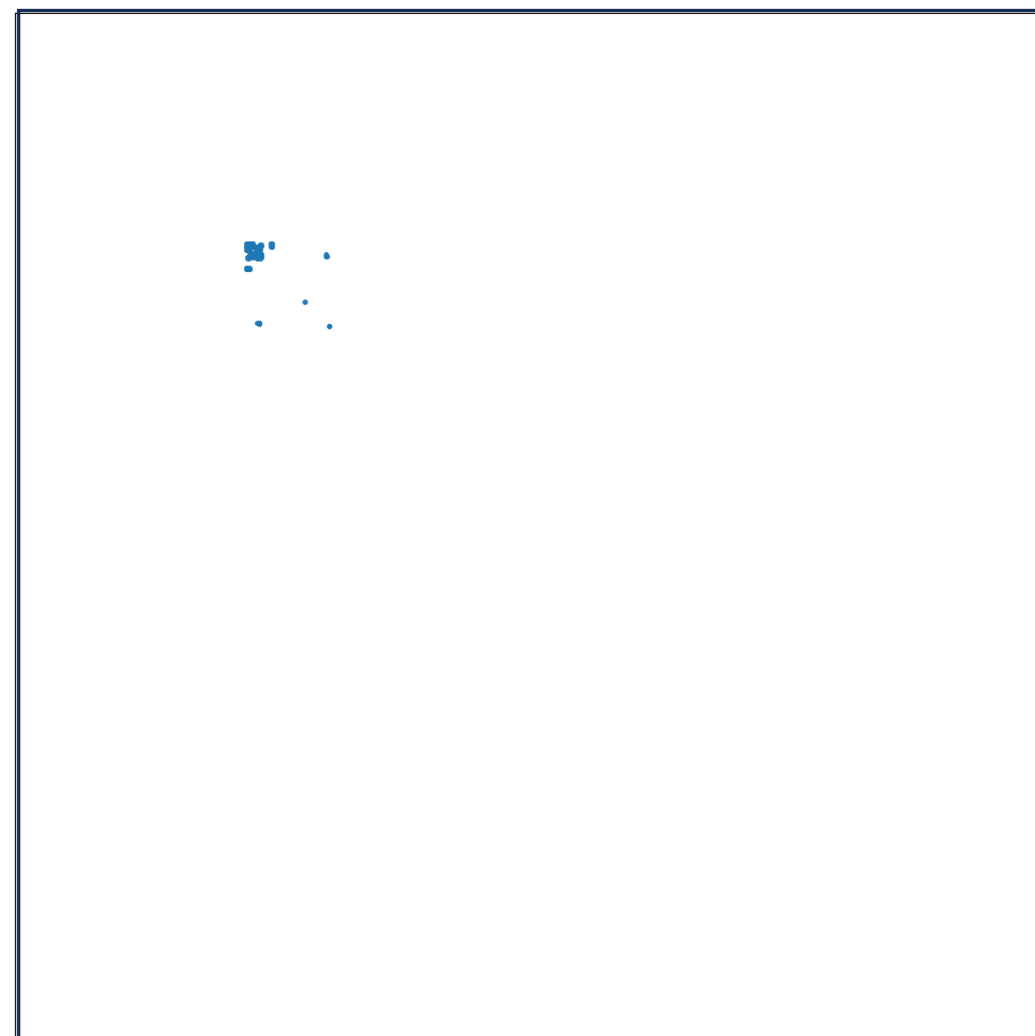


Gradual change in Sparsity

Nodes Involved in Contact

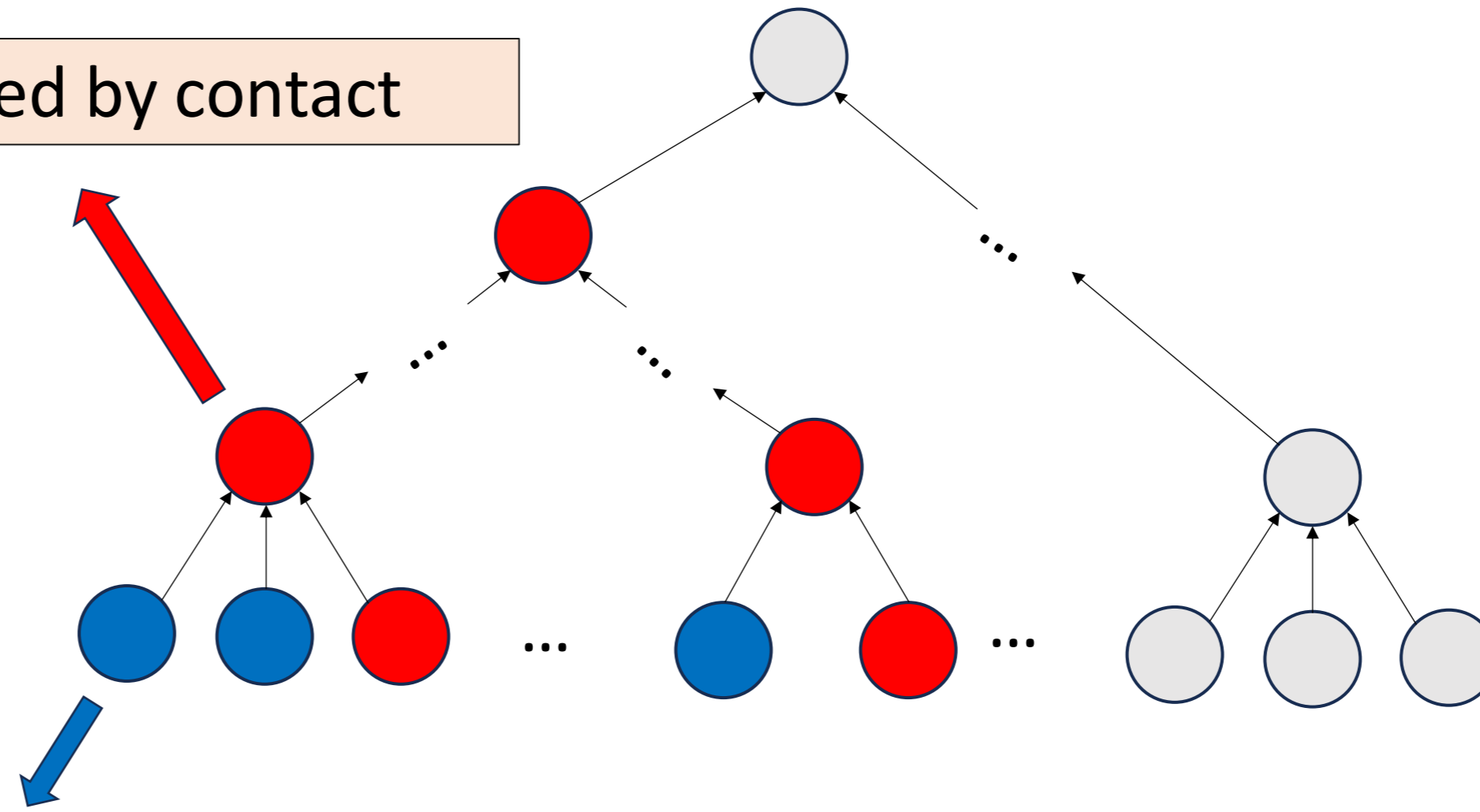


General Re-ordering

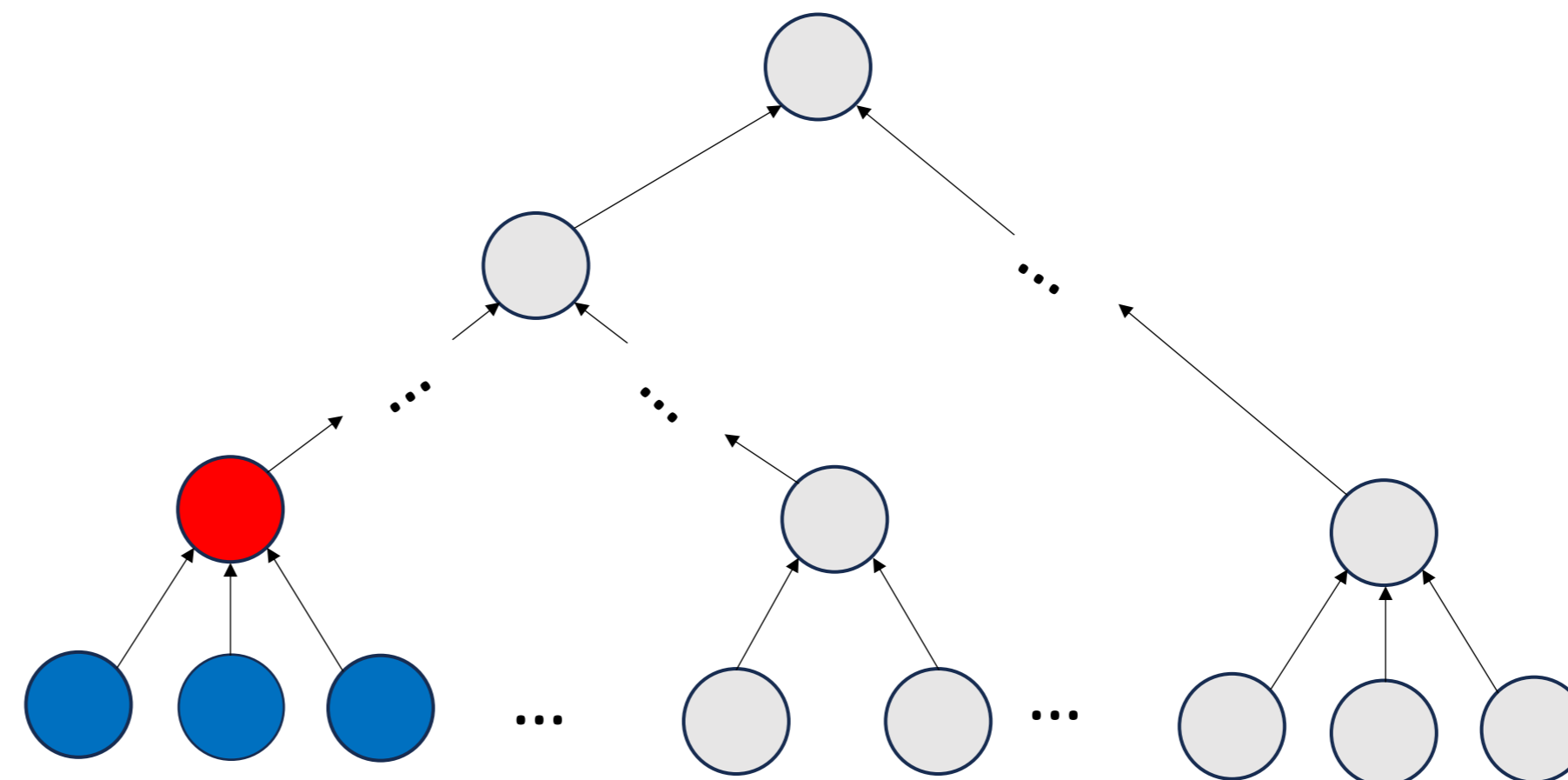


Parth (contact aware) Re-ordering

Nodes affected by contact



Nodes involved in contact

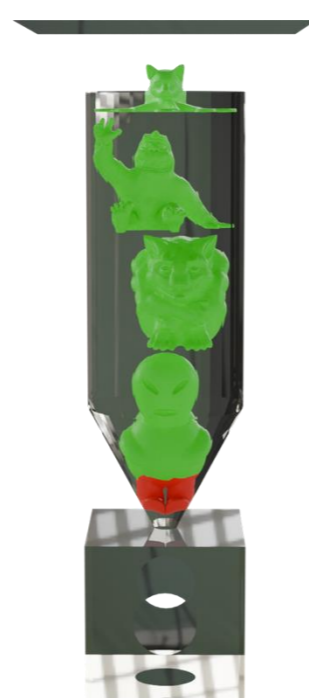


Data Dependence Graph

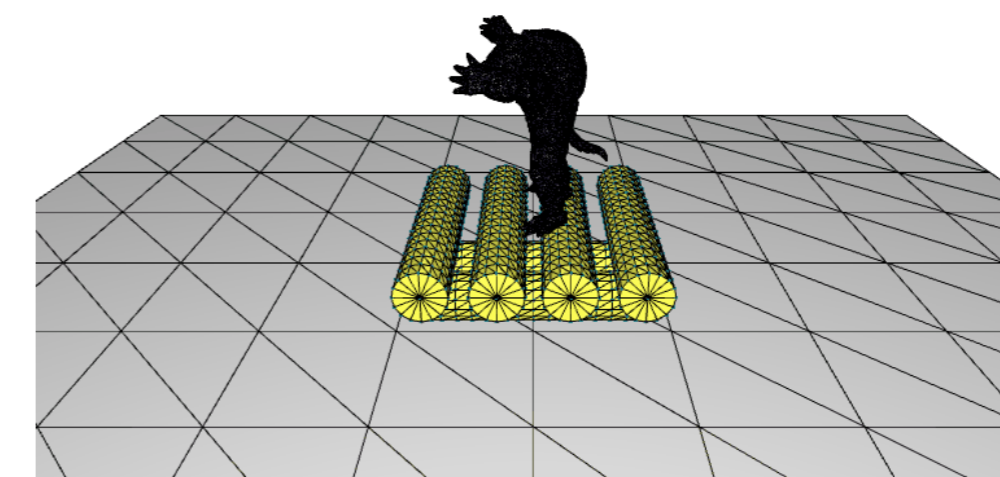
Parth: A Geometry and Re-ordering Aware Mesh Decomposition



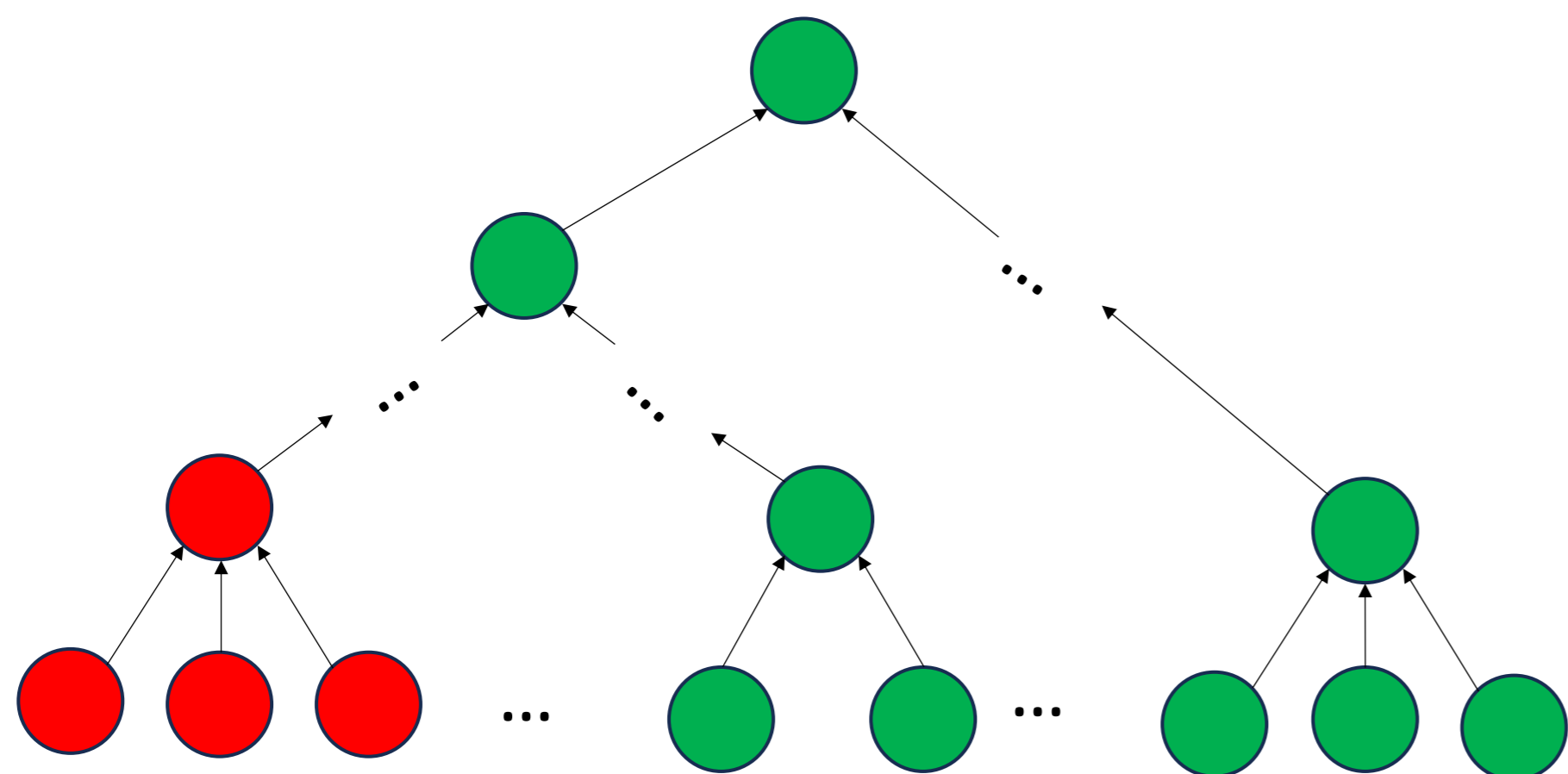
a) Parth decomposition



b) Parth localization



Arma Roller Simulation

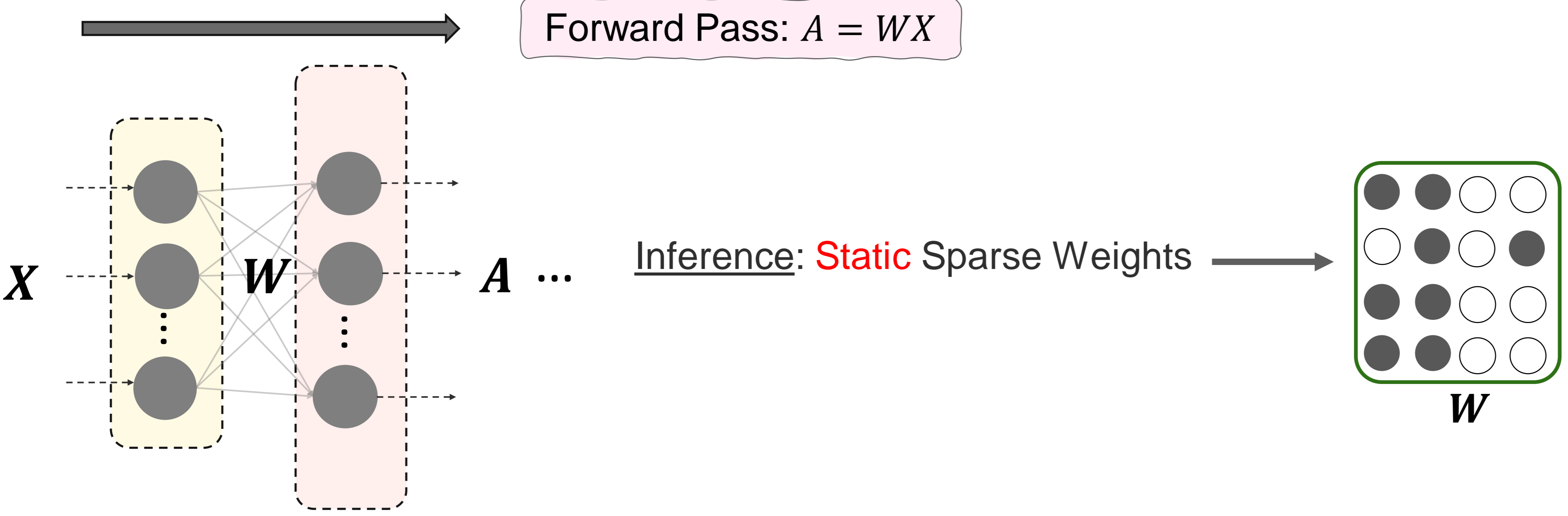


c) Data dependence graph with reuse (colored in green)

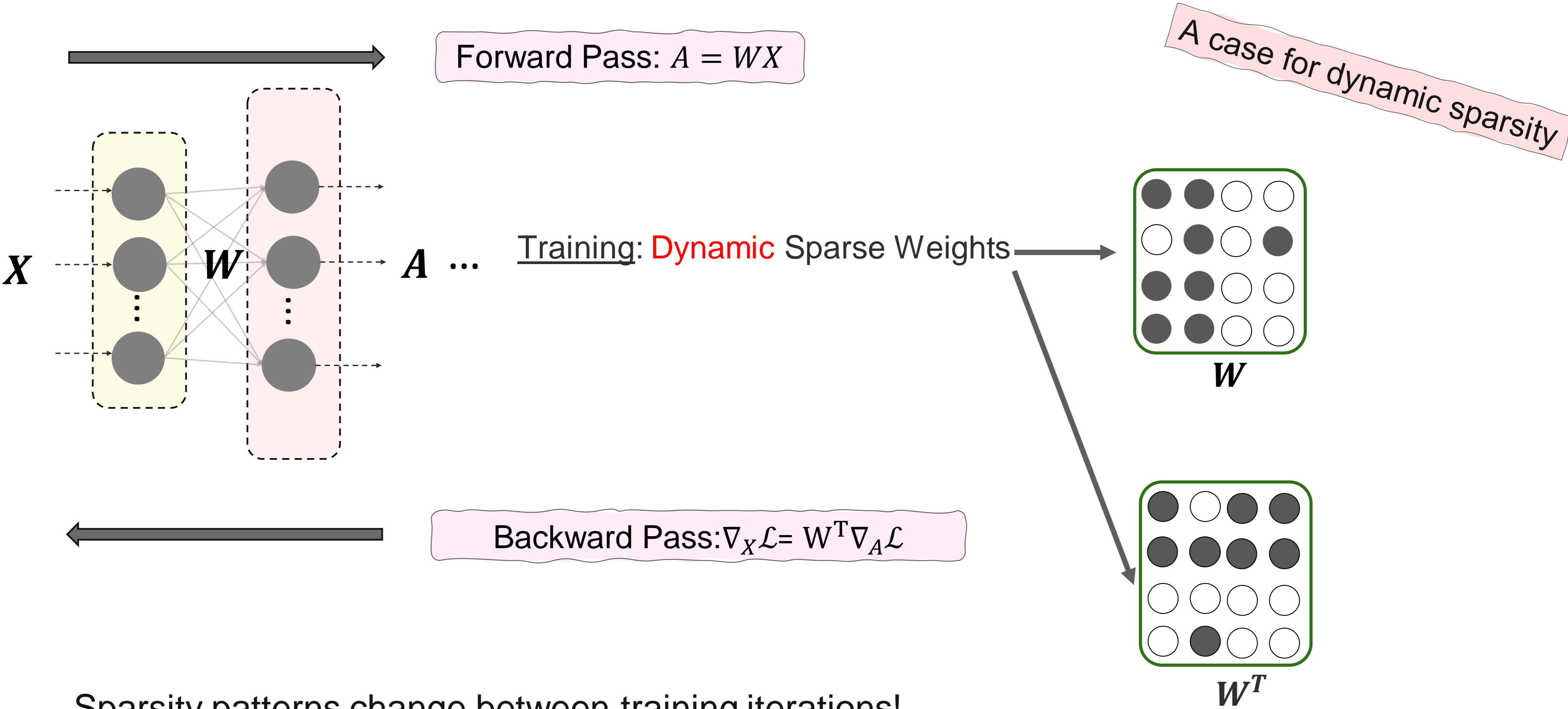
IPC Simulations	Apple Accelerate	Intel MKL	CHOLMOD
Dolphin Funnel	2.95x	2.15x	1.97x
Ball Mesh Roller	2.43x	1.87x	1.66x
Mat On Board	2.08x	1.93x	1.72x
Rods Twist	2.07x	1.87x	1.67x
Squeeze Out	2.62x	2.28x	1.99x
Arma Roller	1.33x	1.61x	1.55x

On Intel Xeon with 20 cores and M2 pro with 12 cores

Sparsity in the Neural Network Inference: Forward Pass

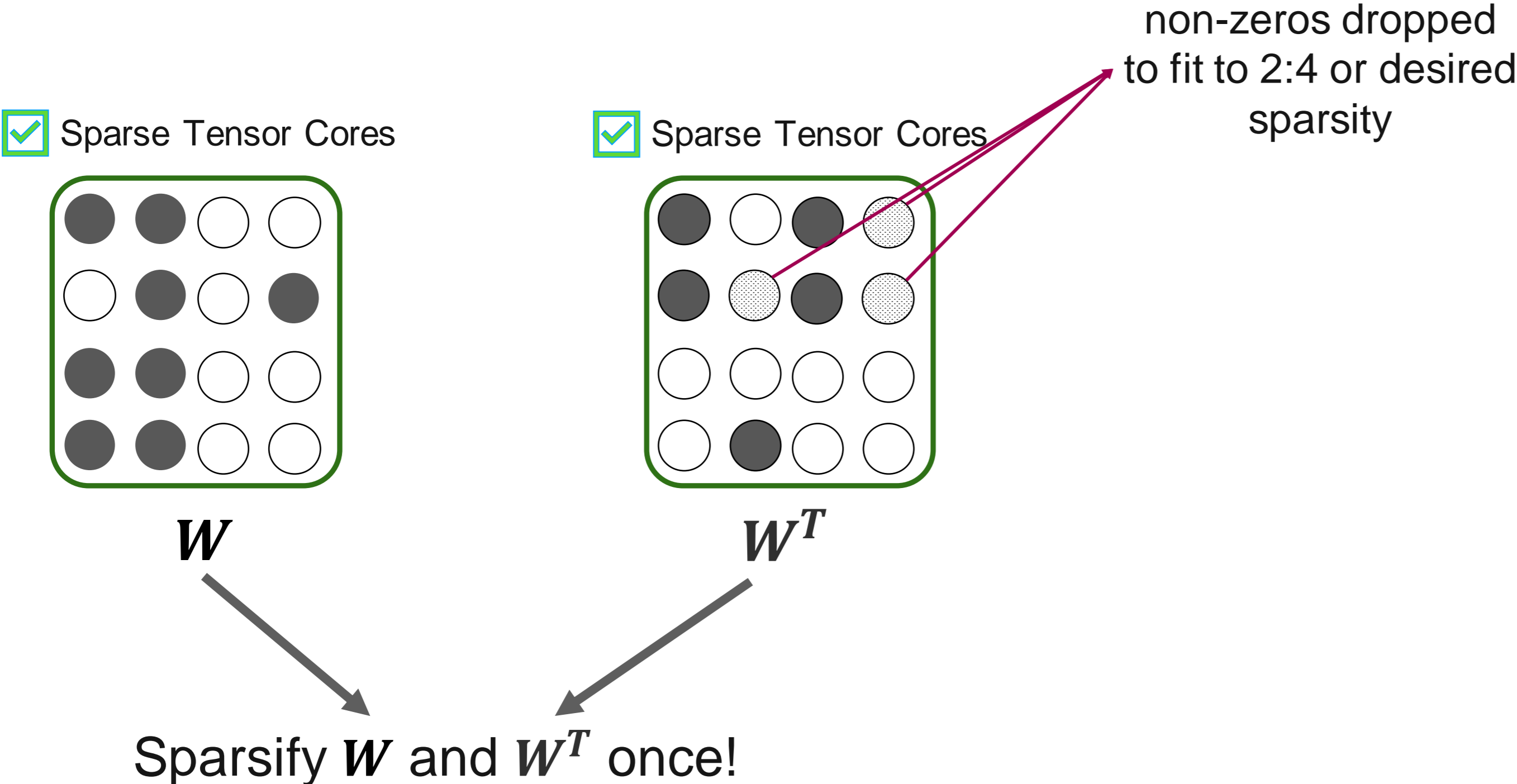


Sparsity in the Neural Network Training: Forward/Backward Pass



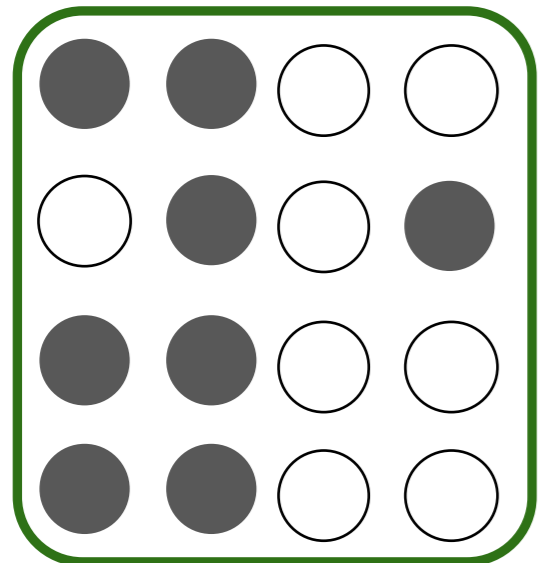
Sparsity patterns change between training iterations!

DNN Training Algorithms with Static Pruning



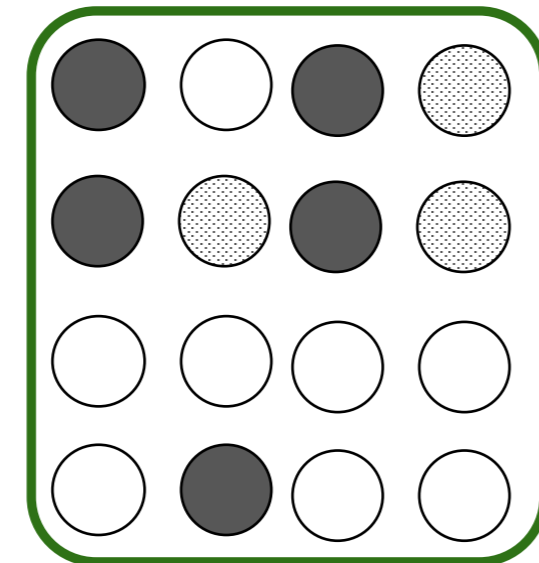
DNN Training Algorithms with Static Pruning

✓ Sparse Tensor Cores



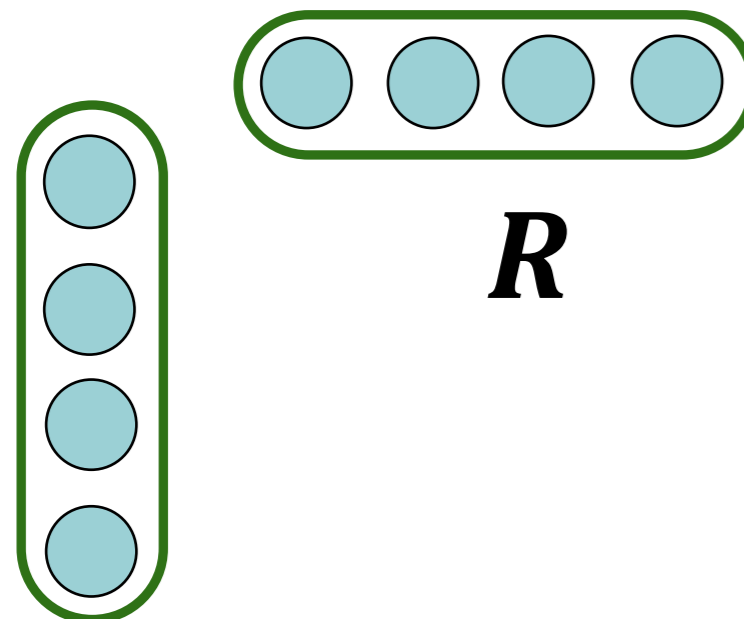
W

✓ Sparse Tensor Cores



W^T

+



R

Use low rank adaptors for accuracy

Forward Pass: $A = WX + LRX$

Model	Bert-Large	GPT2-Xlarge	GPT3-Large	GPT3-Xlarge	GPT3-2.7B
Speedup (x)	1.09	1.12	1.13	1.11	1.10

Dataset	Dense	$r = 0$	$r = 4$	$r = 16$	$r = 64$
SQuAD v1.1	90.4	89.1	89.1	89.2	89.5
GLUE	80.2	77.4	77.7	77.8	78.2

BERT-Large-Uncased Accuracy Results, r is rank

DNN Training Algorithms with Static Pruning

Second-Order optimizers at backward pass

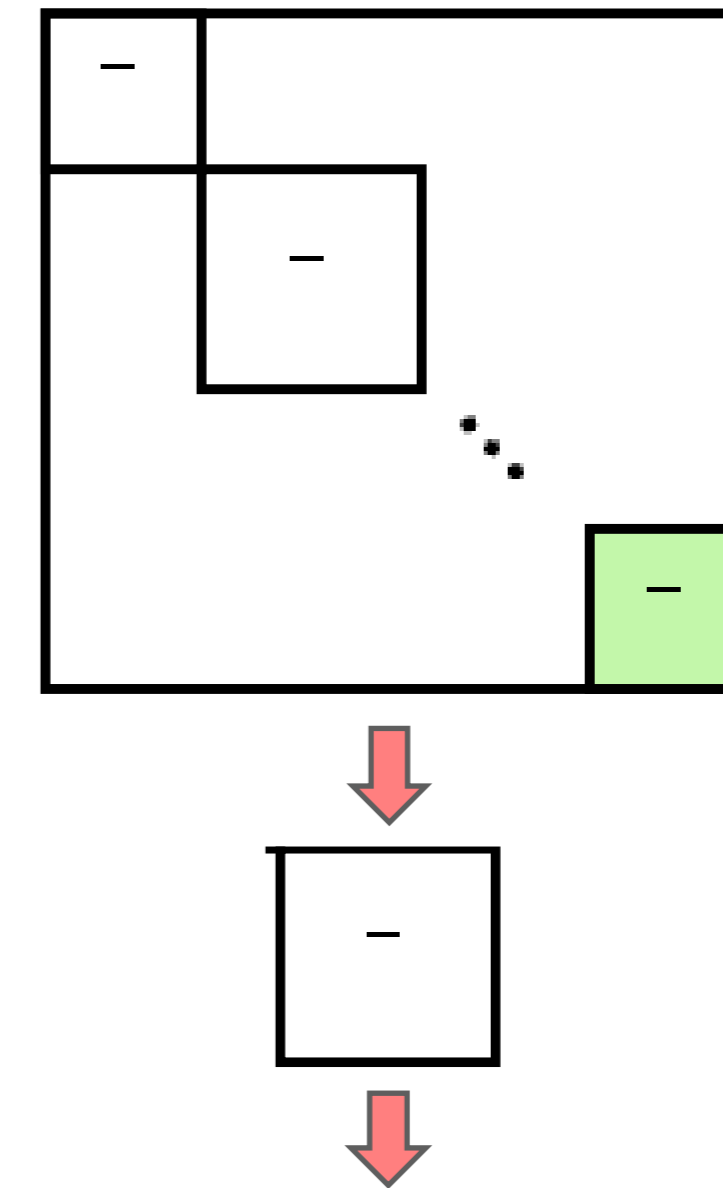
Stochastic Gradient Decent:

$$W \leftarrow W - \alpha \nabla \mathcal{L}(W)$$



Natural Gradient Decent:

$$W \leftarrow W - \alpha \mathbf{F}^{-1} \nabla \mathcal{L}(W)$$



Sparsification + quantization + rank-1 updates

MKOR reduces Bert training time vs. LAMB **2.57** times

Second-Order optimizers for backward pass to control sparsity patterns:

- *HyLo: A Hybrid Low-Rank Natural Gradient Descent Method*, Mu, Soori, Dehnavi [SC'22]
- *MKOR: Momentum-Enabled Kronecker-Factor-Based Optimizer Using Rank-1 Updates*, Mozaffari, Zhang, Dehnavi [NeurIPS'23]

Programmability!

Not another programming language!



Input to the Domain-Specific Compiler

Numerical Method

Sparse matrix-vector multiplication (SpMV)

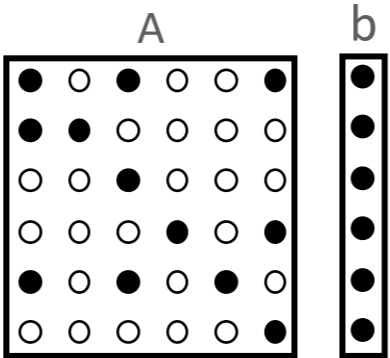
Sparsity info: Data Format, etc.

e.g. A is in Compressed Row (CSR) format

Programmability!

Input Program (NPB Benchmark)

```
for(j = 0; j < lastrow - firstrow + 1; j++){  
  sum = 0.0;  
  for(k = rowstr[j]; k < rowstr[j+1]; k++){  
    sum = sum + A[k]*b[colidx[k]];  
  }  
  y[j] = sum;  
}
```



Input to the Domain-Specific Compiler

Numerical Method

Sparse matrix-vector multiplication (SpMV)

≠

Sparsity info: Data Format, etc.

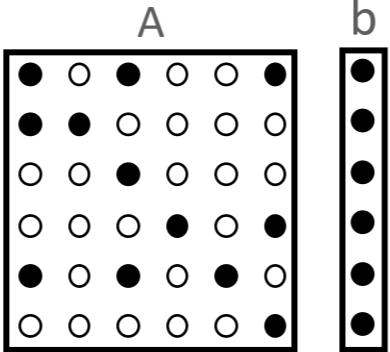
e.g. A is in Compressed Row (CSR) format



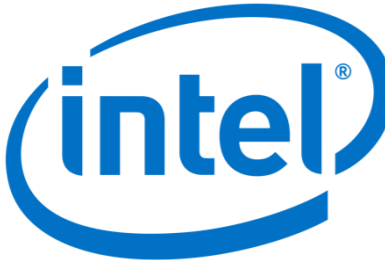
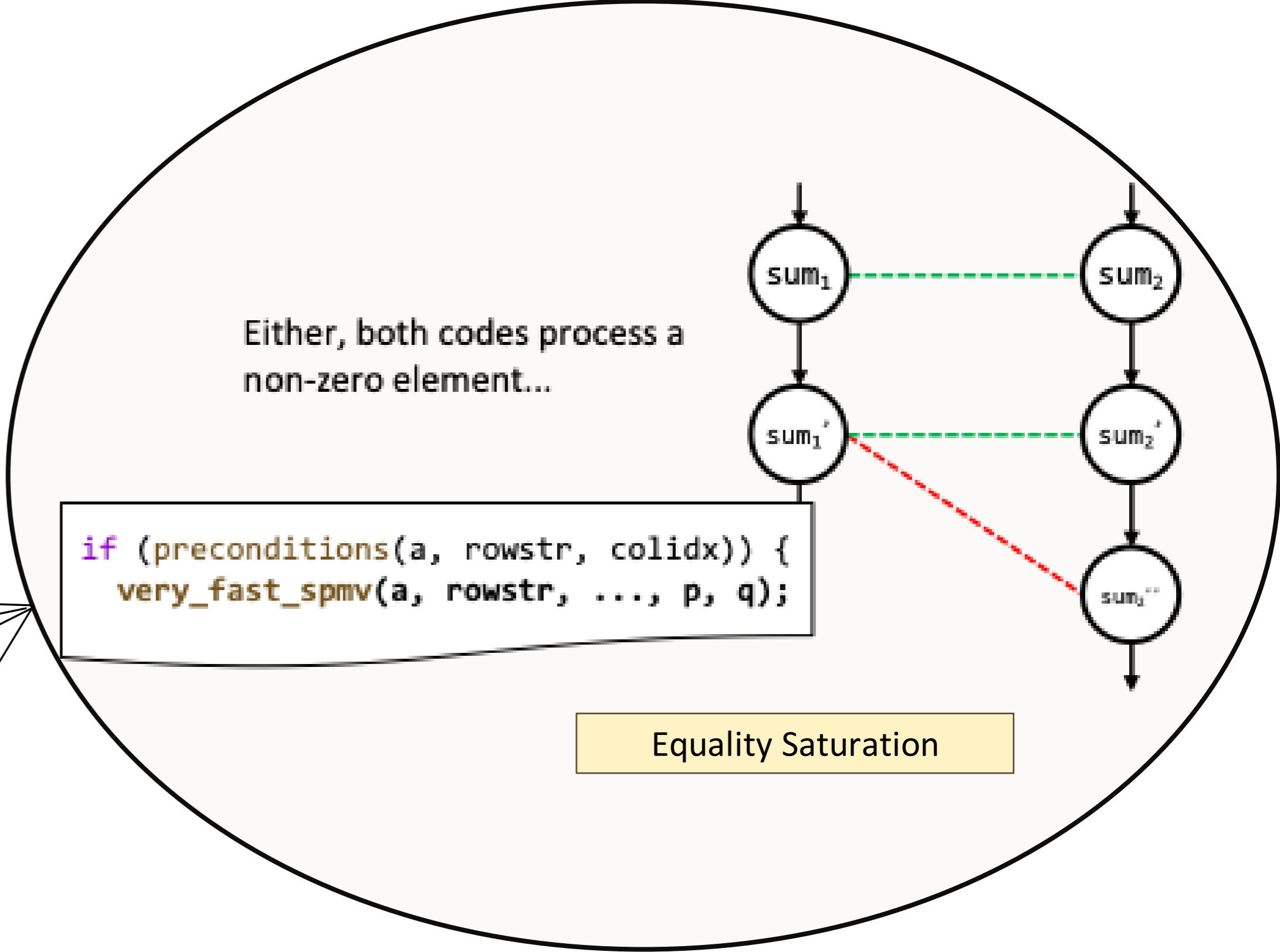
Automatically Translating Non-affine (e.g. sparse) Codes

Input Program (NPB Benchmark)

```
for(j = 0; j < lastrow - firstrow + 1; j++){  
  sum = 0.0;  
  for(k = rowstr[j]; k < rowstr[j+1]; k++){  
    sum = sum + A[k]*b[colidx[k]];  
  }  
  y[j] = sum;  
}
```



Automatic Translator



Backend

```
mk1_sparse_d_mv(...);
```



Backend

```
cusparseSpMV(...);
```

Backend

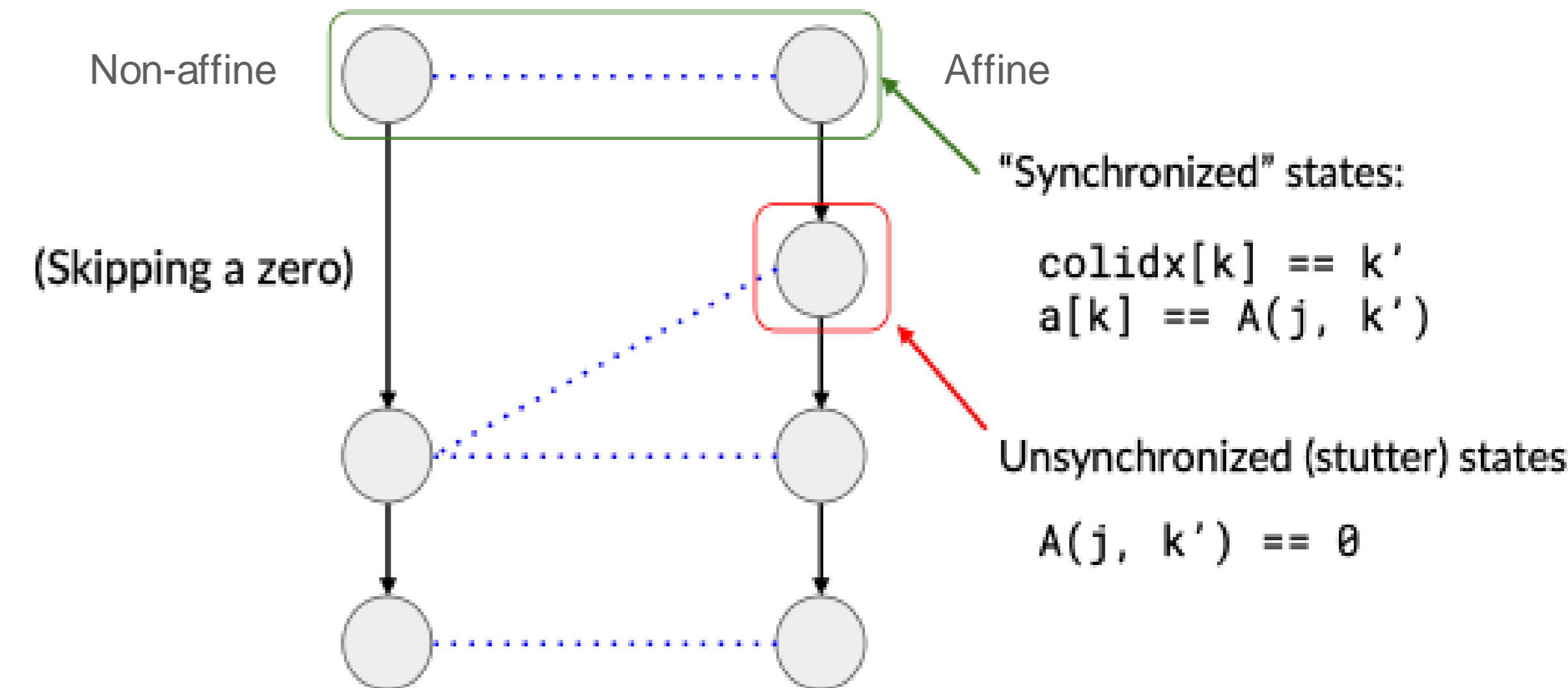


```
#pragma omp parallel for ...
```

DOMAIN SPECIFIC LANGUAGE (SYMPILER, TACO, REV,...)

Automatically Translating Non-affine Code

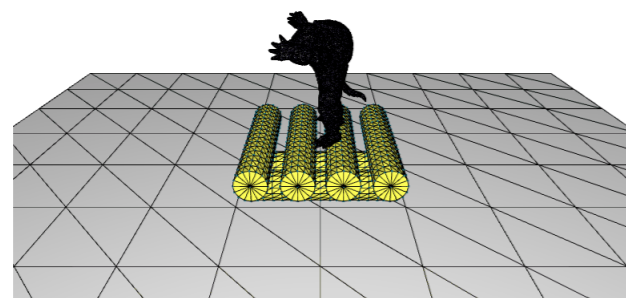
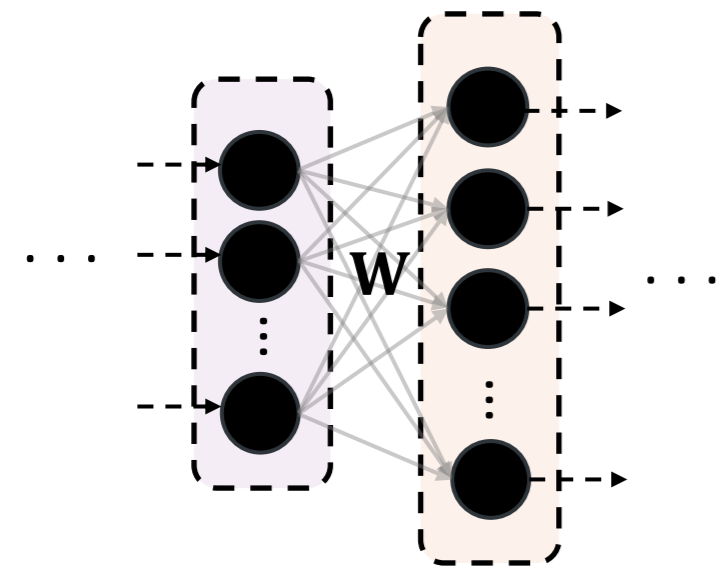
Rev: *Automatically translating* non-affine codes with branching bisimulation, *Laird, Liu, Bjorner, Dehnavi, to appear at [PLDI'24]*.



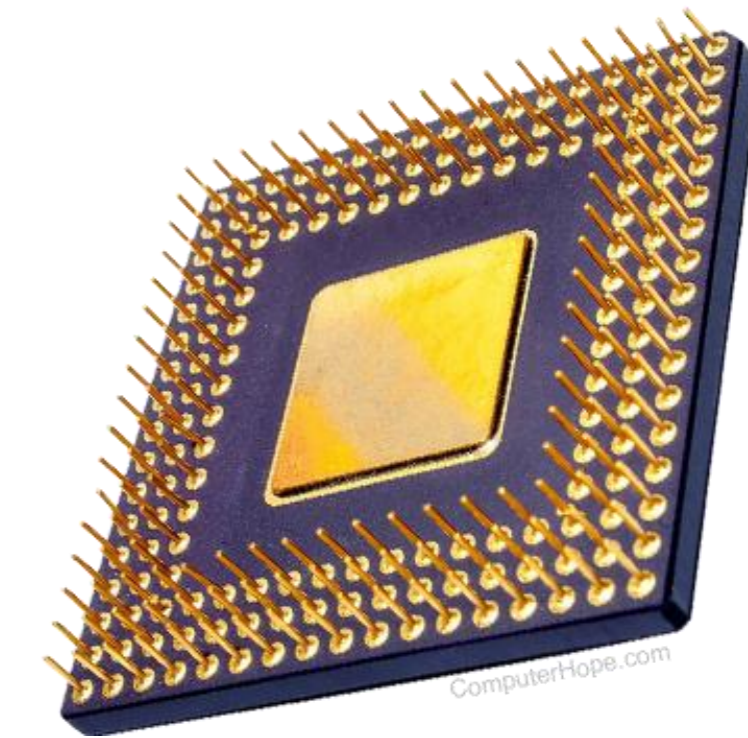
Sparse computation data dependence simplification for efficient *compiler-generated inspectors*, *Mohammadi, et. al.* [PLDI'19].

Take Aways

Sparse Algorithms & Applications

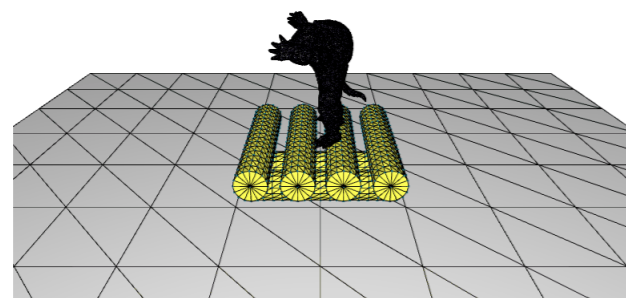
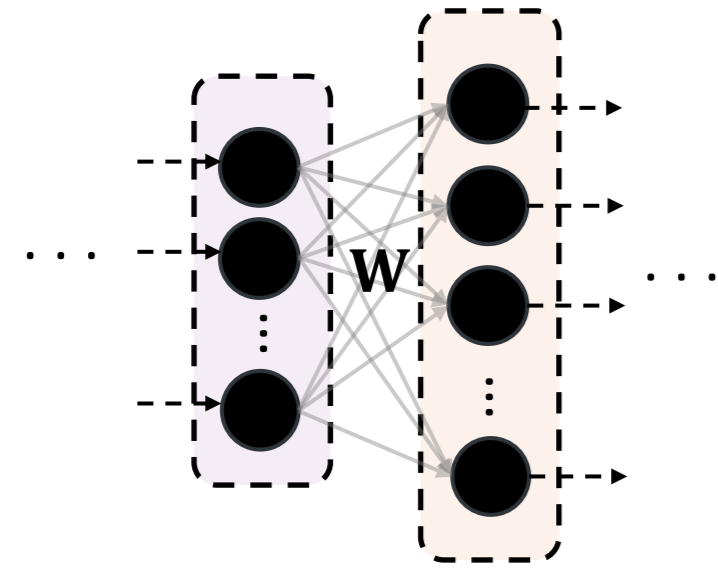


Hardware

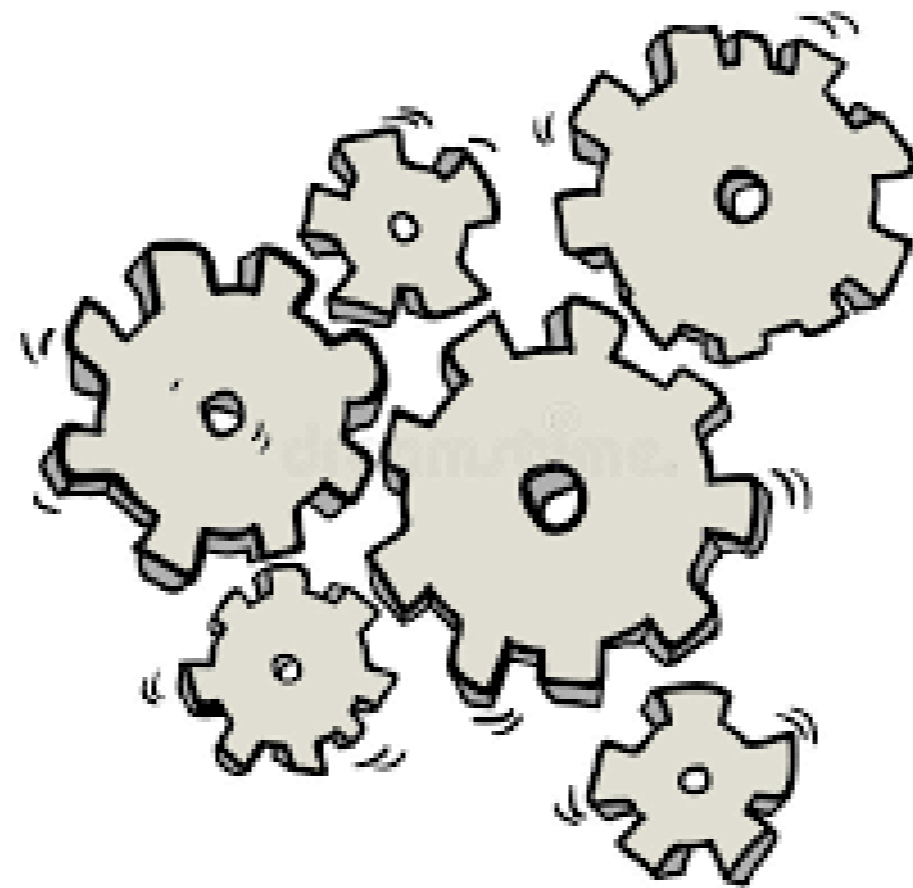


Take Aways

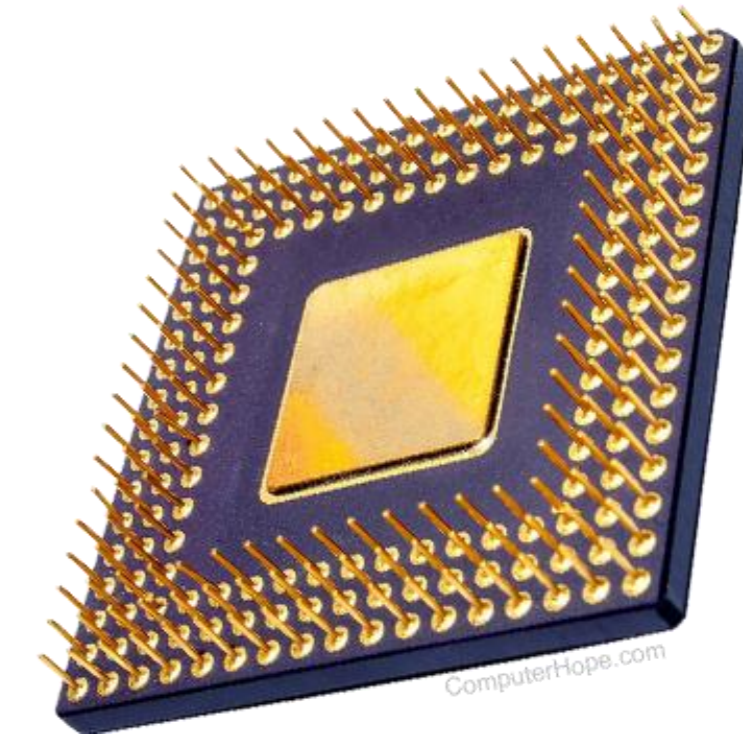
Sparse Algorithms & Applications



Compilers for sparsity



Hardware



Future!

Sparse Algorithms & Applications

