# DISTAL: Compiling Tensor Algebra to Distributed Machines
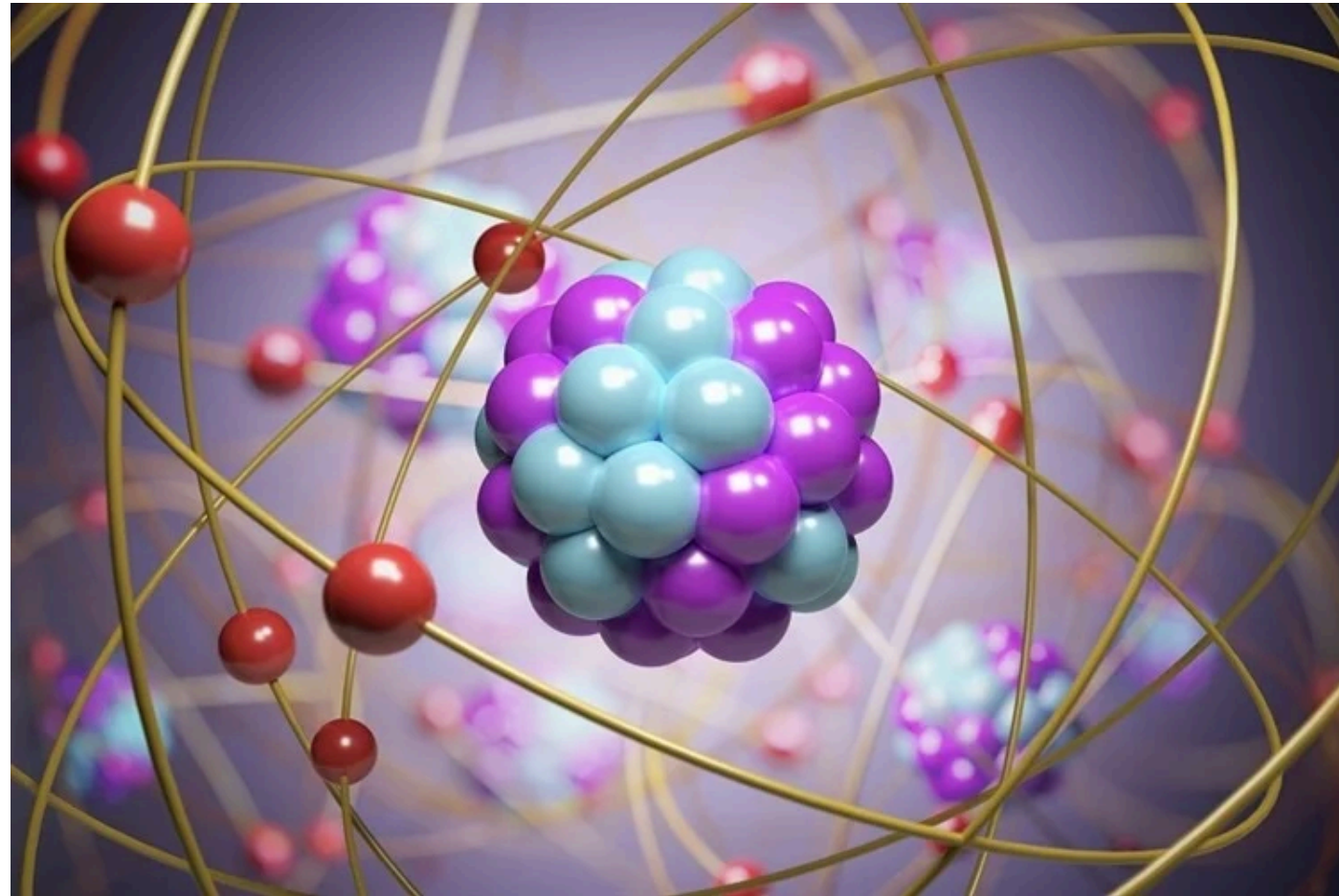
Rohan Yadav, Fred Kjolstad, Alex Aiken
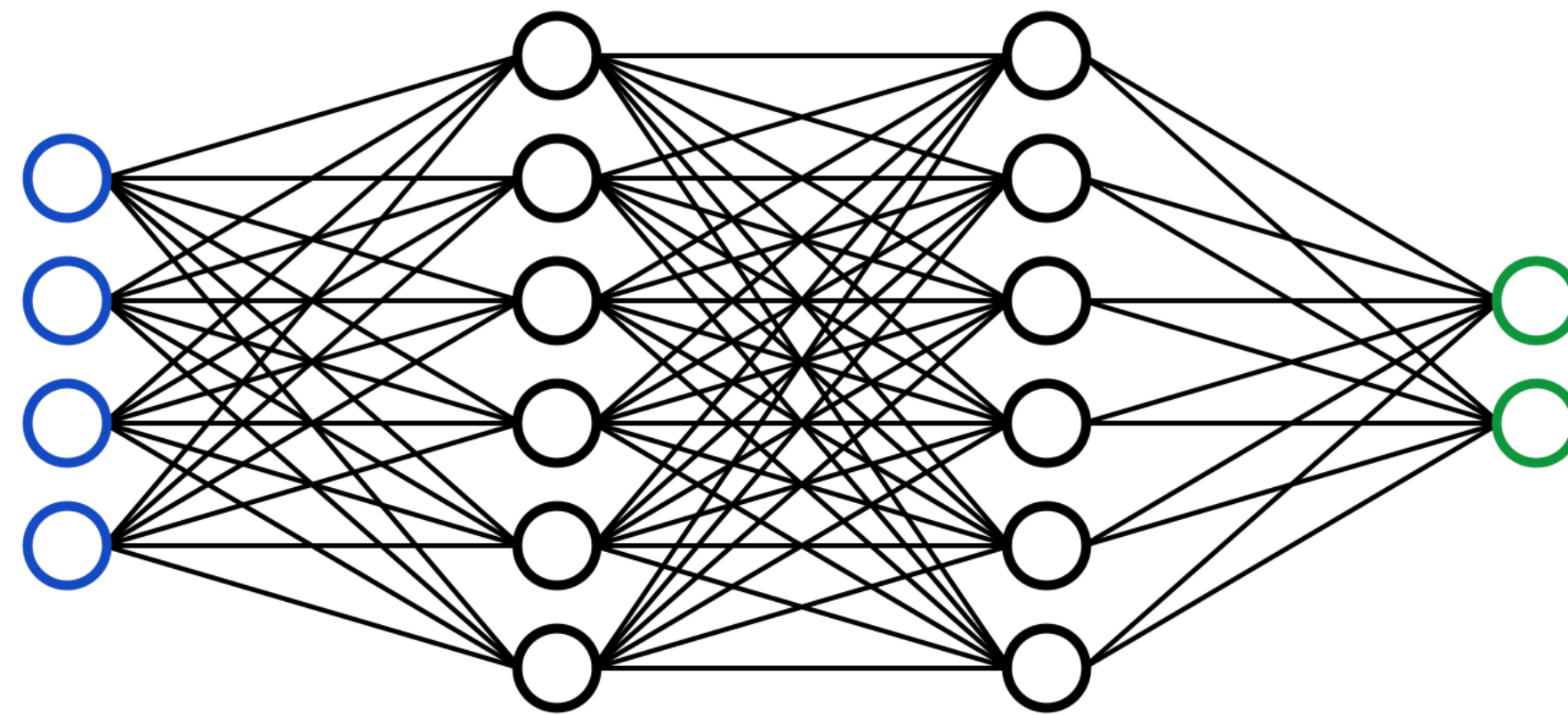
$$A(i, j) = B(i, j, k) \cdot c(k) \longrightarrow$$

# Tensor computations are ubiquitous



Scientific Computing



Machine Learning



Data Analytics

# Key Challenge:

Correctness  Productivity  Performance

# Why is achieving good performance hard?

## Optimizations are intertwined with correctness



## Performance doesn't compose



Energy Cost of Data Movement

1st key observation: **62.7%** of the total system energy is spent on **data movement**



## Variability of target architectures



4

# DSLs and compilation-based approaches have been successful

Halide, TVM, TACO (and many more!)

# How can compilers help?

Optimizations are intertwined with correctness

Performance doesn't compose

Variability of target architectures



*  Com...............ed
   abstr...........
*  Algo....................hrough
   *sche*...........
*  Optim.........................e
   corre........

*  ...............optimizations

Energy Cost of Data Movement

1st key observation: 62.7% of the total system energy is spent on data movement

*  ...................eddle for each
   b......

# What about distributed systems?

# Distributed systems exacerbate prior problems. We need compilers for this case too!

## Optimizations are intertwined with correctness

* Inter-address space communication

## Performance doesn't compose

* Expensive data movement
* Reorganize data to fit library interfaces

## Variability of target architectures

* Machines are heterogenous

# What are the right abstractions for distributed tensor compilation?

Simpler for end users

Capture existing algorithms

Generalize to all tensor programs

# DISTAL: The Distributed Tensor Algebra Compiler

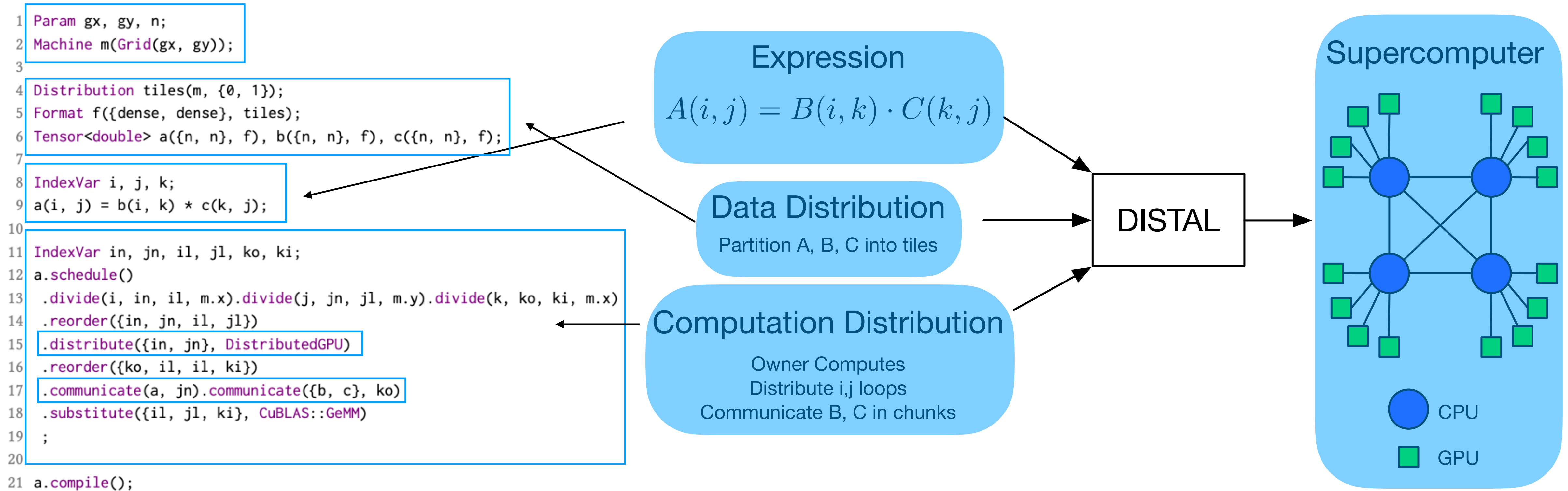Decouple computation, performance optimizations, and data distribution

```
1  Param gx, gy, n;
2  Machine m(Grid(gx, gy));
3
4  Distribution tiles(m, {0, 1});
5  Format f({dense, dense}, tiles);
6  Tensor<double> a({n, n}, f), b({n, n}, f), c({n, n}, f);
7
8  IndexVar i, j, k;
9  a(i, j) = b(i, k) * c(k, j);
10
11 IndexVar in, jn, il, jl, ko, ki;
12 a.schedule()
13  .divide(i, in, il, m.x).divide(j, jn, jl, m.y).divide(k, ko, ki, m.x)
14  .reorder({in, jn, il, jl})
15  .distribute({in, jn}, DistributedGPU)
16  .reorder({ko, il, il, ki})
17  .communicate(a, jn).communicate({b, c}, ko)
18  .substitute({il, jl, ki}, CuBLAS::GeMM)
19  ;
20
21 a.compile();
```

**Expression**

$$A(i,j) = B(i,k) \cdot C(k,j)$$

**Data Distribution**

Partition A, B, C into tiles

**Computation Distribution**

Owner Computes
Distribute i,j loops
Communicate B, C in chunks

DISTAL

**Supercomputer**

CPU

GPU

# Modeling Machines

### Expression

$$A(i,j) = B(i,k) \cdot C(k,j)$$

$$A(i,l) = B(i,j,k) \cdot C(j,l) \cdot D(k,l)$$

$$a = B(i,j,k) \cdot C(i,j,k)$$

$$A(i,j,l) = B(i,j,k) \cdot C(k,l)$$

$$A(i,j) = B(i,j,k) \cdot c(k)$$

### Data Distribution

Partition A into tiles
Replicate B onto all nodes
Place C onto only some nodes

### Computation Distribution

Owner Computes
Distribute i,j loops
Communicate in chunks

DISTAL

## Supercomputer

CPU

GPU

View machines as hyper-rectangular grids of processors, where each processor has a local memory

Expose locality in the physical machine

Structure machine like target computations

# Distributing Data

## Expression

$$A(i, j) = B(i, k) \cdot C(k, j)$$

$$A(i, l) = B(i, j, k) \cdot C(j, l) \cdot D(k, l)$$

$$a = B(i, j, k) \cdot C(i, j, k)$$

$$A(i, j, l) = B(i, j, k) \cdot C(k, l)$$
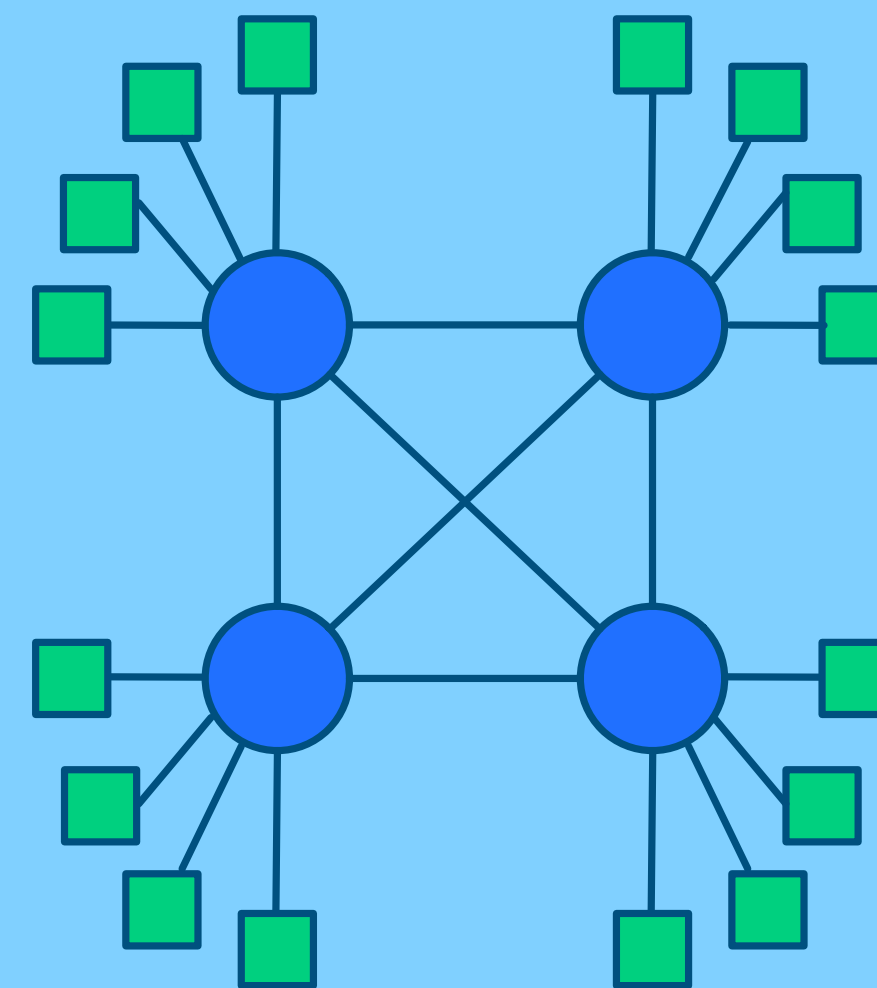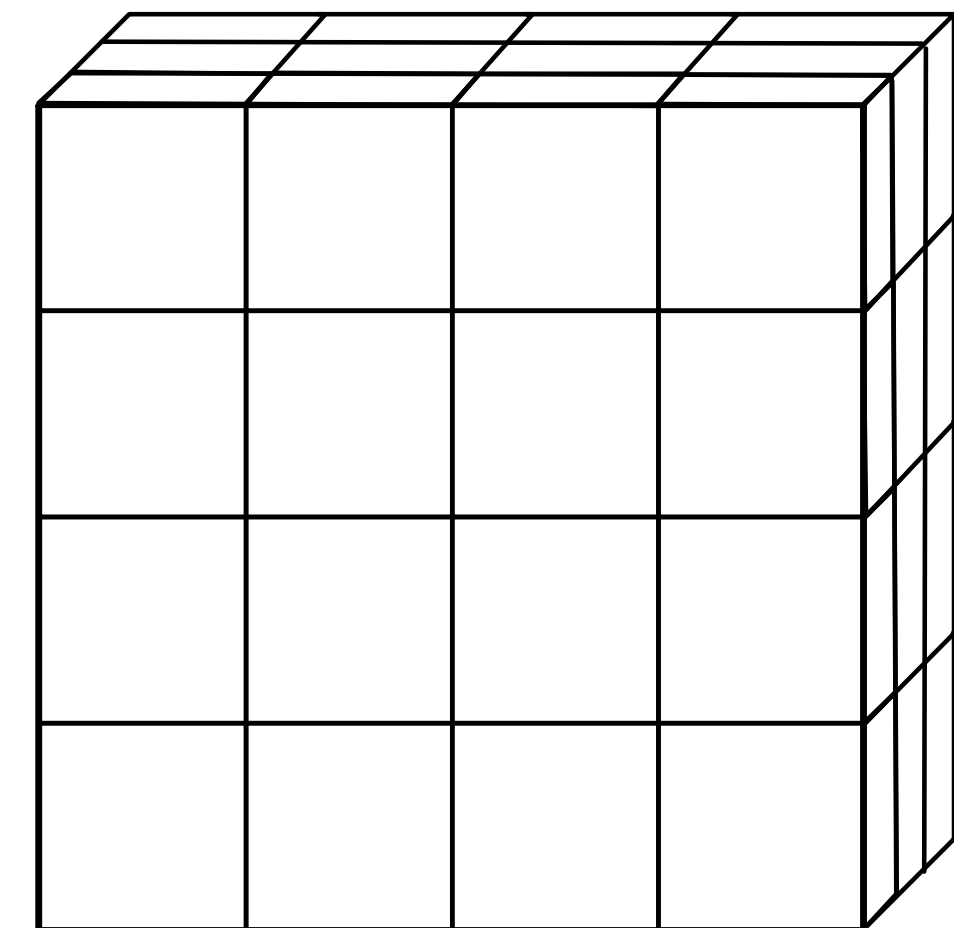
$$A(i, j) = B(i, j, k) \cdot c(k)$$
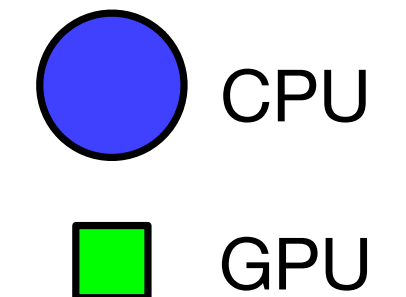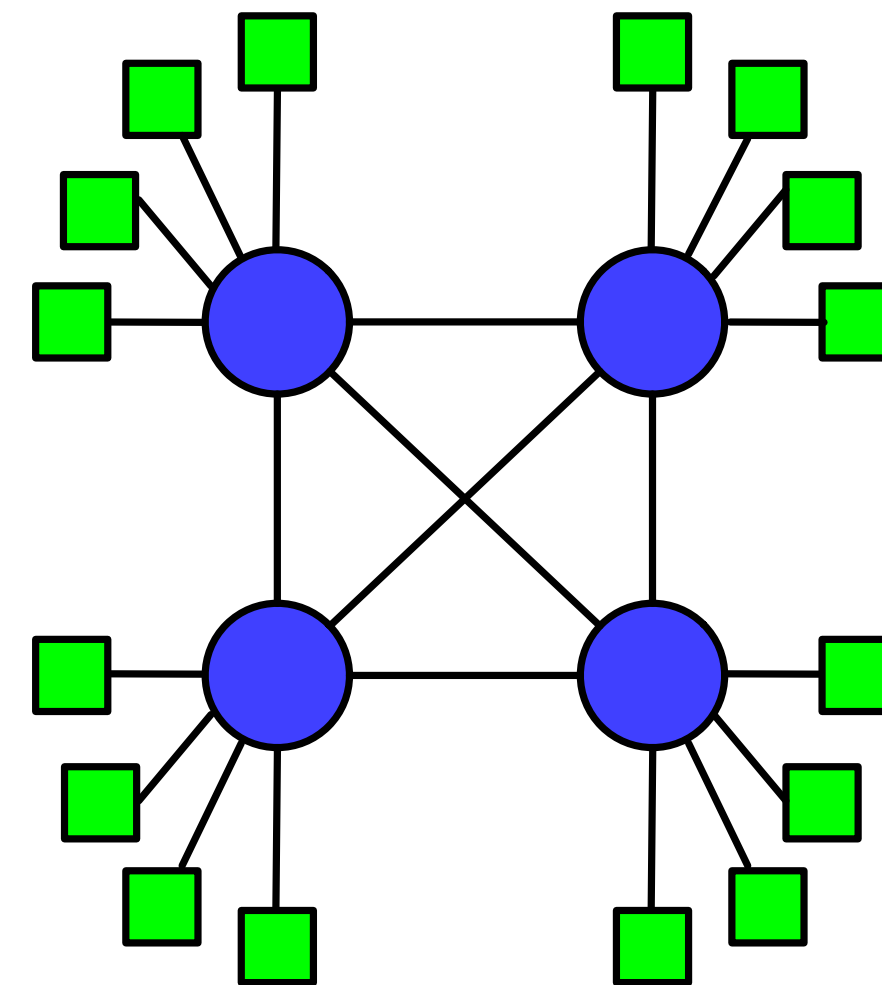
**Data Distribution**
Partition A into tiles
Replicate B onto all nodes
Place C onto only some nodes

## Computation Distribution
Owner Computes
Distribute i,j loops
Communicate in chunks

DISTAL

## Supercomputer



CPU

GPU

# Tensor Distribution Notation

Describe how dimensions of a tensor $\mathcal{T}$ map onto a machine $\mathcal{M}$

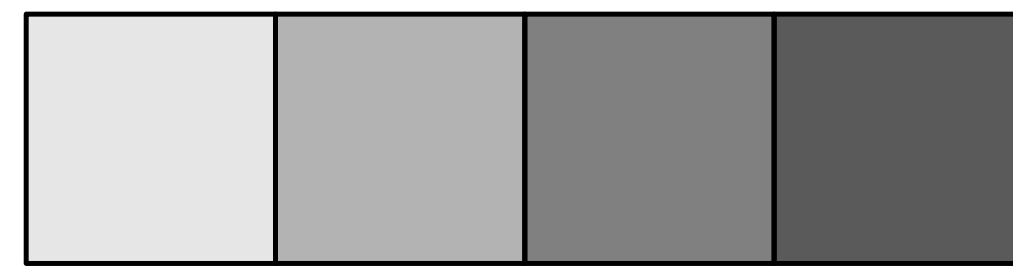$$\mathcal{T}_{\mathrm{xy}} \mapsto_{\mathrm{x}} \mathcal{M}$$

Name each dimension of $\mathcal{T}$ and $\mathcal{M}$

Dimensions of $\mathcal{T}$ are partitioned by dimensions of $\mathcal{M}$ with the same name
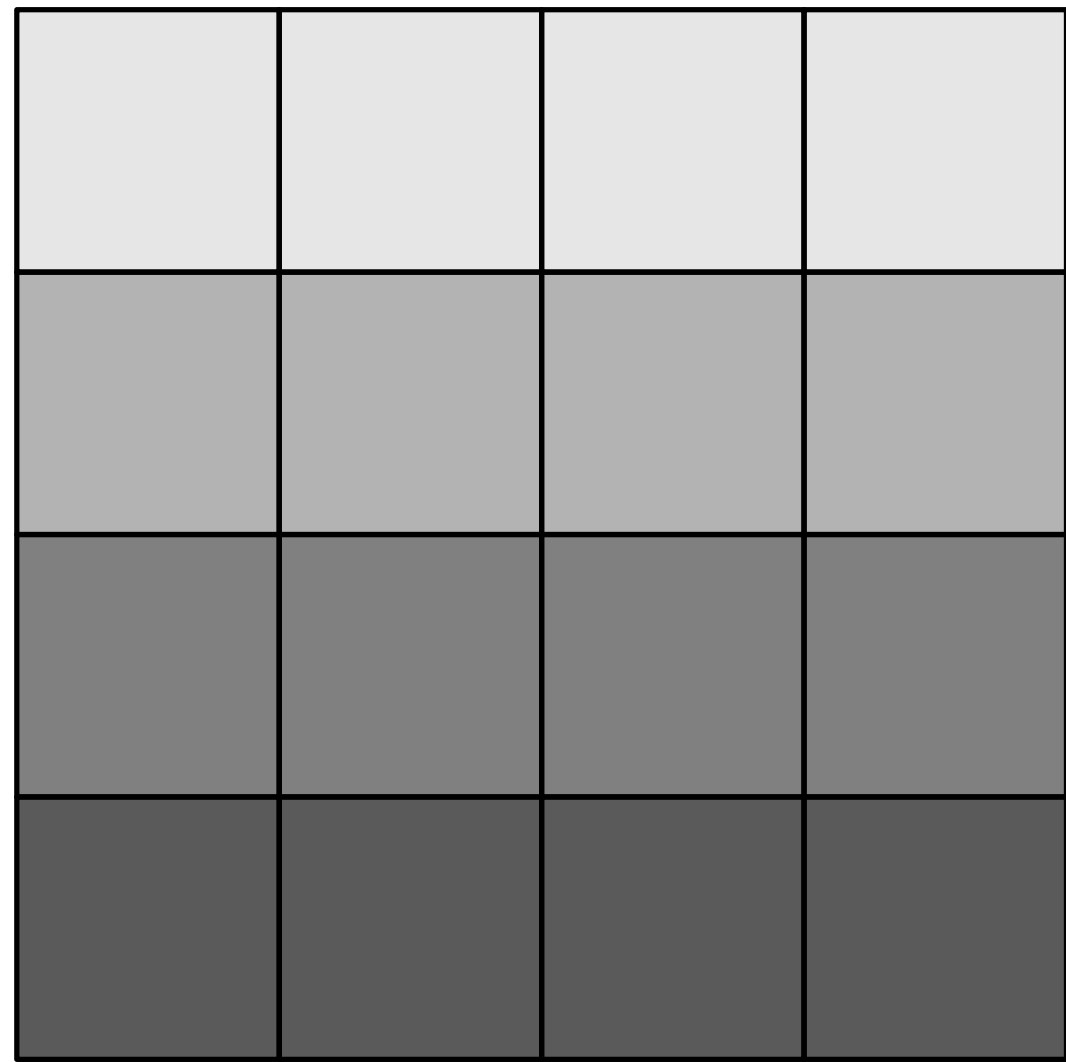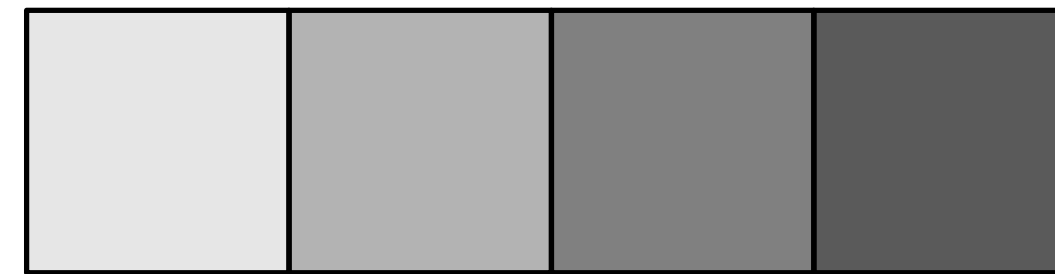
$$\mathcal{T} \;_{\mathrm{x}}\!\!\longmapsto_{\mathrm{x}} \mathcal{M}$$

$$\mathcal{T}$$

$$\mathcal{M}$$

$$\mathcal{T} \underset{\mathrm{xy} \mapsto \mathrm{x}}{} \mathcal{M}$$

$\mathcal{T}$

$\mathcal{M}$

$$\mathcal{T} \, _{\mathrm{xy}} \!\mapsto_{\mathrm{xy}} \mathcal{M}$$

$\mathcal{T}$

$\mathcal{M}$

$$\mathcal{T}_{\mathrm{xyz}\mapsto\mathrm{xy}}\mathcal{M}$$

$\mathcal{T}$

$\mathcal{M}$

$$\mathcal{T}_{\text{xy}} \mapsto_{\text{xy}*} \mathcal{M}$$

$\mathcal{T}$

$\mathcal{M}$

$$\mathcal{T} \underset{\mathrm{xy} \mapsto \mathrm{xy0}}{} \mathcal{M}$$

$$\mathcal{T}$$

$$\mathcal{M}$$

# Distributing Computation

Expression

$$A(i,j) = B(i,k) \cdot C(k,j)$$

$$A(i,l) = B(i,j,k) \cdot C(j,l) \cdot D(k,l)$$

$$a = B(i,j,k) \cdot C(i,j,k)$$

$$A(i,j,l) = B(i,j,k) \cdot C(k,l)$$

$$A(i,j) = B(i,j,k) \cdot c(k)$$

Data Distribution

Partition A into tiles
Replicate B onto all nodes
Place C onto only some nodes

Computation Distribution
Owner Computes
Distribute i,j loops
Communicate in chunks

DISTAL

Supercomputer



⬤ CPU

🟩 GPU

# Iteration Spaces
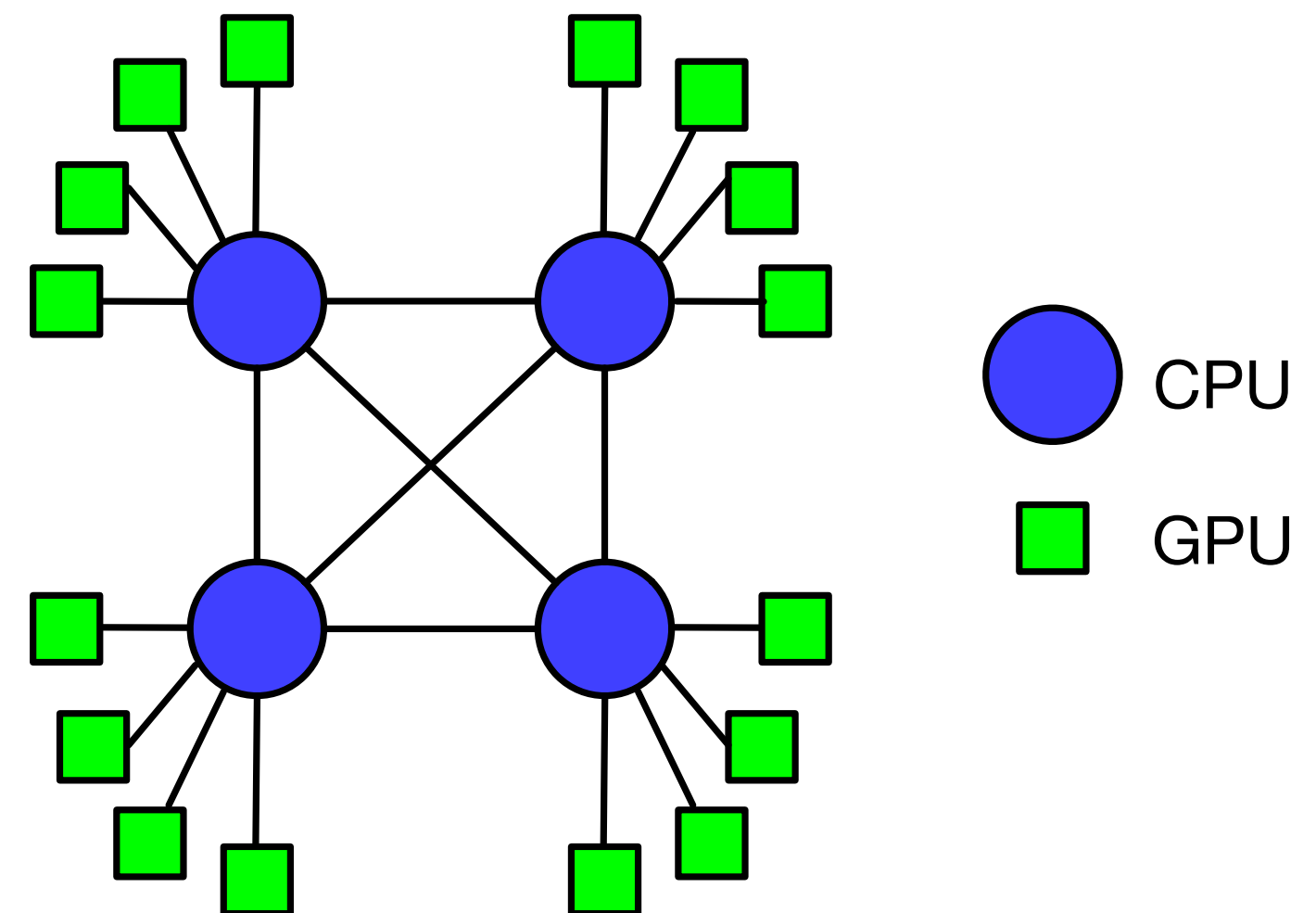
Hyper-rectangular grid of points representing each point in a set of nested loops

$$\forall_i \; a(i) = \sum_j b(j) \;\; (i,j) \in \{0, 1, 2\} \times \{0, 1, 2\}$$

```
for i in (0, len(a)):
  for j in (0, len(b)):
    a[i] += b[j]
```

# Execution Spaces

Space of all processors in $\mathcal{M}$ cross a time dimension

$$\forall_i \ a(i) = \sum_j b(j) \ \ (i,j) \in \{0,1,2\} \times \{0,1,2\}$$

$P_0$

(0, 0) (0, 1) (0, 2) (1, 0) (1, 1) (1, 2) (2, 0) (2, 1) (2, 2)
○ ○ ○ ○ ○ ○ ○ ○ ○
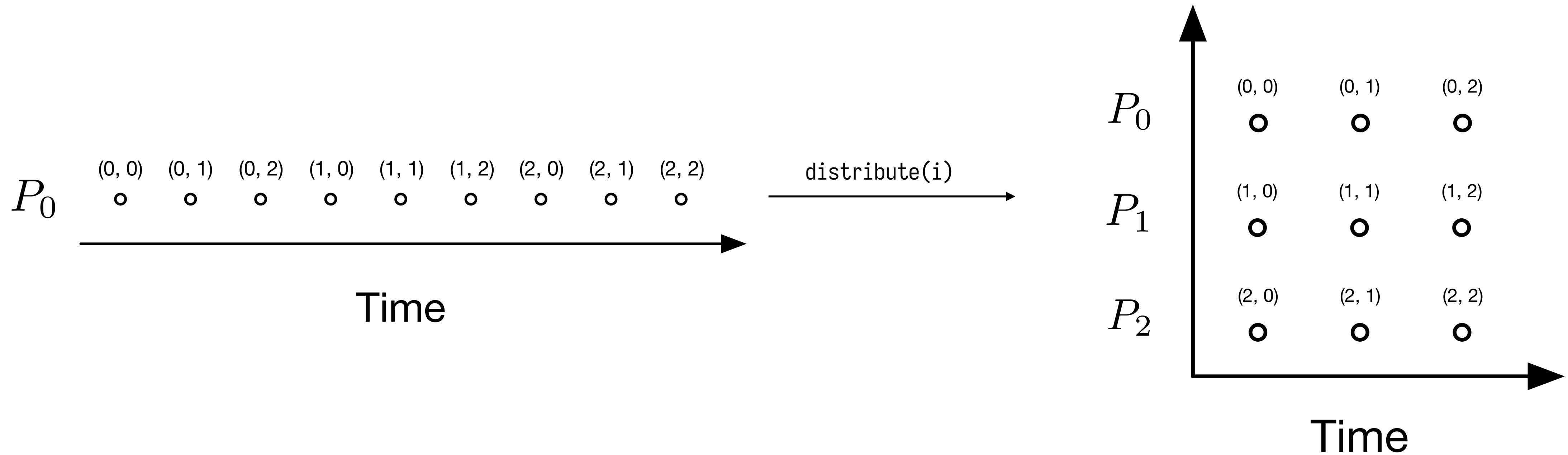
Time

# Scheduling

Change execution of iteration space through scheduling transformations
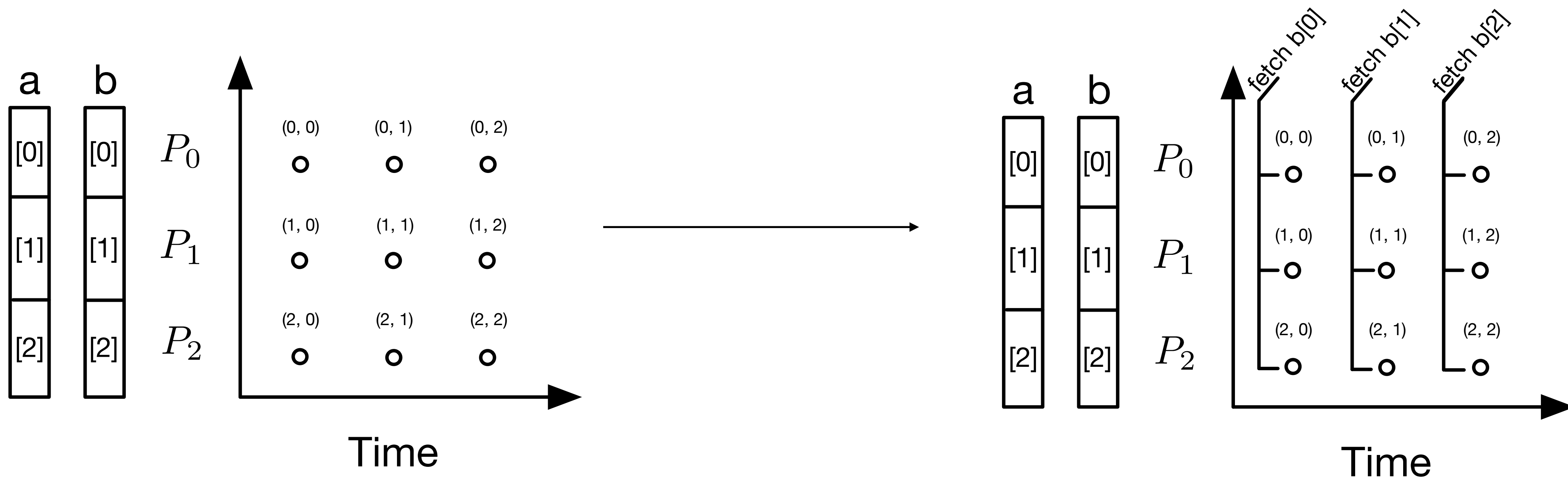
# distribute

$$\forall_i \ a(i) = \sum_j b(j) \qquad (i, j) \in \{0, 1, 2\} \times \{0, 1, 2\}$$

# What about communication?

$$\forall_i \ a(i) = \sum_j b(j) \qquad \text{s.t.} \qquad a \ _{\text{x}}\!\mapsto_{\text{x}} \mathcal{M} \quad b \ _{\text{x}}\!\mapsto_{\text{x}} \mathcal{M}$$
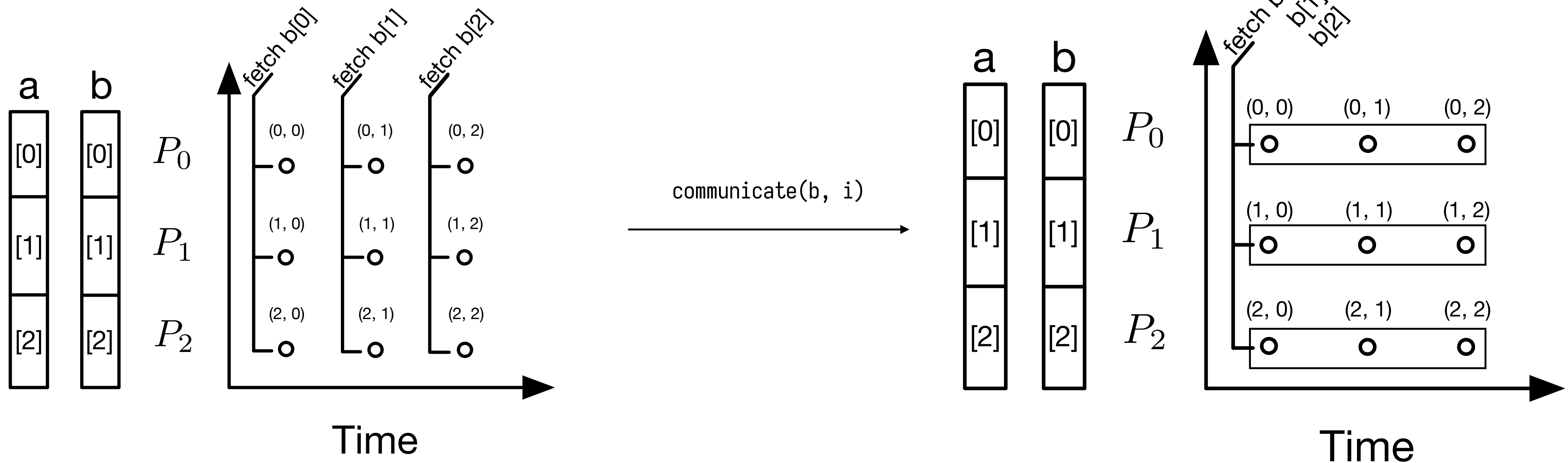
Insert communication at each iteration space point, as needed

# communicate

$$\forall_i \ a(i) = \sum_j b(j)$$

s.t. $\quad a_{\text{x}} \mapsto_{\text{x}} \mathcal{M} \quad b_{\text{x}} \mapsto_{\text{x}} \mathcal{M}$

Tradeoff: memory vs communication frequency!



communicate(b, i)

# Breaking symmetry with rotate

$$\forall_i \; a(i) = \sum_j b(j)$$

s.t. $\quad a \;_{\text{x}}\!\mapsto_{\text{x}} \mathcal{M} \quad b \;_{\text{x}}\!\mapsto_{\text{x}} \mathcal{M}$

# rotate

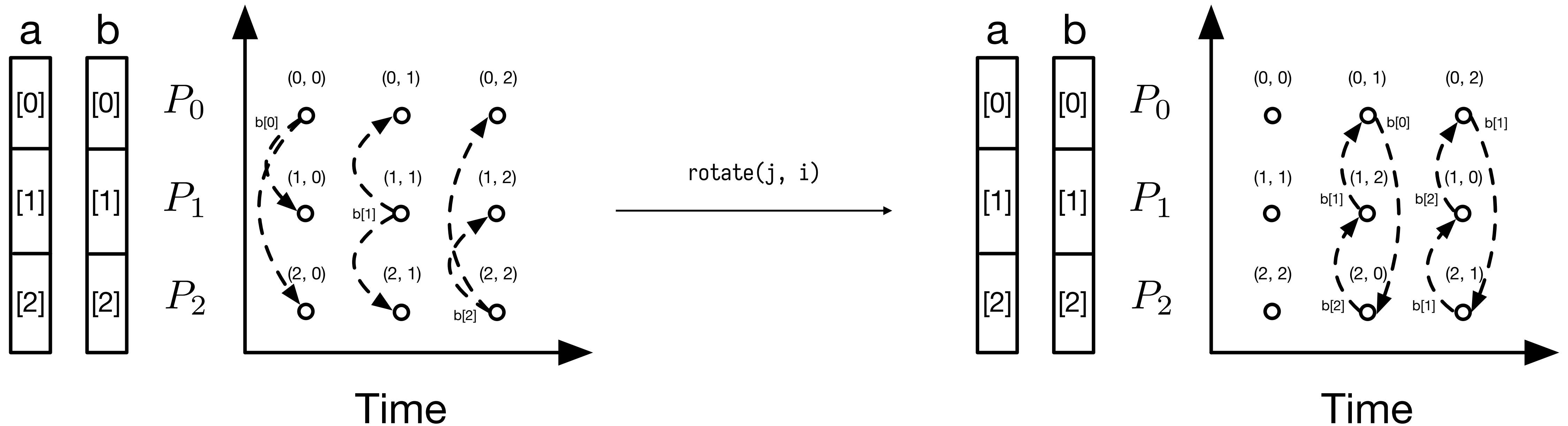$$\forall_i \; a(i) = \sum_j b(j) \qquad \text{s.t.} \qquad a_{\;x} \mapsto_x \mathcal{M} \quad b_{\;x} \mapsto_x \mathcal{M}$$

# rotate



rotate(j, i)

a    b

[0]  [0]  $P_0$

[1]  [1]  $P_1$

[2]  [2]  $P_2$

(0, 0)    (0, 1)    (0, 2)

b[0]      b[1]

(1, 1)  (1, 2)    (1, 0)

b[1]      b[2]

(2, 2)  (2, 0)    (2, 1)

b[2]      b[1]

Time

Use the modulus operator!

```
for j in (0, extent(j)):          for j' in (0, extent(j)):
                                      j = j' + i mod extent(j)
  …                                 …
```

# How expressive are these abstractions?

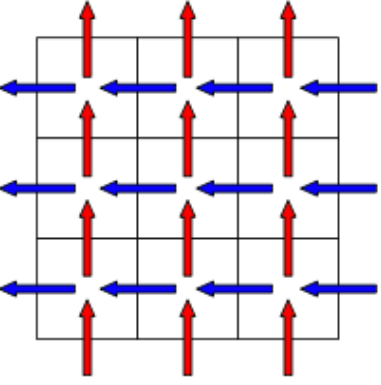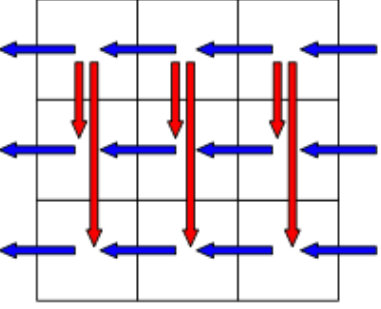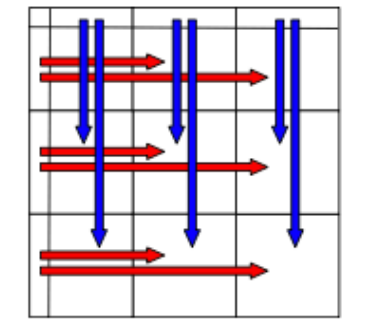| | Comm. Pattern | Target Machine | Data Distribution | Schedule |
|---|---|---|---|---|
| Cannon's Algorithm |  | $\mathcal{M}(gx, gy)$ | $A_{ij} \mapsto_{ij} \mathcal{M}$<br>$B_{ij} \mapsto_{ij} \mathcal{M}$<br>$C_{ij} \mapsto_{ij} \mathcal{M}$ | ```.distribute({i, j}, {in, jn}, {il, jl}, Grid(gx, gy))``` <br> ```.divide(k, ko, ki, gx)``` <br> ```.reorder({ko, il, jl, ki})``` <br> ```.rotate(ko, {in, jn}, kos)``` <br> ```.communicate(A, jn)``` <br> ```.communicate({B, C}, kos)``` |
| PUMMA |  | $\mathcal{M}(gx, gy)$ | $A_{ij} \mapsto_{ij} \mathcal{M}$<br>$B_{ij} \mapsto_{ij} \mathcal{M}$<br>$C_{ij} \mapsto_{ij} \mathcal{M}$ | ```.distribute({i, j}, {in, jn}, {il, jl}, Grid(gx, gy))``` <br> ```.divide(k, ko, ki, gx)``` <br> ```.reorder({ko, il, jl, ki})``` <br> ```.rotate(ko, {in}, kos)``` <br> ```.communicate(A, jn)``` <br> ```.communicate({B, C}, kos)``` |
| SUMMA |  | $\mathcal{M}(gx, gy)$ | $A_{ij} \mapsto_{ij} \mathcal{M}$<br>$B_{ij} \mapsto_{ij} \mathcal{M}$<br>$C_{ij} \mapsto_{ij} \mathcal{M}$ | ```.distribute({i, j}, {in, jn}, {il, jl}, Grid(gx, gy))``` <br> ```.split(k, ko, ki, chunkSize)``` <br> ```.reorder({ko, il, jl, ki})``` <br> ```.communicate(A, jn)``` <br> ```.communicate({B, C}, ko)``` |
| Johnson's Algorithm |  | $\mathcal{M}(\sqrt[3]{p}, \sqrt[3]{p}, \sqrt[3]{p})$ | $A_{ij} \mapsto_{ij0} \mathcal{M}$<br>$B_{ik} \mapsto_{i0k} \mathcal{M}$<br>$C_{kj} \mapsto_{0jk} \mathcal{M}$ | ```.distribute({i, j, k}, {in, jn, kn},```<br>```         {il, jl, kl}, Grid(∛p, ∛p, ∛p))``` <br> ```.communicate({A, B, C}, kn)``` |
| Solomonik's Algorithm |  | $\mathcal{M}(\sqrt{\frac{p}{c}}, \sqrt{\frac{p}{c}}, c)$ | $A_{ij} \mapsto_{ij0} \mathcal{M}$<br>$B_{ij} \mapsto_{ij0} \mathcal{M}$<br>$C_{ij} \mapsto_{ij0} \mathcal{M}$ | ```.distribute({i, j, k}, {in, jn, kn},```<br>```         {il, jl, kl}, Grid(√(p/c), √(p/c), c))``` <br> ```.divide(kl, k1, k2, √(p/c³))``` <br> ```.reorder({k1, il, jl, k2})``` <br> ```.rotate(k1, {in, jn}, k1s)``` <br> ```.communicate(A, jn)``` <br> ```.communicate({B, C}, k1s)``` |
| COSMA |  | induced by schedule | induced by schedule | ```// gx, gy, gz, numSteps computed by COSMA scheduler.``` <br> ```.distribute({i, j, k}, {in, jn, kn}```<br>```         {il, jl, kl}, Grid(gx, gy, gz))``` <br> ```.divide(kl, klo, kli, numSteps)``` <br> ```.reorder({klo, il, jl, kli})``` <br> ```.communicate(A, kn)``` <br> ```.communicate({B, C}, klo)``` |

# SUMMA Algorithm

A(i, j) = B(i, k) * C(k, j)

```
# Arrange p processors into a 2D grid.
# Assign a tile of A, B, C to each processor.
for all P_{ij} in parallel:
  for kc in (0, k, chunkSize):
    B_l = row broadcast the kc to kc+chunkSize columns of B
    C_l = col broadcast the kc to kc+chunkSize rows of C
    A += B_l * C_l
```

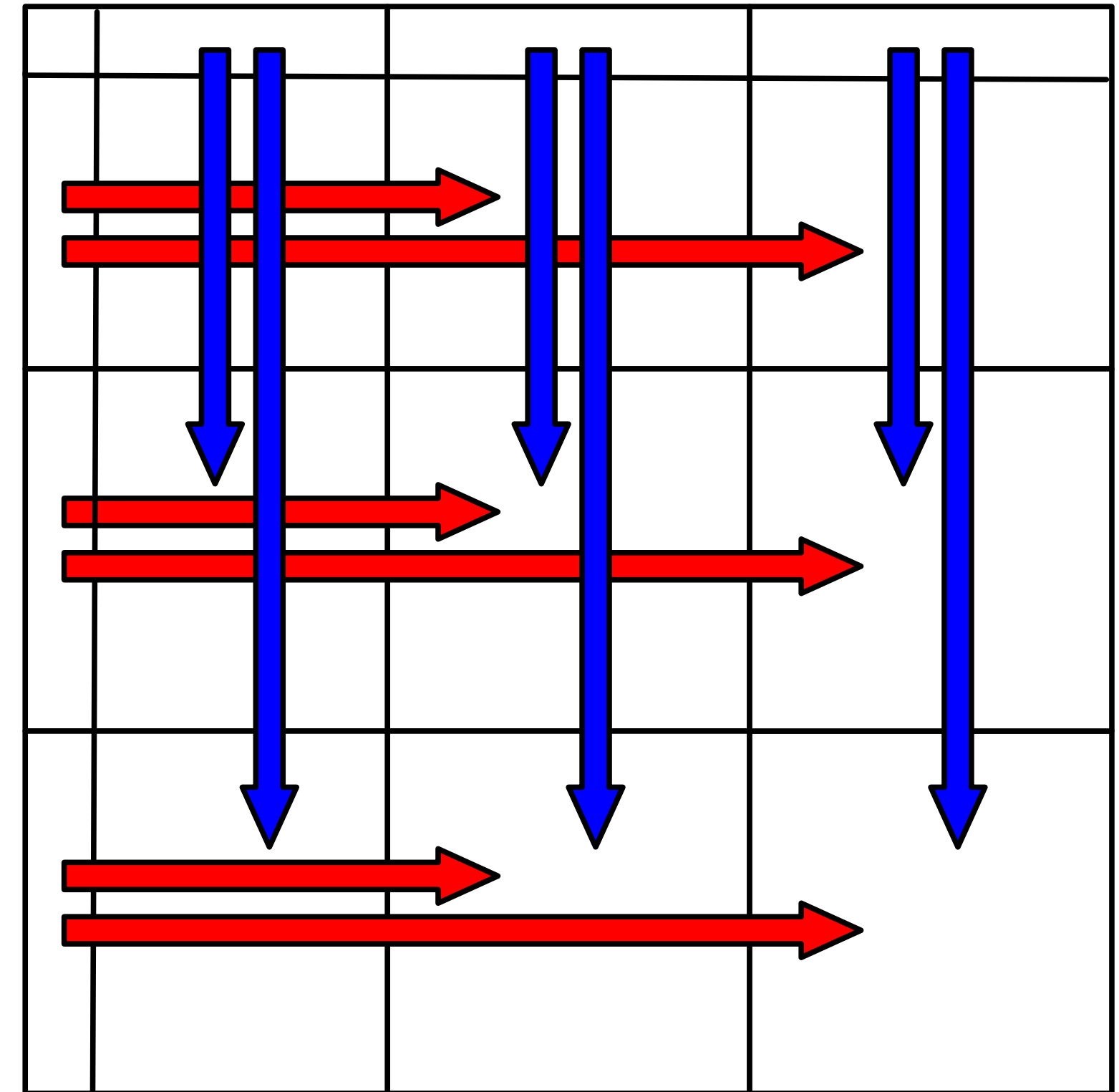$$\mathcal{M} = \mathsf{Grid}(gx, gy)$$

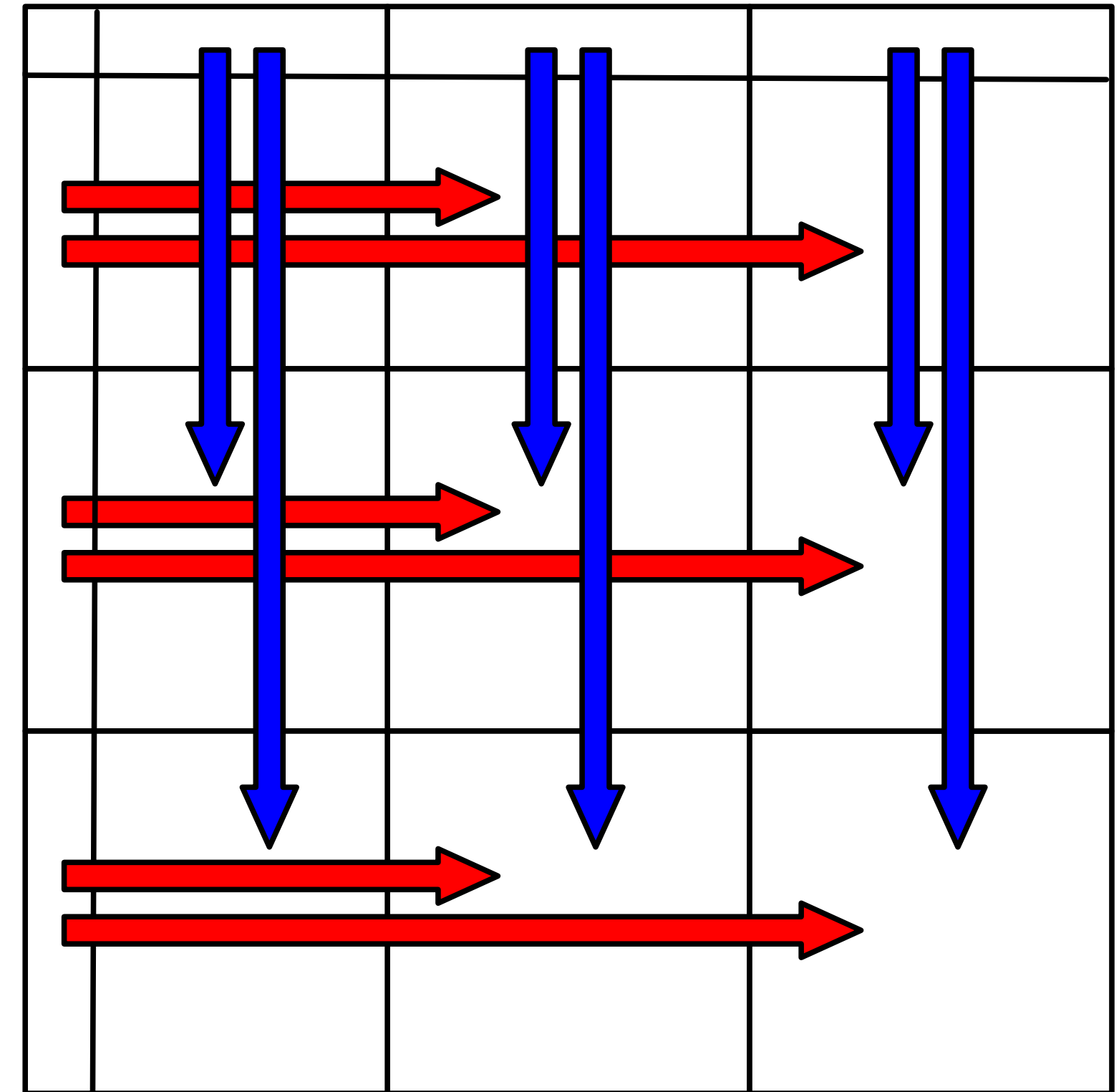$$A \mathrel{{}_{\mathrm{xy}}\!\!\mapsto_{\mathrm{xy}}} \mathcal{M}$$

$$B \mathrel{{}_{\mathrm{xy}}\!\!\mapsto_{\mathrm{xy}}} \mathcal{M}$$

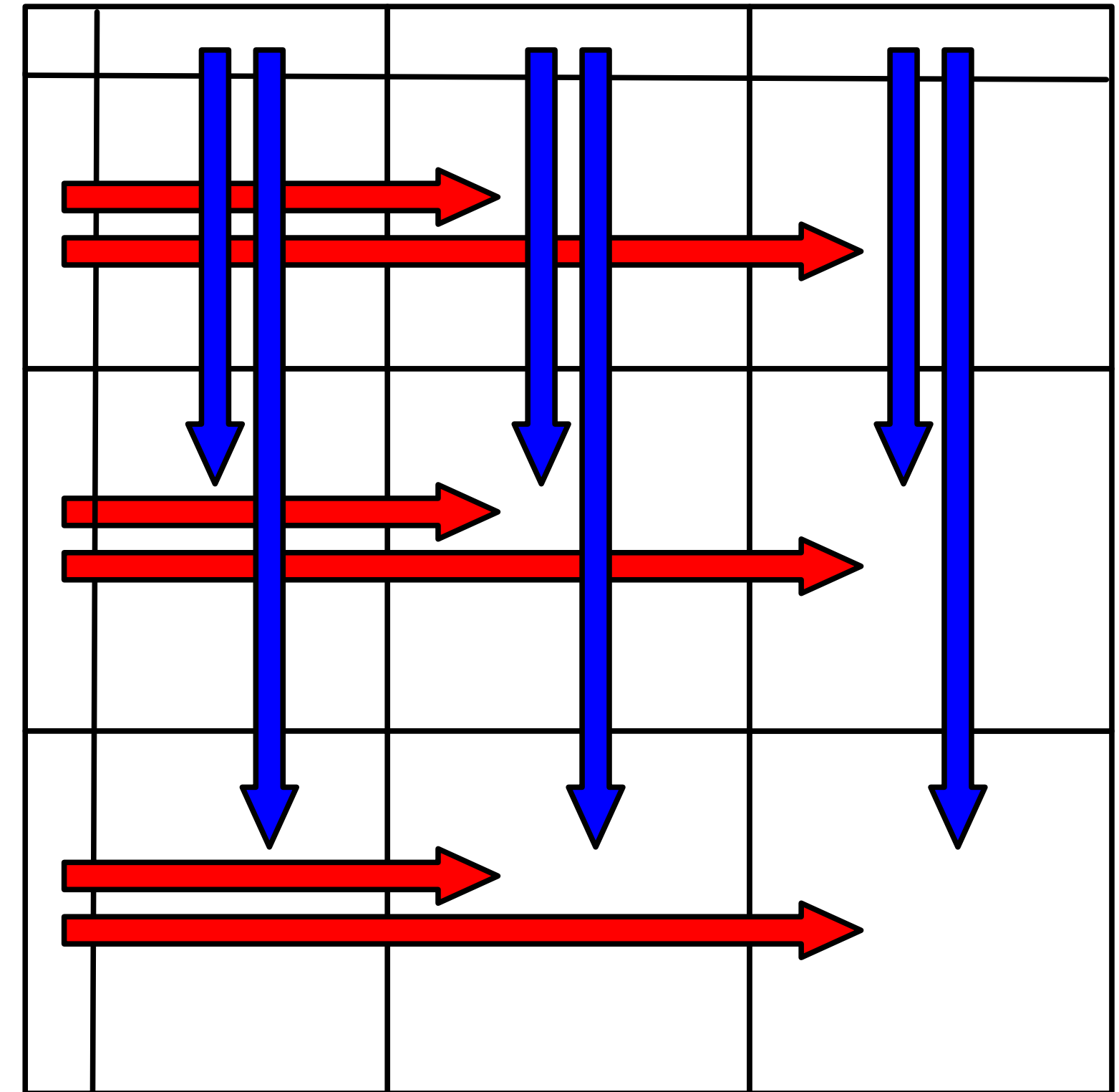$$C \mathrel{{}_{\mathrm{xy}}\!\!\mapsto_{\mathrm{xy}}} \mathcal{M}$$

```
for i:
  for j:
    for k:
      A(i, j) += B(i, k) * C(k, j)
```

divide(i, il, in, gx)
divide(j, jl, jn, gy)
reorder({in, jn, il, jl})

```
for in:
  for jn:
    for il:
      for jl:
        for k:
          A(i, j) += B(i, k) * C(k, j)
```
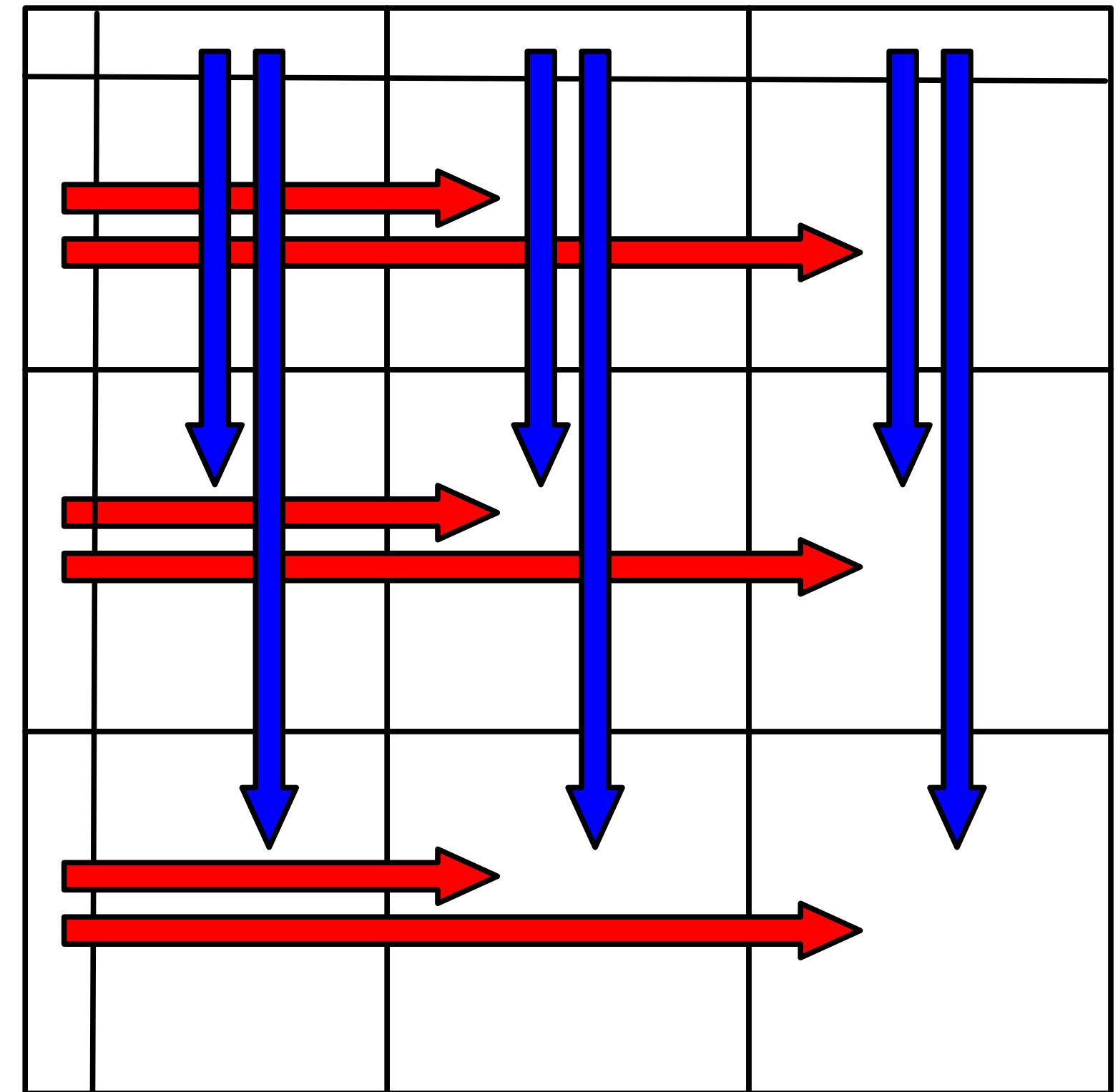
```
divide(i, il, in, gx)
divide(j, jl, jn, gx)
reorder({in, jn, il, jl})
split(k, ko, ki, chunkSize)
reorder(ko, il, jl, ki)
```

```
for in:
  for jn:
    for ko:
      for il:
        for jl:
          for ki:
            A(i, j) += B(i, k) * C(k, j)
```
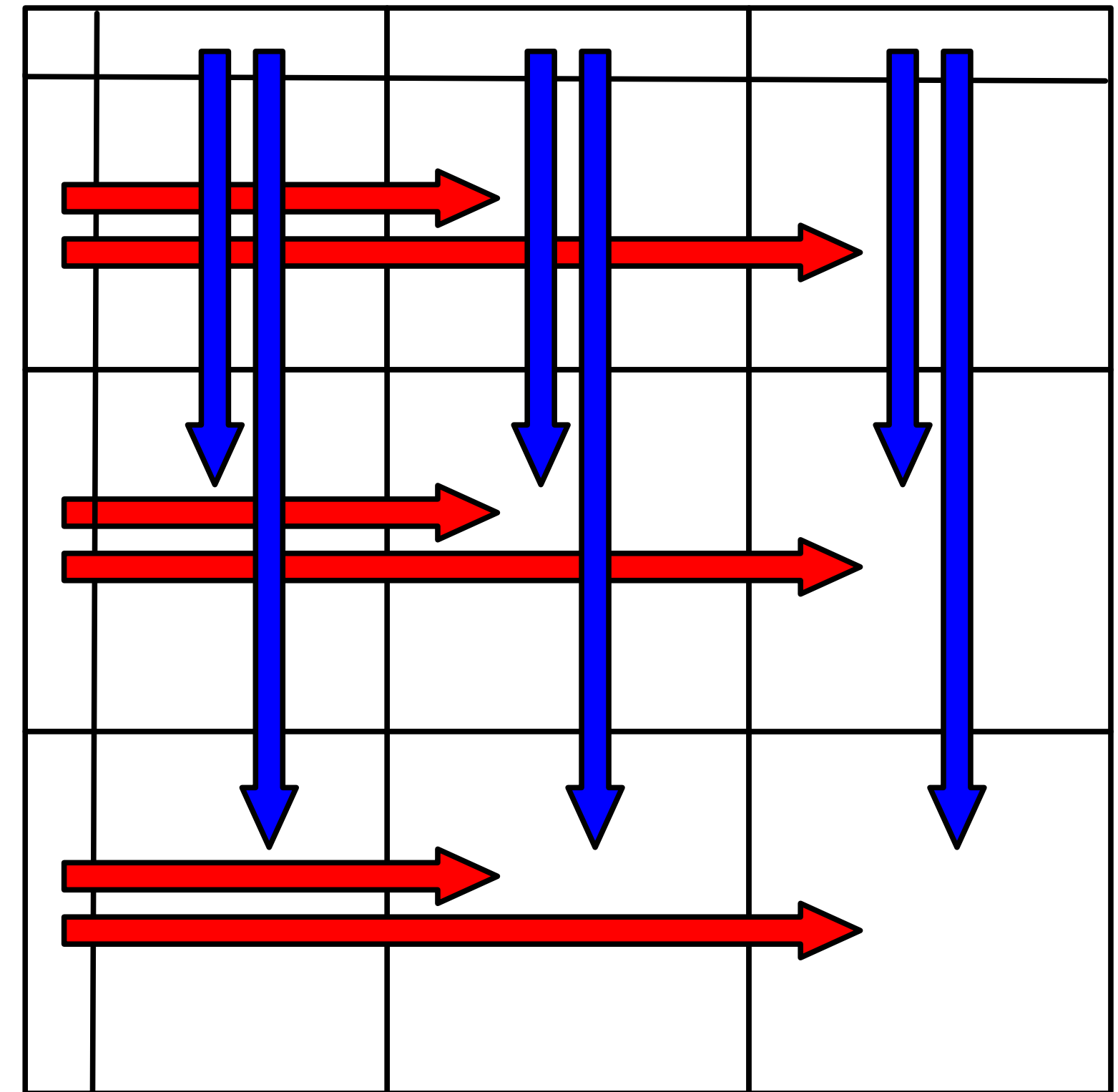
```
divide(i, il, in, gx)
divide(j, jl, jn, gx)
reorder({in, jn, il, jl})
split(k, ko, ki, chunkSize)
reorder(ko, il, jl, ki)
distribute(in, jn)
```

```
distributed for in, jn:
  for ko:
    for il:
      for jl:
        for ki:
          A(i, j) += B(i, k) * C(k, j)
```
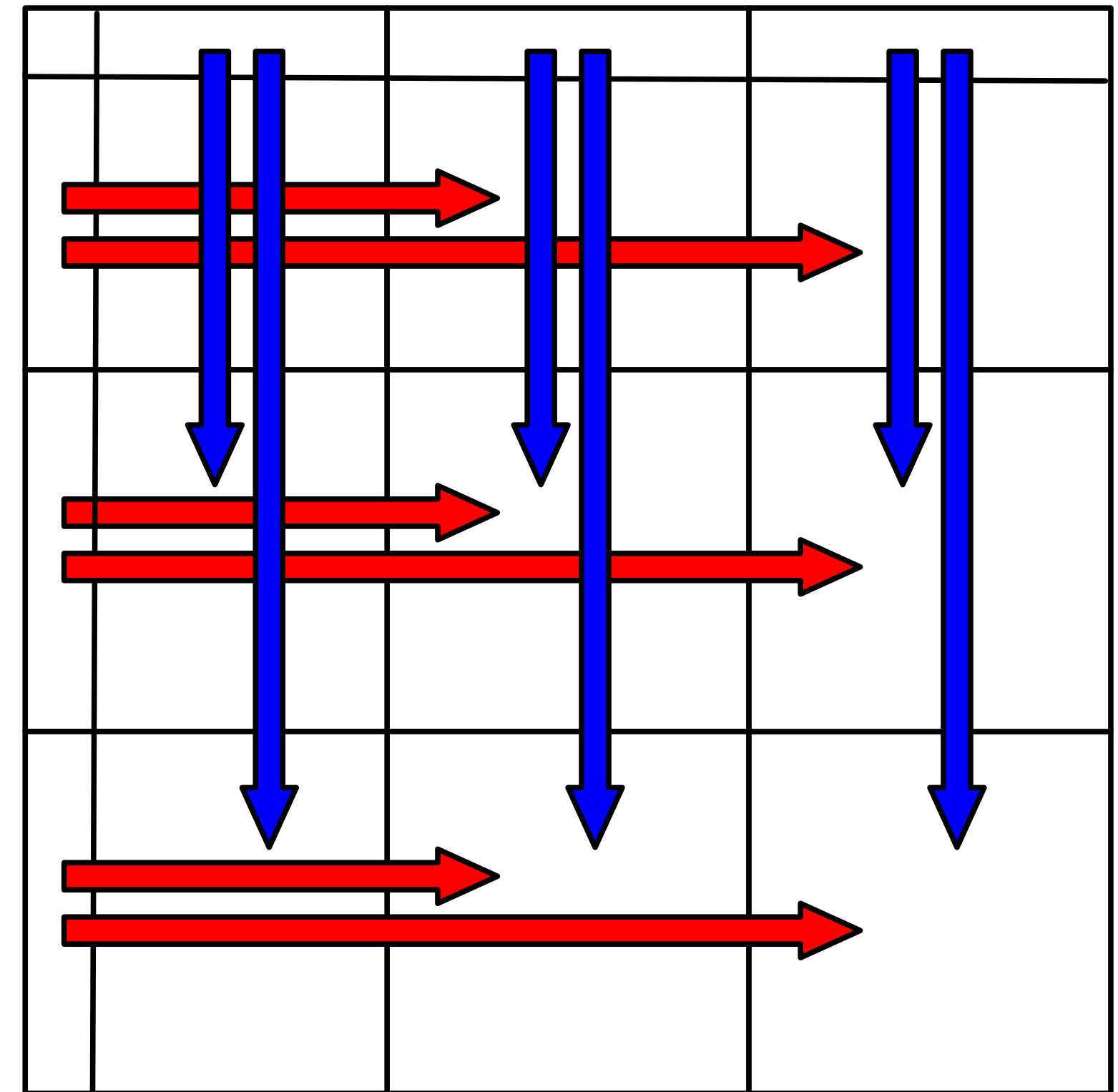
```
distributed for in, jn:
  communicate A
  for ko:
    communicate B, C
    for il:
      for jl:
        for ki:
          A(i, j) += B(i, k) * C(k, j)
```

```
divide(i, il, in, gx)
divide(j, jl, jn, gx)
reorder({in, jn, il, jl})
split(k, ko, ki, chunkSize)
reorder(ko, il, jl, ki)
distribute(in, jn)
communicate(A, jn)
communicate({B, C}, ko)
```

# Compilation Process

## Expression

$$A(i, j) = B(i, k) \cdot C(k, j)$$
$$A(i, l) = B(i, j, k) \cdot C(j, l) \cdot D(k, l)$$
$$a = B(i, j, k) \cdot C(i, j, k)$$
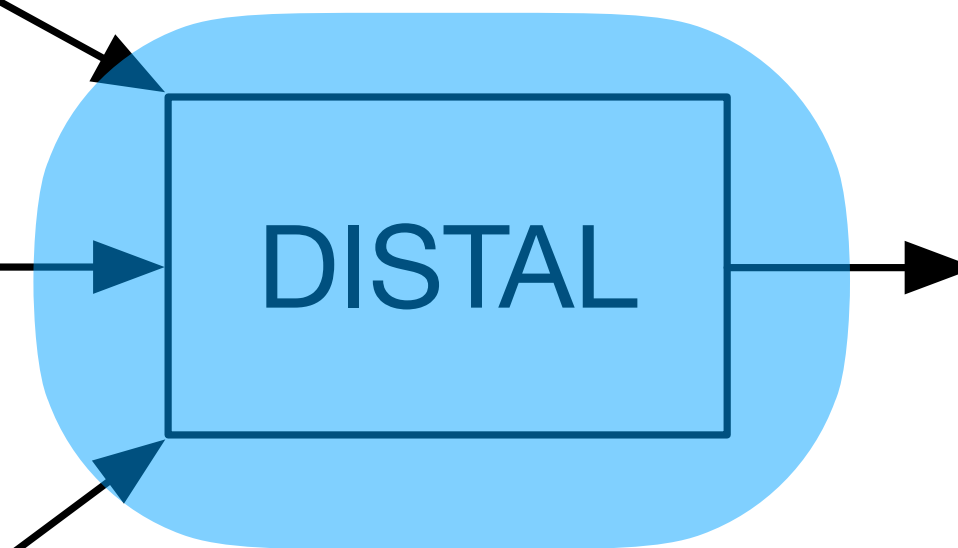$$A(i, j, l) = B(i, j, k) \cdot C(k, l)$$
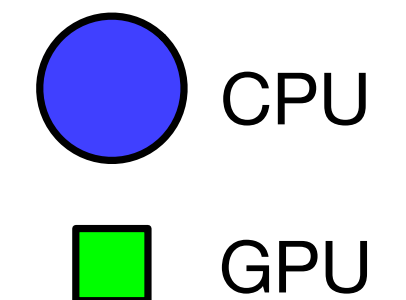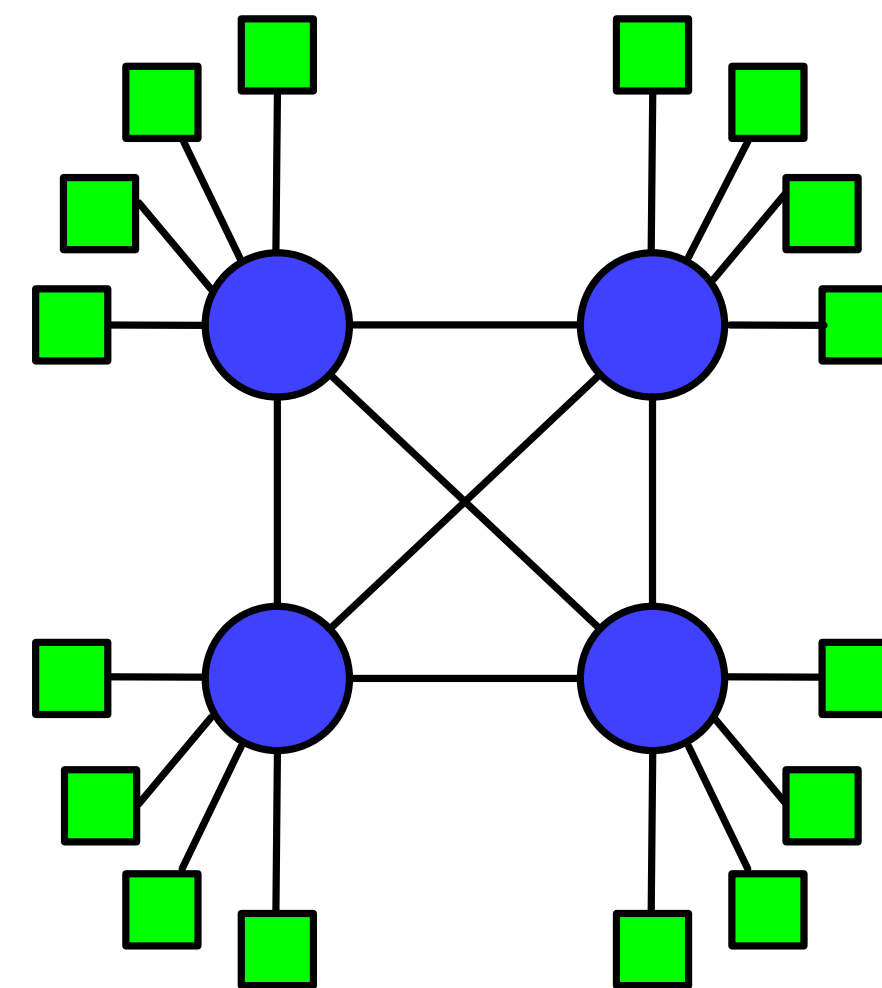$$A(i, j) = B(i, j, k) \cdot c(k)$$

## Data Distribution
Partition A into tiles
Replicate B onto all nodes
Place C onto only some nodes

## Computation Distribution
Owner Computes
Distribute i,j loops
Communicate in chunks

DISTAL

Supercomputer

CPU

GPU

40

# Translate to Concrete Index Notation

$$A(i,j) = B(i,j,k) \cdot c(k) \quad \longrightarrow \quad \forall_i \forall_j \forall_k \ A(i,j) \mathrel{+}= B(i,j,k) \cdot c(k)$$

Iteration order unspecified                Specified Iteration Order

# Scheduling operations rewrite Concrete Index Notation

$$\ldots \forall_i S \xrightarrow{\text{divide}(i,i_o,i_i,c)} \ldots \forall_{i_o} \forall_{i_i} S \text{ s.t. divide}(i, i_o, i_i, c)$$

$$\ldots \forall_i S \xrightarrow{\text{distribute}(i)} \ldots \forall_i S \text{ s.t. distribute}(i)$$

$$\ldots \forall_I \forall_t S \xrightarrow{\text{rotate}(t,I,r)} \ldots \forall_I \forall_r S \text{ s.t. rotate}(t, I, r)$$

$$\ldots \forall_i S \xrightarrow{\text{communicate}(\mathcal{T},i)} \ldots \forall_i S \text{ s.t. communicate}(\mathcal{T}, i)$$

Target specific backend handles these constructs now!

# Compiling Tensor Distribution Notation

$$\mathcal{T}_{\ xy \mapsto_x} \mathcal{M}$$

$$\downarrow$$

$$\forall_x \forall_y \ \mathcal{T}(x, y)$$

$$\downarrow$$

$$\forall_{xo} \forall_{xi} \forall_y \ \mathcal{T}(x, y) \text{ s.t. divide}(x, xo, xi, gx)$$

$$\downarrow$$

$$\forall_{xo} \forall_{xi} \forall_y \ \mathcal{T}(x, y) \text{ s.t. divide}(x, xo, xi, gx), \text{distribute}(xo), \text{comm.}(\mathcal{T}, xo)$$

# What's the backend?

# Legion

Distributed task-based runtime system

   Tasks operate on bulk data

   System moves memory between processors for tasks to use

`distribute` → launch a set of tasks

`communicate` → tell Legion the data to transfer

`rotate` → perform a loop transformation with modulus

# Evaluation

# Comparisons

Distributed GEMM — evaluate performance on a highly optimized kernel

Compare against COSMA, Cyclops Tensor Framework, ScaLAPACK

Higher order tensor kernels — evaluate performance on the long tail

Compare against Cyclops Tensor Framework

# Results (Methodology)

Experiments run on up to 256 nodes of Lassen (4 V100 GPUs/node, 40 Power9 CPUs/node, IB interconnect)

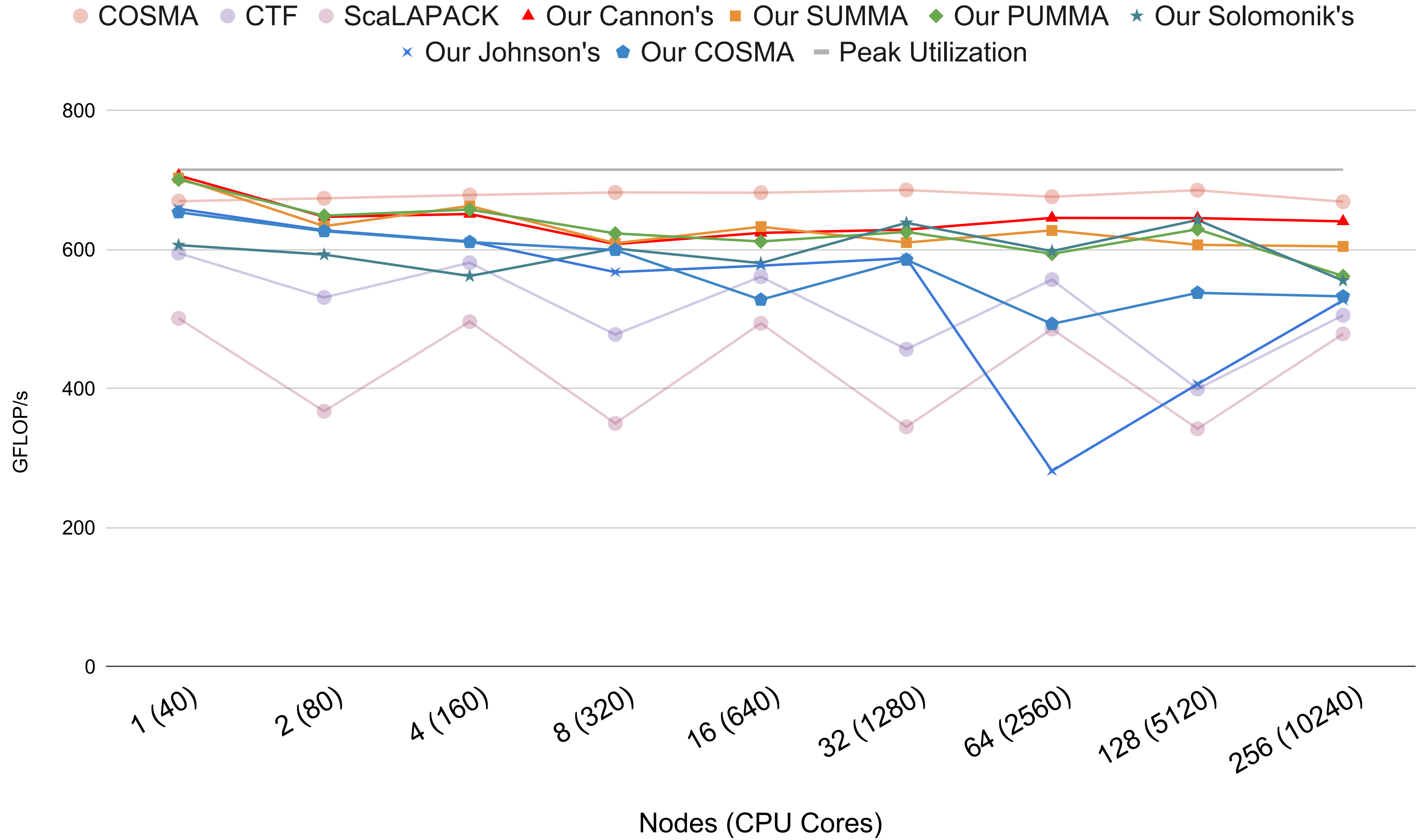All systems configured to use the same BLAS / CuBLAS for GEMM

All experiments are weak-scaling (memory / node stays constant)

Results reported in GFLOP/s (compute bound) and GB/s (bandwidth bound) per node
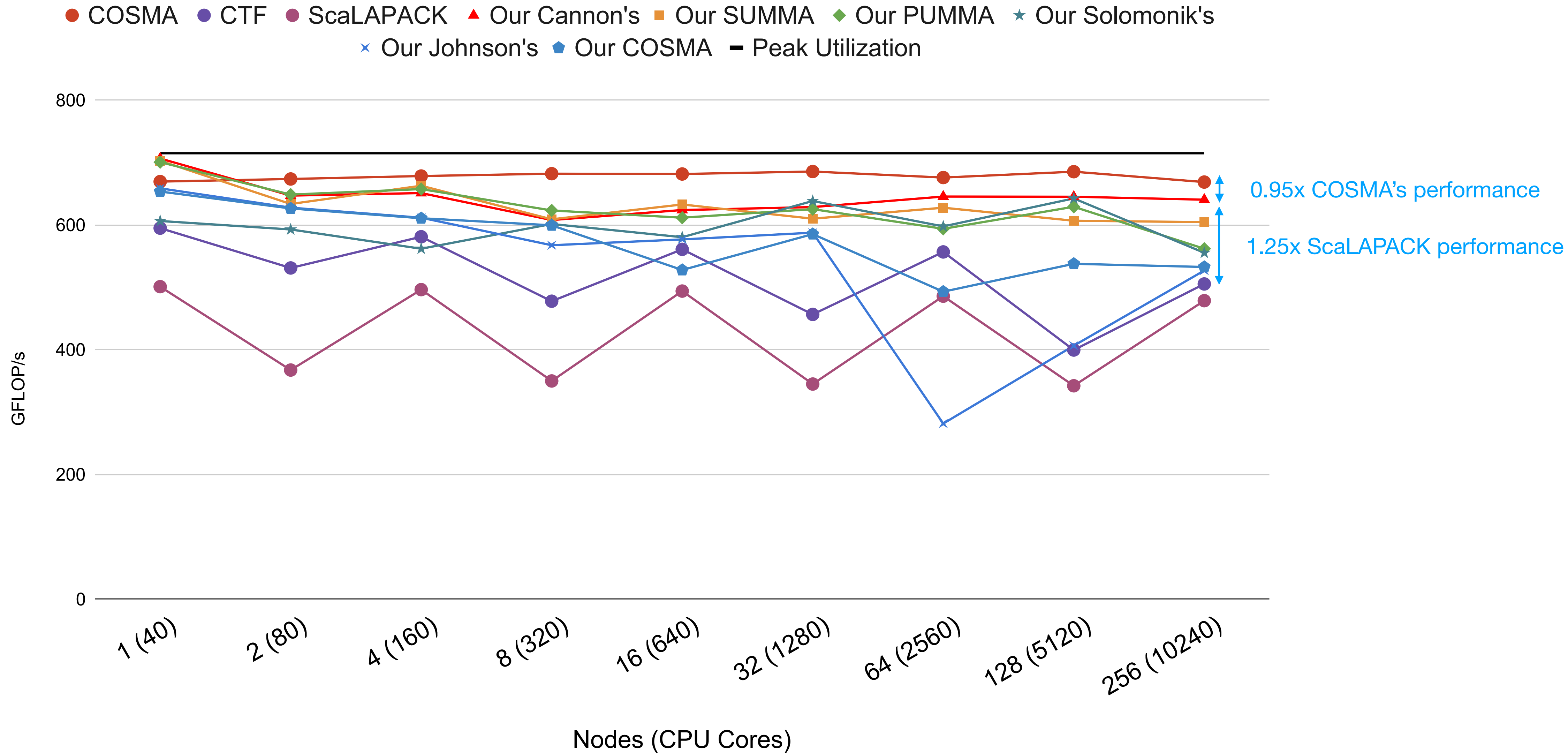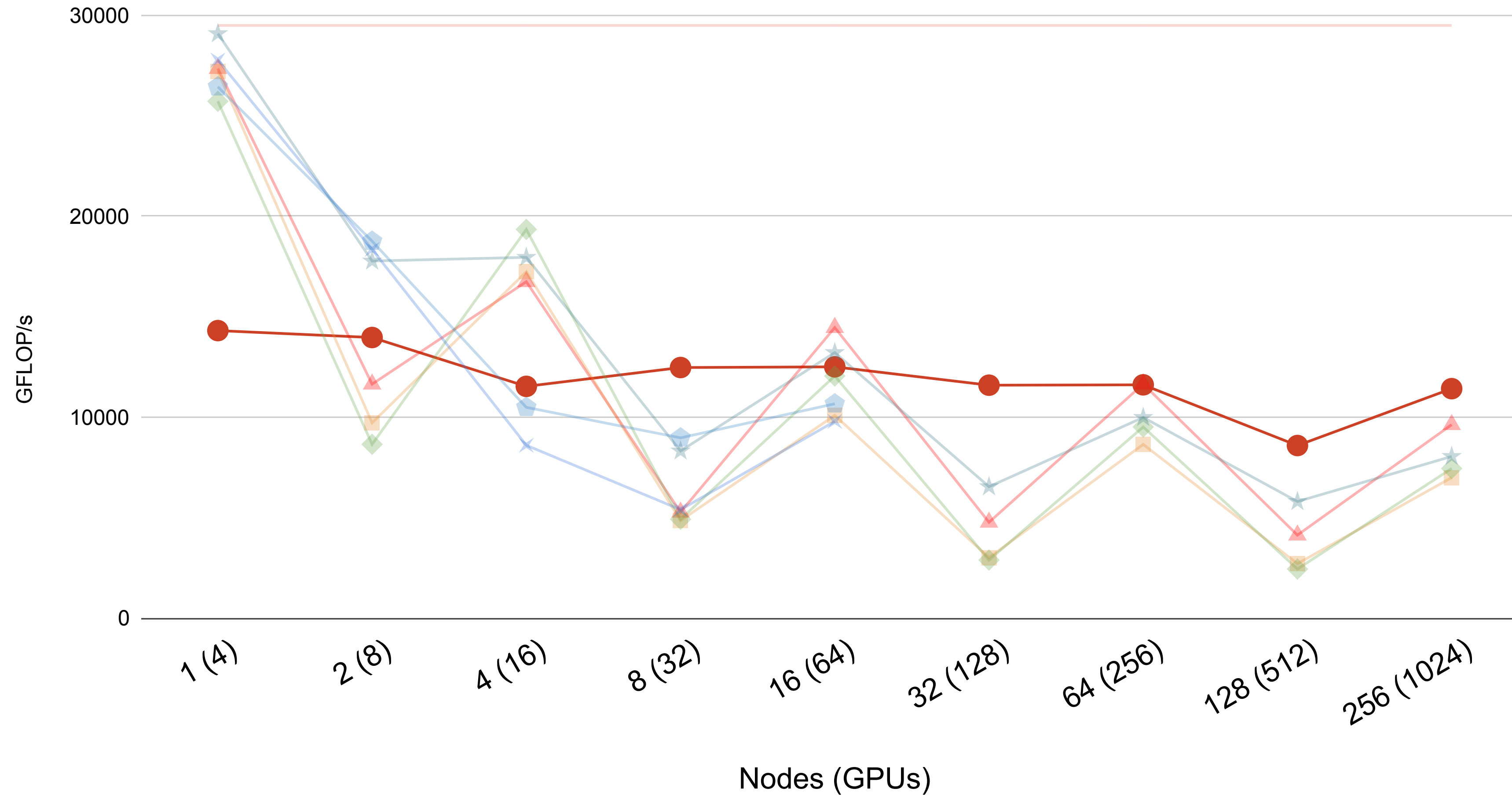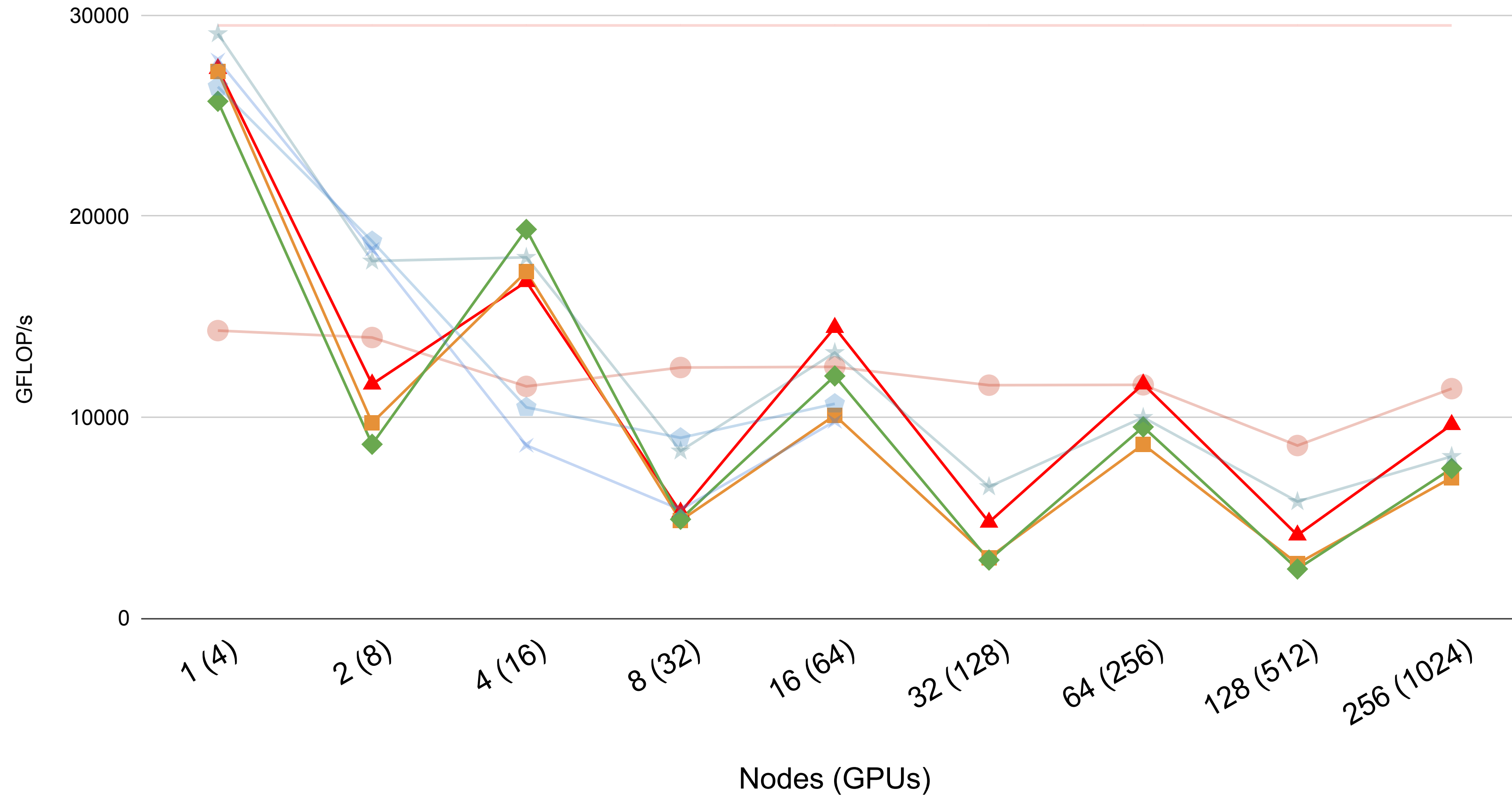
# GEMM (CPU)

# GEMM (CPU)

# GEMM (CPU)

# GEMM (GPU)



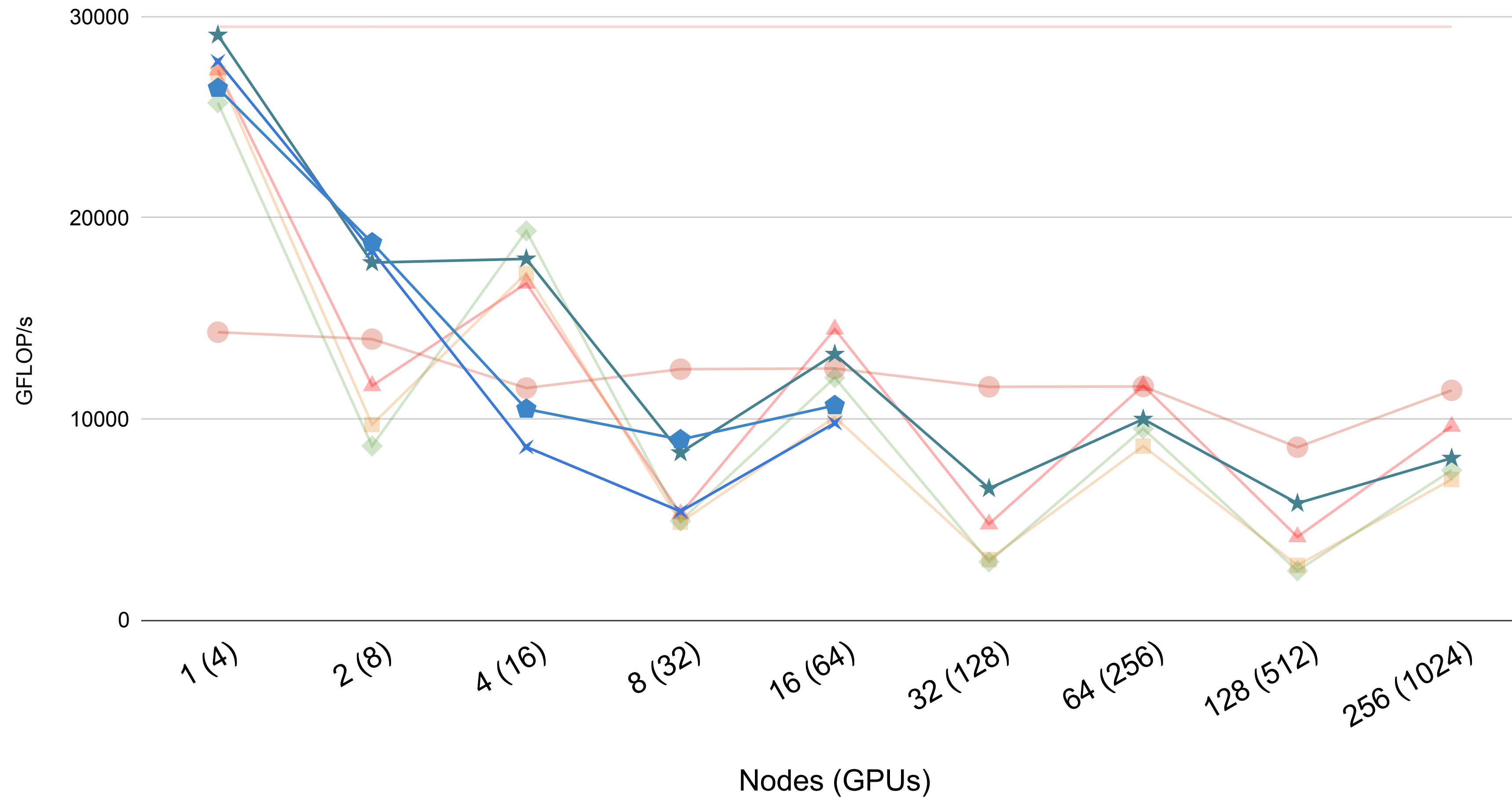COSMA ● Our Cannon's ▲ Our SUMMA ■ Our PUMMA ◆ Our Solomonik's ★ Our Johnson's ✕ Our Cosma ⬠ Peak Utilization ━

GFLOP/s

Nodes (GPUs)

# GEMM (GPU)

# GEMM (GPU)

# GEMM (GPU)



Legend: COSMA, Our Cannon's, Our SUMMA, Our PUMMA, Our Solomonik's, Our Johnson's, Our Cosma, Peak Utilization

Y-axis: GFLOP/s (0, 10000, 20000, 30000)

X-axis: Nodes (GPUs) — 1 (4), 2 (8), 4 (16), 8 (32), 16 (64), 32 (128), 64 (256), 128 (512), 256 (1024)
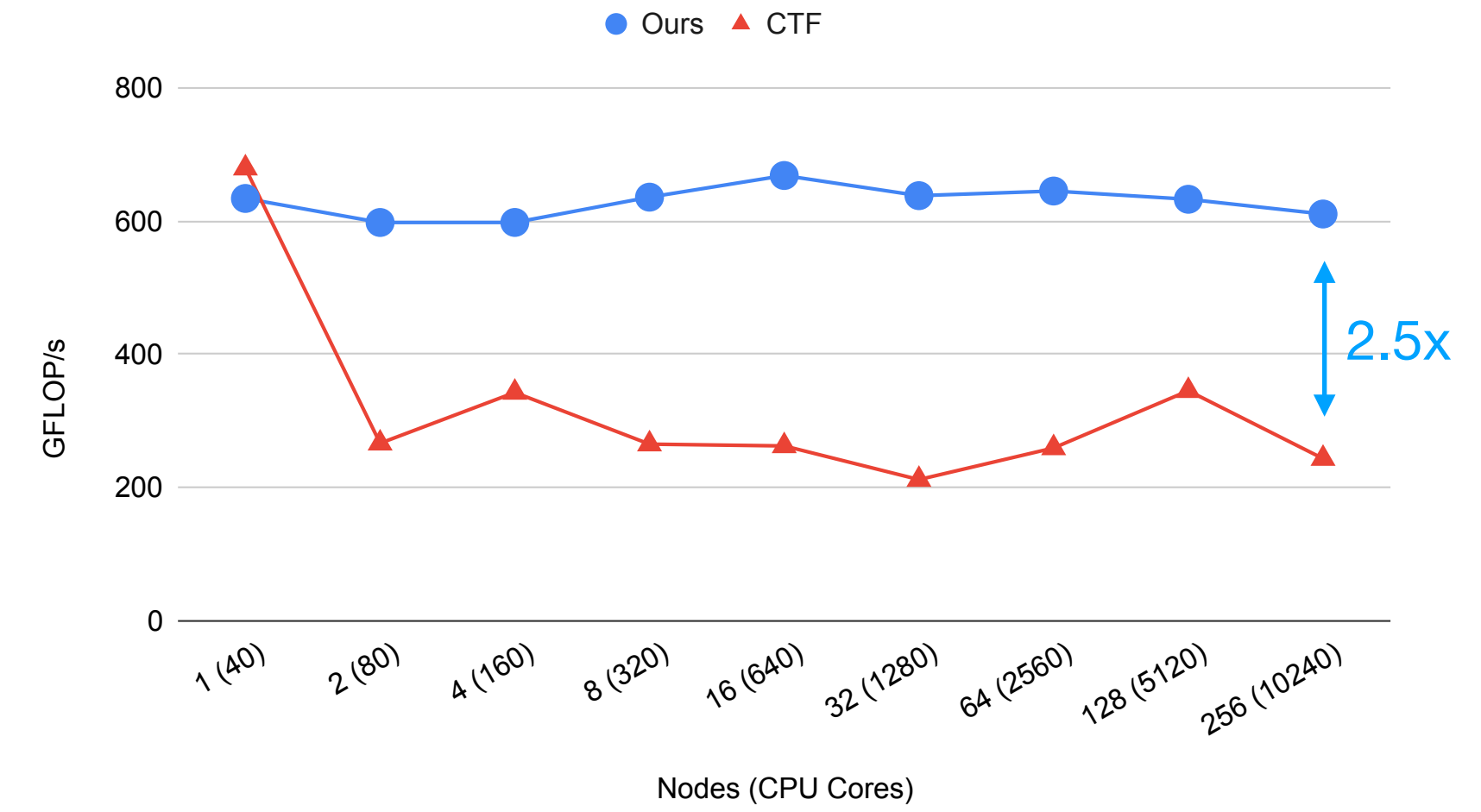
0.85x COSMA's performance

# Higher Order Tensor Operations (CPU)
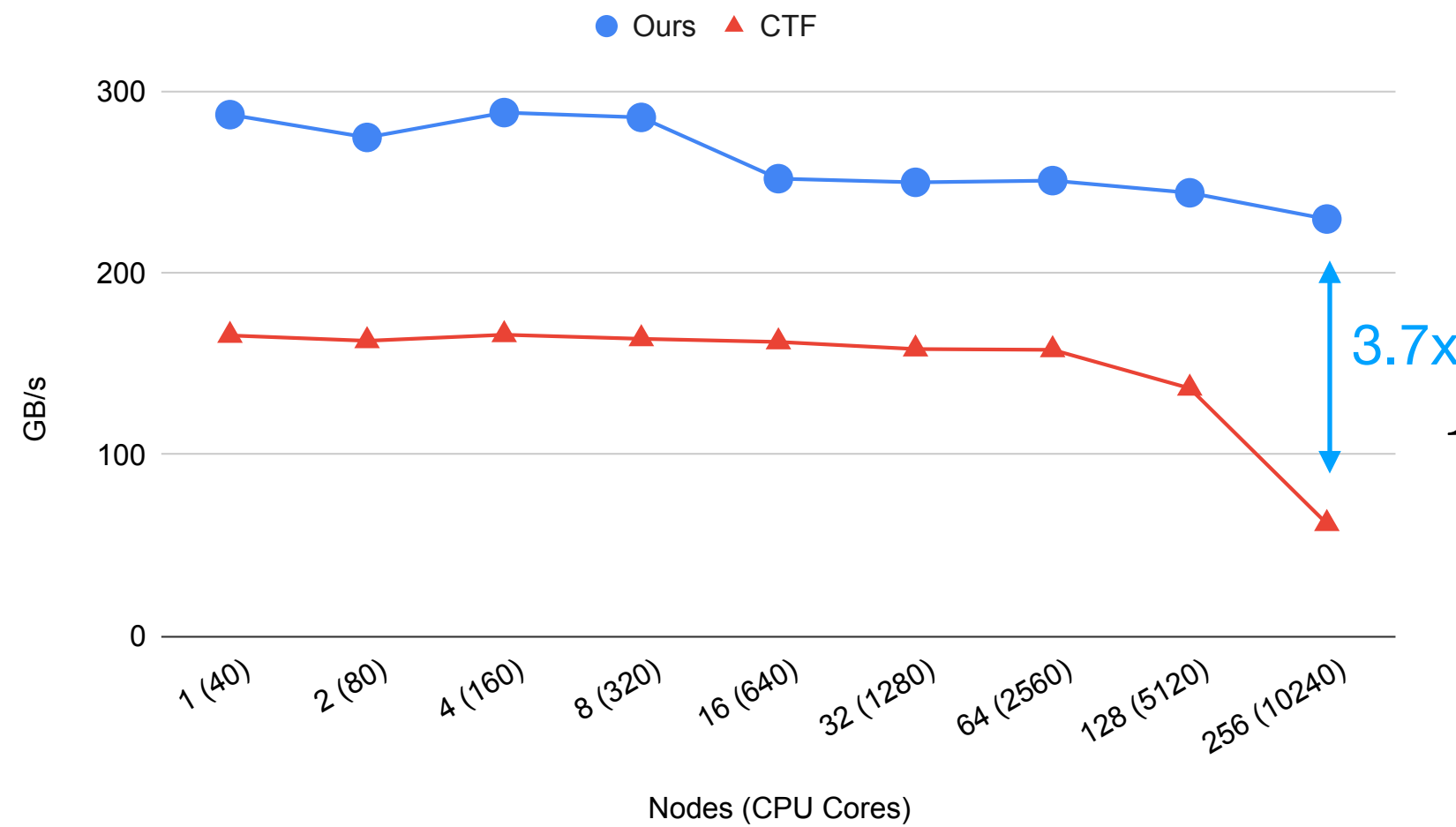
TTV

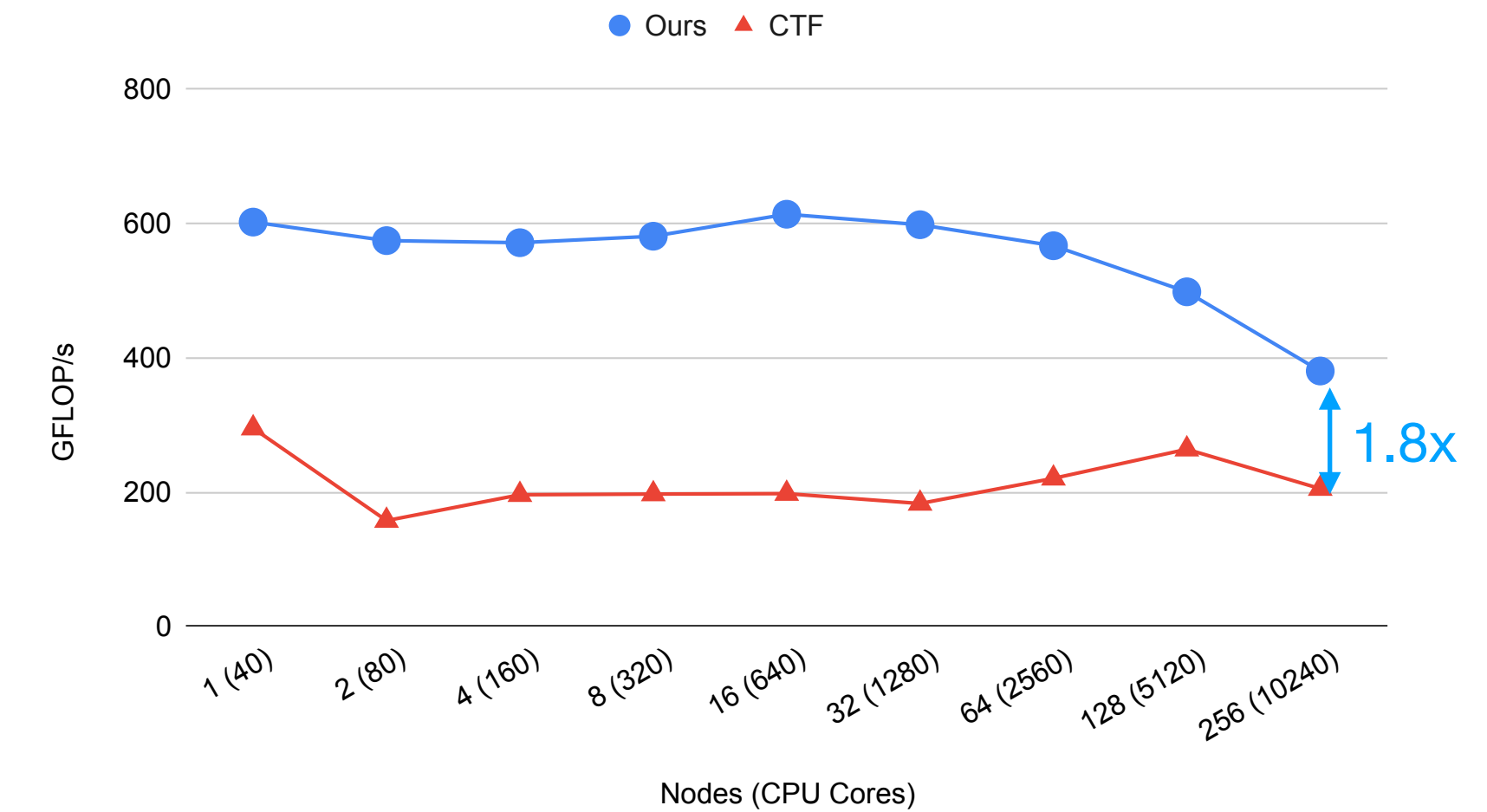$$A_{ij} = B_{ijk} \cdot C_k$$



TTM

$$A_{ijl} = B_{ijk} \cdot C_{kl}$$



InnerProd

$$a = B_{ijk} \cdot C_{ijk}$$



MTTKRP

$$A_{il} = B_{ijk} \cdot C_{jl} \cdot D_{kl}$$



56

# Higher Order Tensor Operations (GPU)

TTV

$$A_{ij} = B_{ijk} \cdot C_k$$



TTM

$$A_{ijl} = B_{ijk} \cdot C_{kl}$$
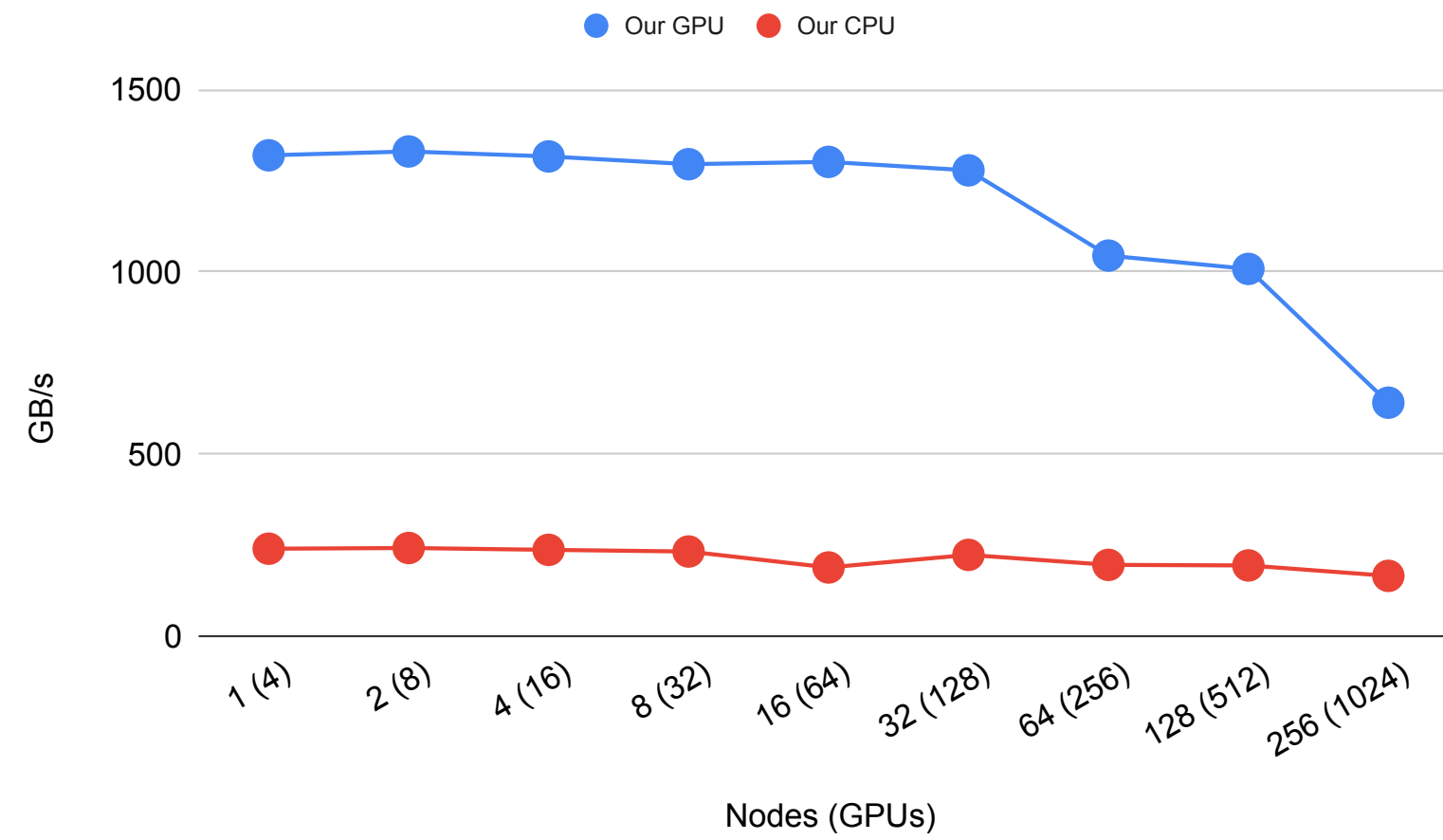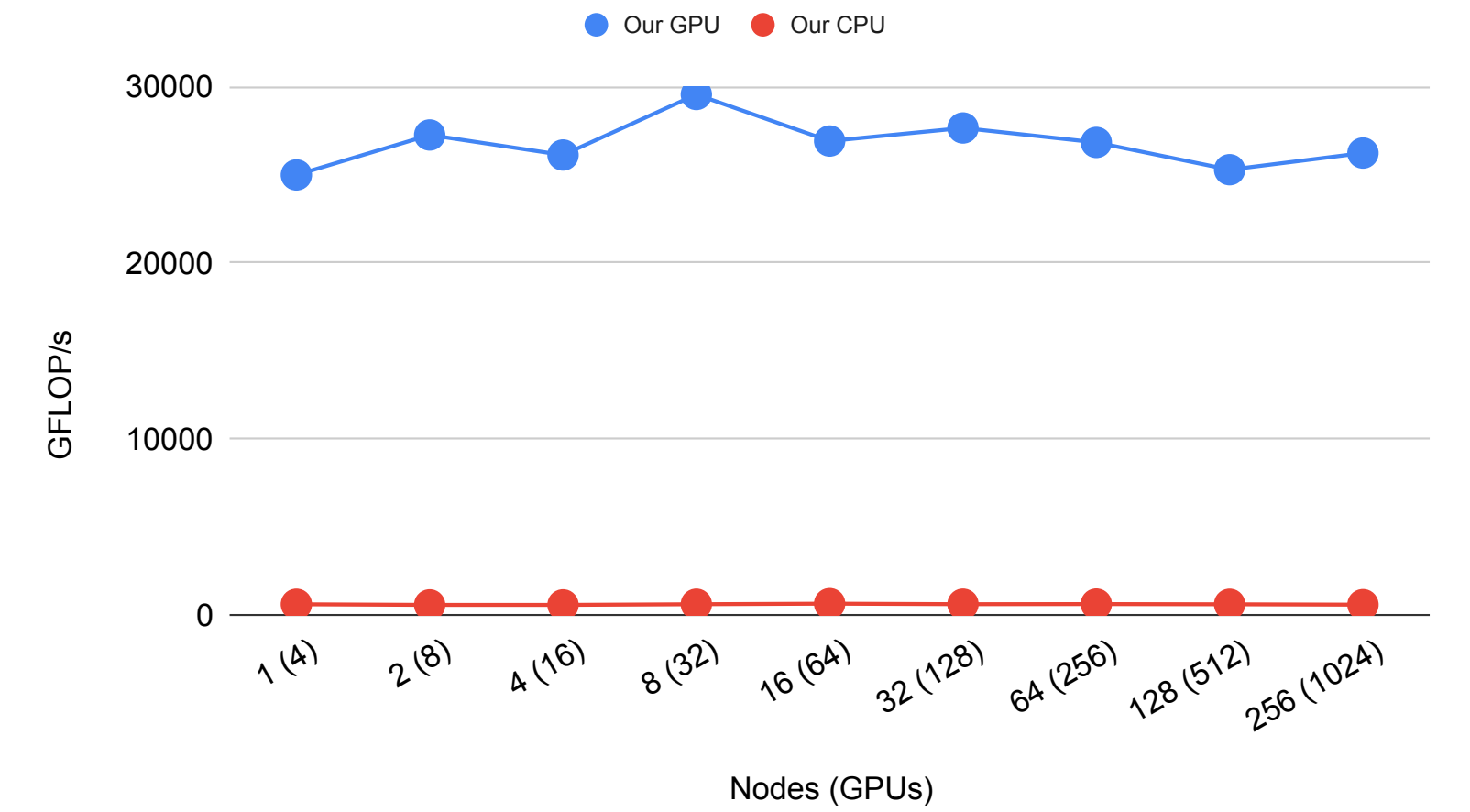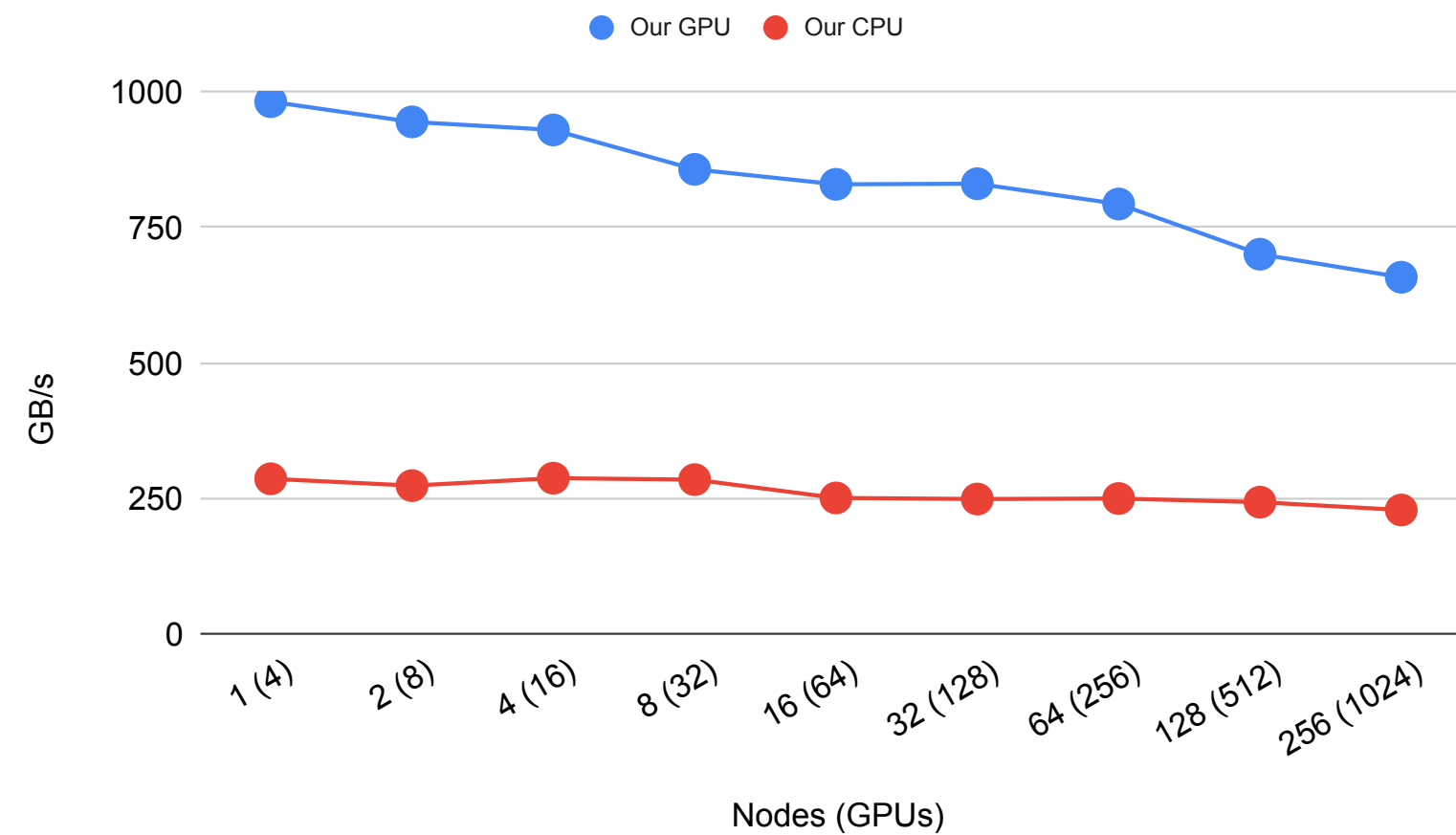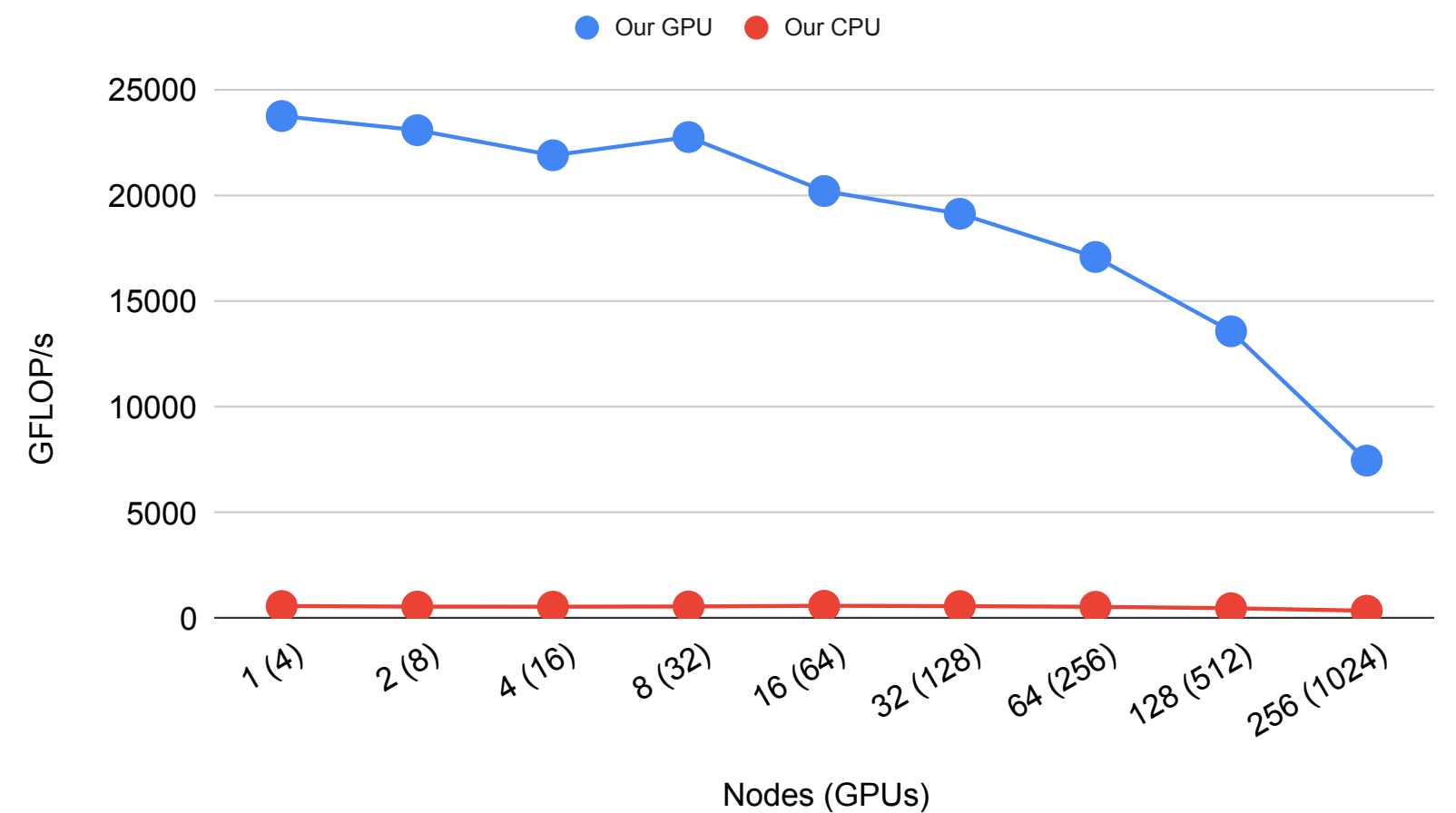


InnerProd

$$a = B_{ijk} \cdot C_{ijk}$$



MTTKRP

$$A_{il} = B_{ijk} \cdot C_{jl} \cdot D_{kl}$$

# Conclusion

DISTAL combines separate specifications of data and computation distribution

DISTAL can represent many existing algorithms

DISTAL can achieve high performance

Future work — extending to sparse tensors.
Vision: distributed implementations of ANY tensor program with ANY tensor formats!

Contact: rohany@cs.stanford.edu

# Extra slides

# DISTAL

- Decouple computation, performance optimizations, and data distribution

- Extension to TACO

Einsum notation programs

✖

Format-based data distribution

✖

Schedule for compute distribution

═

Too many to code by hand!

### Expression

$$A(i,j) = B(i,k) \cdot C(k,j)$$
$$A(i,l) = B(i,j,k) \cdot C(j,l) \cdot D(k,l)$$
$$a = B(i,j,k) \cdot C(i,j,k)$$
$$A(i,j,l) = B(i,j,k) \cdot C(k,l)$$
$$A(i,j) = B(i,j,k) \cdot c(k)$$

### Data Distribution
Partition A into tiles
Replicate B onto all nodes
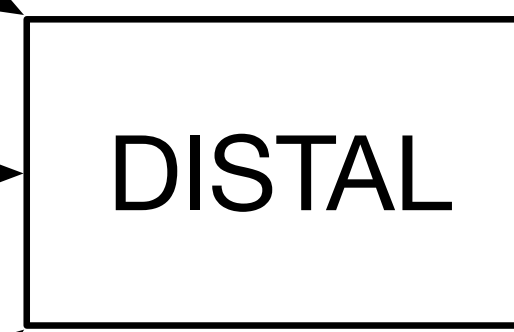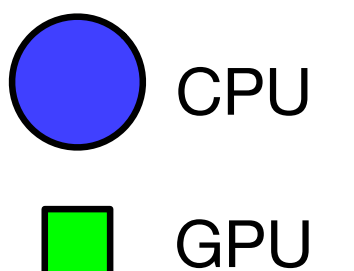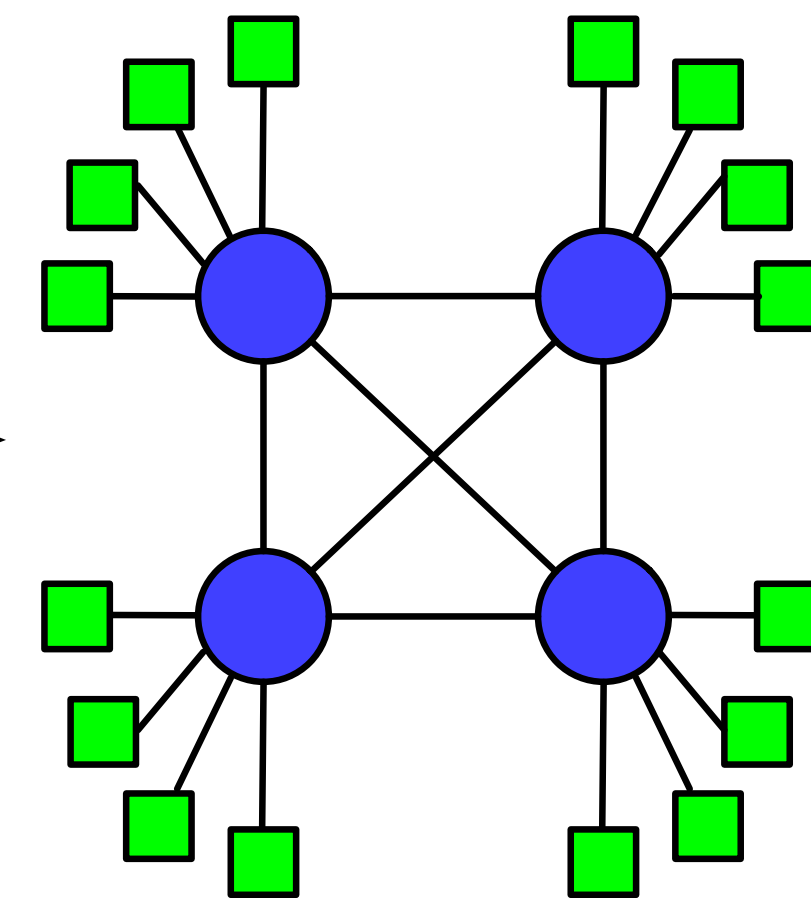Place C onto only some nodes

### Computation Distribution
Owner Computes
Distribute i,j loops
Communicate in chunks

DISTAL

### Supercomputer

● CPU

■ GPU

60

# Legion

Handles many features necessary for performance on modern machines
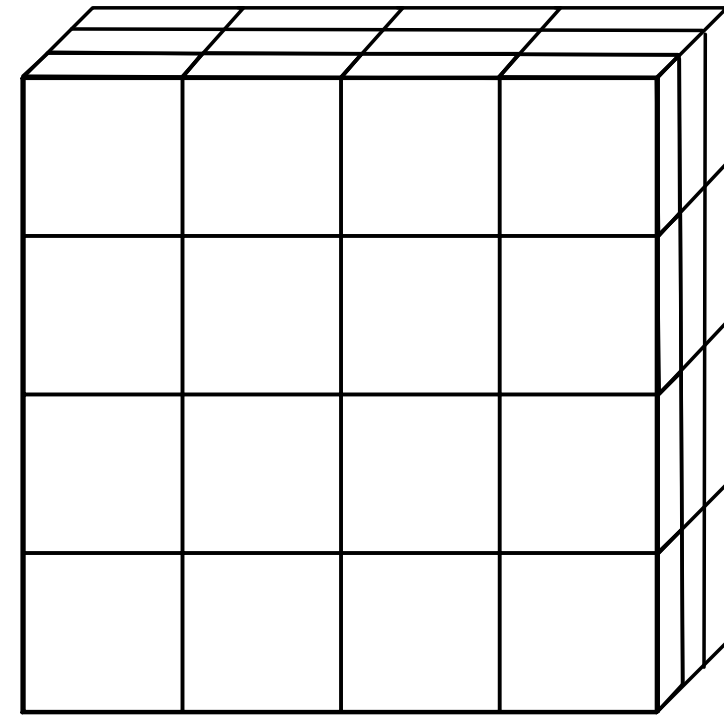
Overlap of communication and computation

Data movement through deep memory hierarchies

Native support for accelerators

Control over placement of computation and data
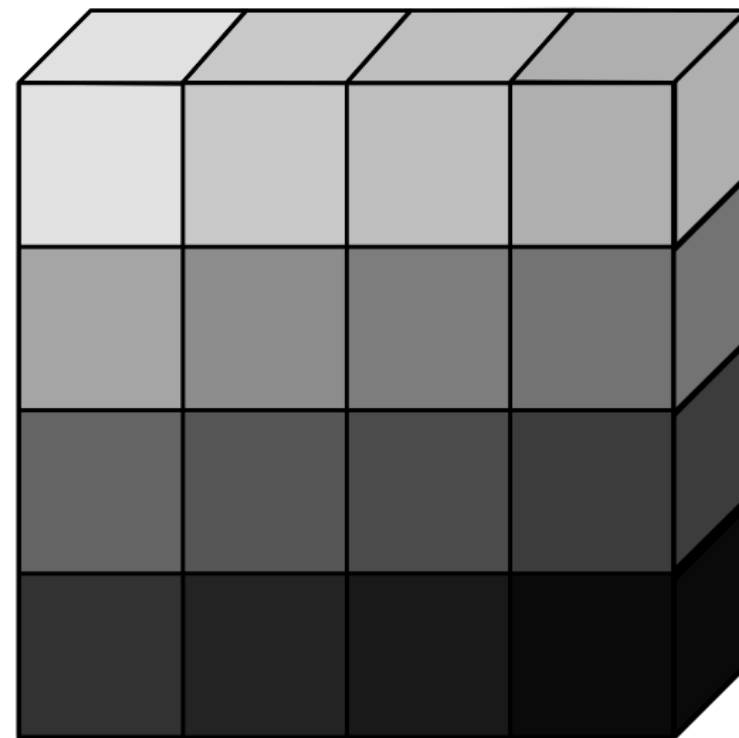
# Legion API



Regions

Tasks

```
void task(r1 : Region, r2 : Region, r3 : Region) { … }
```



Partitions

$$\mathcal{M}$$

Mapping

# Lowering Concrete Index Notation to Legion

$$\forall_i \ \forall_j \ A_{ij} \mathrel{+}= B_{ij}$$

$$\forall_{in}\forall_{jn}\forall_{il}\forall_{jl} \ A_{ij} \mathrel{+}= B_{ij} \ \text{s.t.} \ \text{div.}(i, in, il, gx), \text{div.}(j, jn, jl, gy), \text{dist.}([in, jn]), \text{comm.}([A, B], jn)$$

```
for in:
  for jn:
    for il:
      for jl:
```

# Lowering Concrete Index Notation to Legion

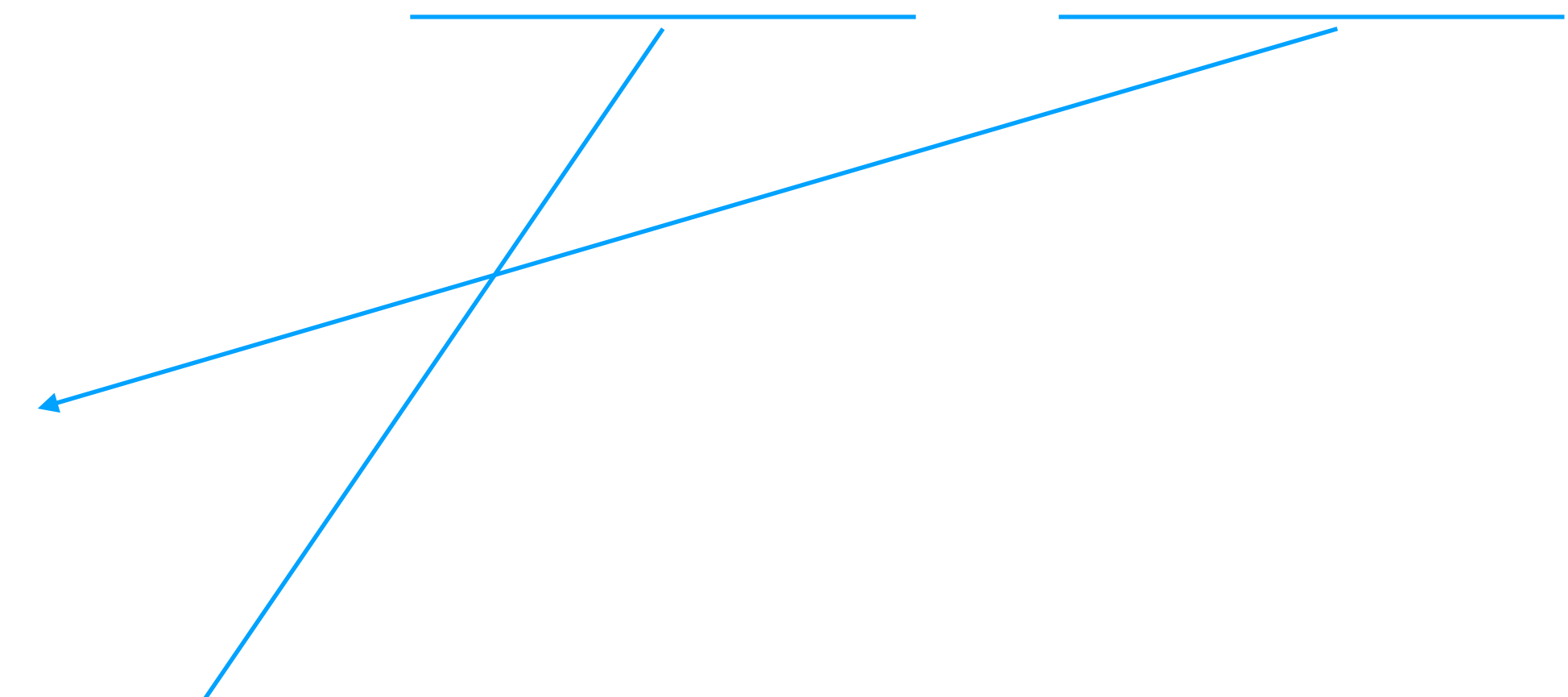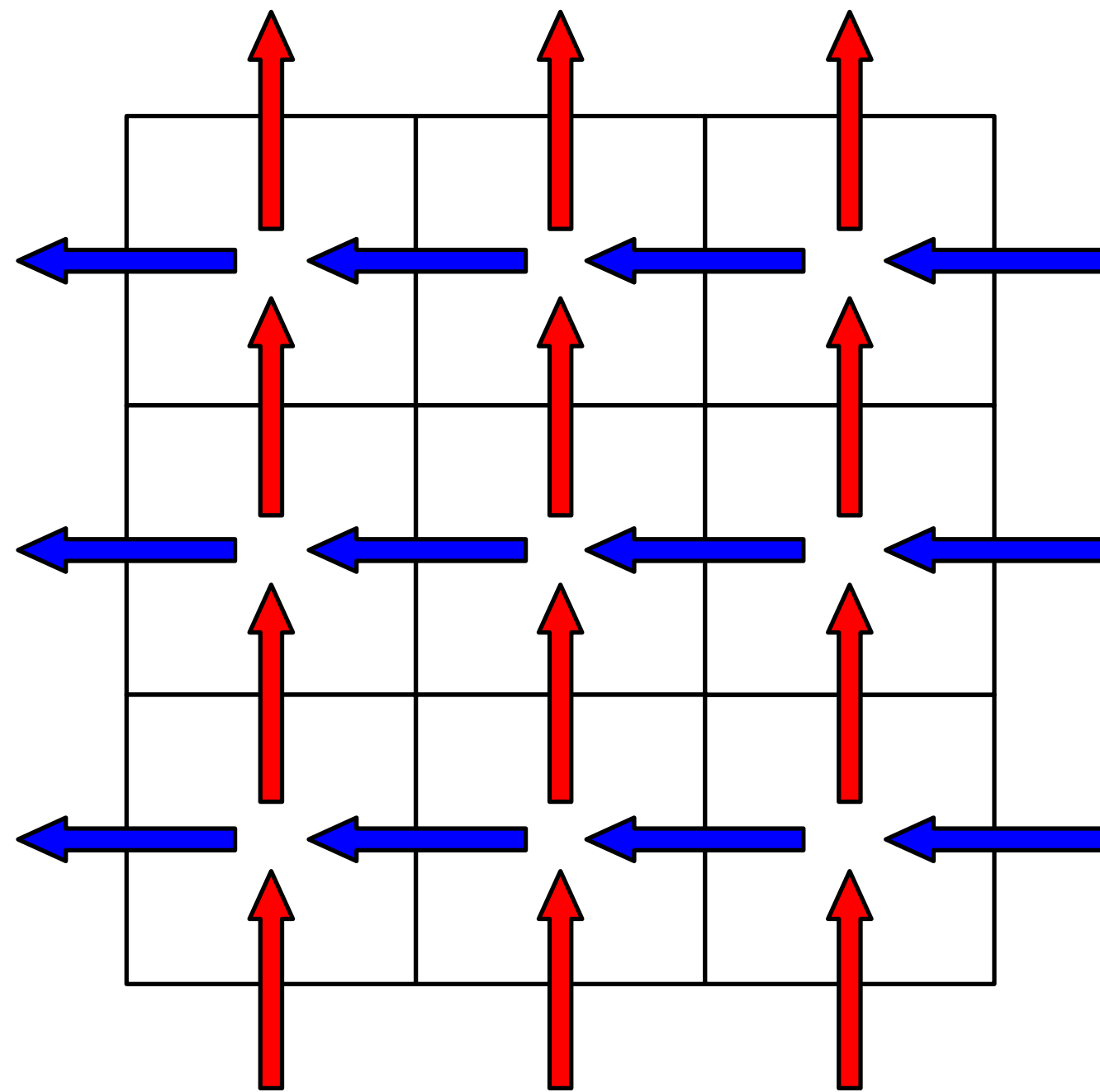$$\forall_{in}\forall_{jn}\forall_{il}\forall_{jl} \; A_{ij} \mathrel{+}= B_{ij} \; \text{s.t.} \; \text{div.}(i, in, il, gx), \text{div.}(j, jn, jl, gy), \text{dist.}([in, jn]), \text{comm.}([A, B], jn)$$

```
ind           for in:                                          art):
  A              for jn:
  B                 for il:
  f                    for jl:
                          i = in * (il.hi - il.lo) + il
                          j = jn * (jl.hi - jl.lo) + jl
                          A(i, j) += B(i, j)
```

# Cannon's Algorithm



```
# Arrange p processors into a 2D grid, √p × √p.
# Assign a tile of A, B, C to each processor.
# Perform an initial data shift.
for all P_ij in parallel:
  shift B_ij i spaces to the left
  shift C_ij j spaces upwards
for all P_ij in parallel:
  for k in (0, √p):
    A_ij += B_ij * C_ij
    shift B_ij to the left
    shift C_ij upwards
```

```
divide(i, il, in, gx)
divide(j, jl, jn, gx)
reorder({in, jn, il, jl})
split(k, ko, ki, chunkSize)
reorder(ko, il, jl, ki)
distribute(in, jn)
communicate(A, jn)
communicate({B, C}, ko)
```
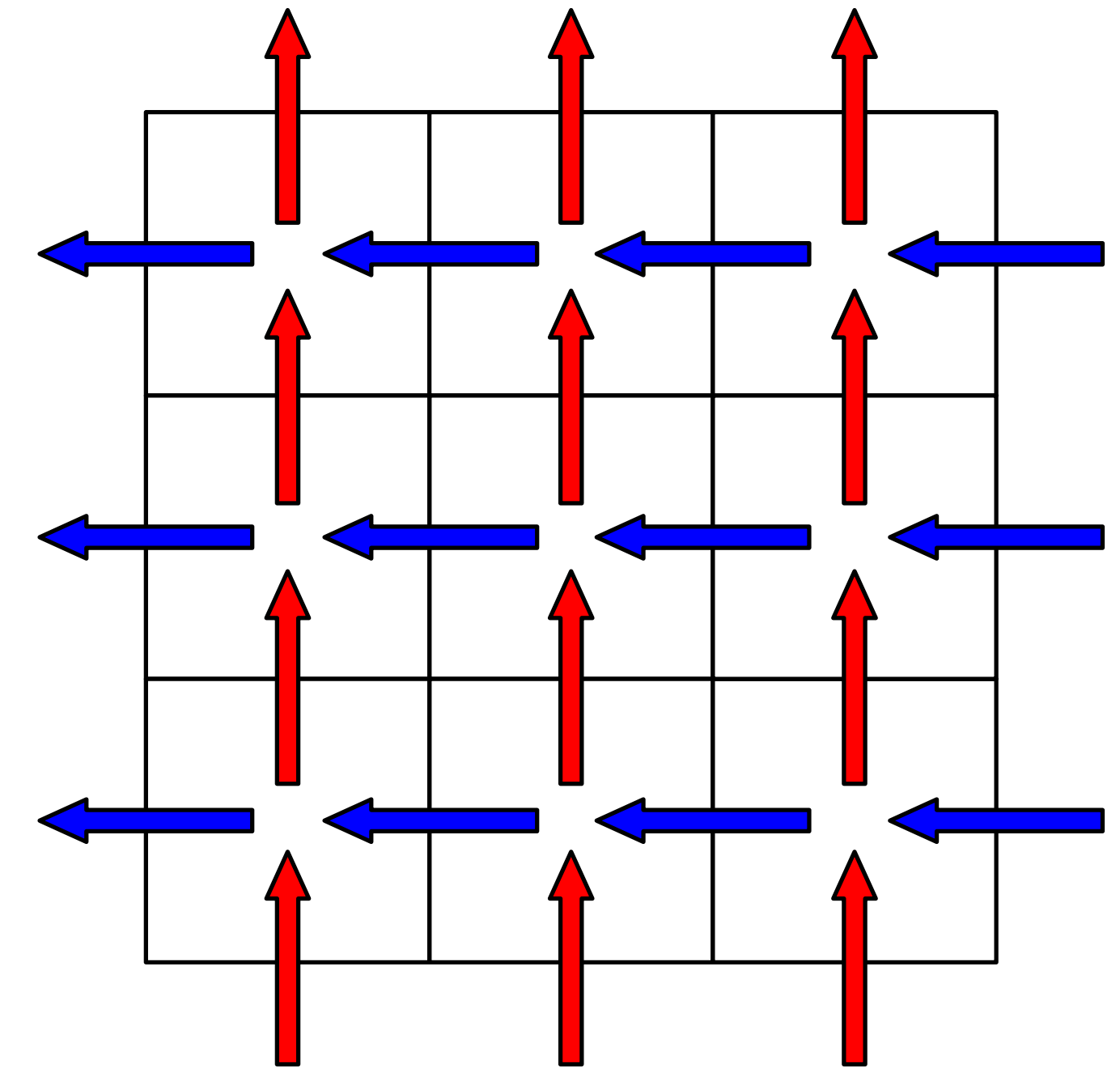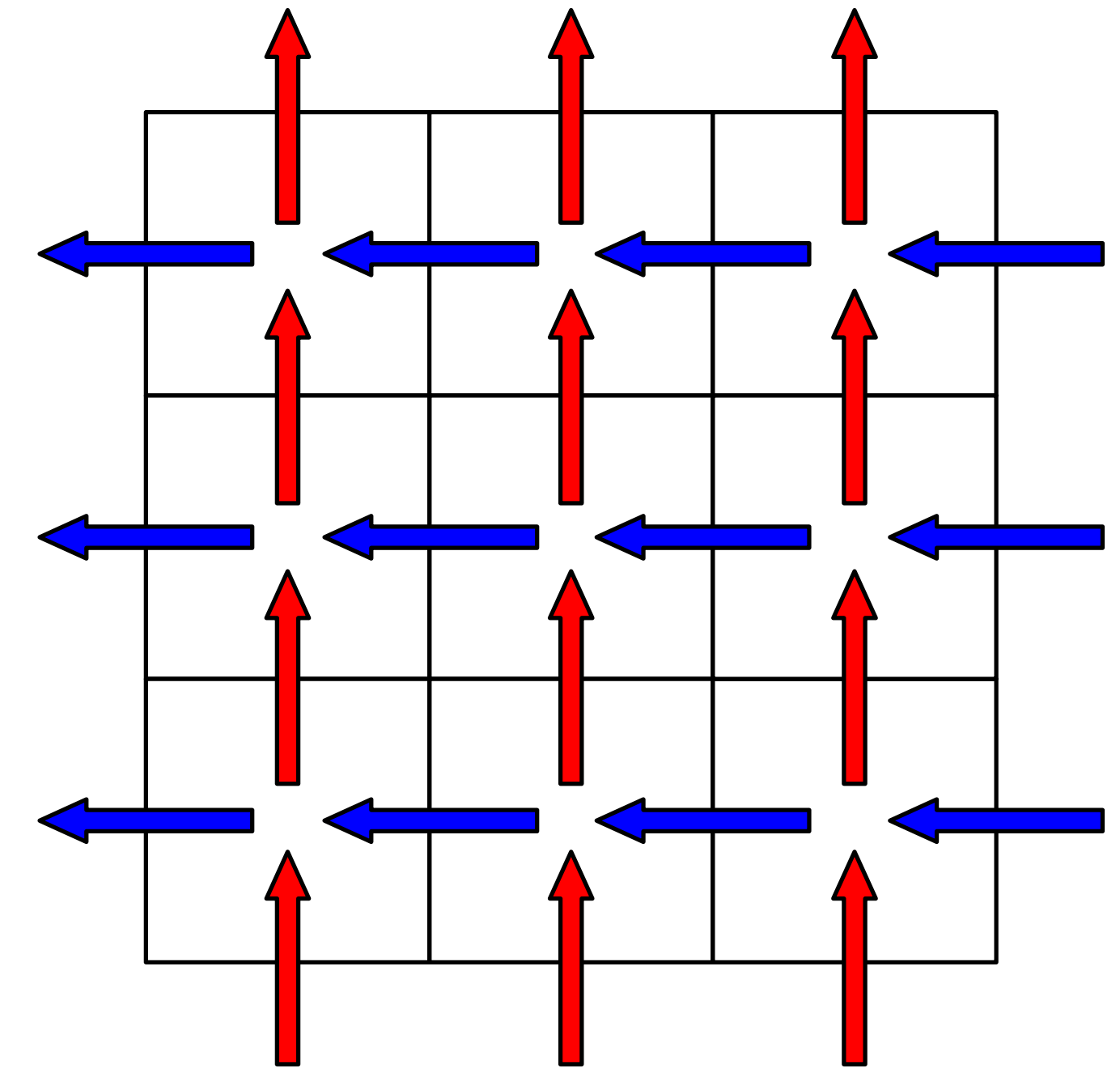
```
distributed for in, jn:
  communicate A
  for ko:
    communicate B, C
    for il:
      for jl:
        for ki:
          A(i, j) += B(i, k) * C(k, j)
```

```
divide(i, in, il, gx)
divide(j, jn, jl, gy)
divide(k, ko, ki, gx)
rotate(ko, {in, jn}, kr)
reorder({in, jn, kr, il, jl, ki})
distribute(in, jn)
communicate(A, jn)
communicate({B, C}, kr)
```

```
distributed for in, jn:
  communicate A
  for kr:
    communicate B, C
    for il:
      for jl:
        for ki:
          A(i, j) += B(i, k) * C(k, j)
```

# Emergence of systolic communication

$$A(i, j) = B(i, k) * C(k, j)$$

$$kr = ko + in + jn \bmod 3$$
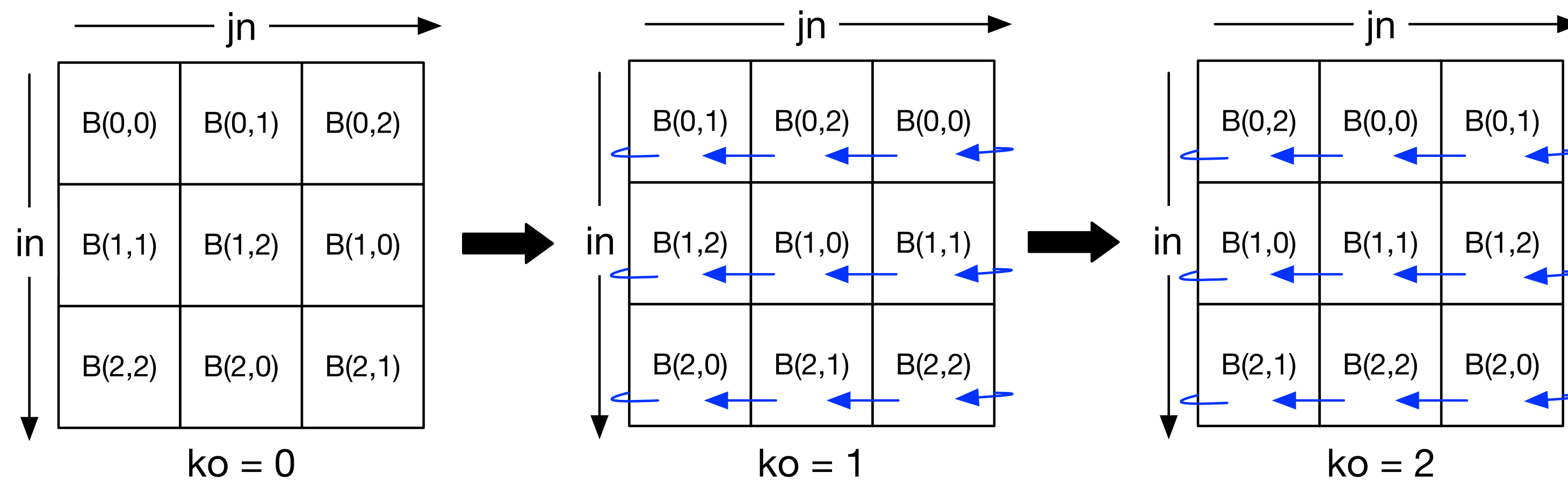
```
distributed for in, jn:
  communicate A
  for kr:
    communicate B, C
    for il:
      for jl:
        for ki:
          A(i, j) += B(i, k) * C(k, j)
```

# GEMM (CPU)



Legend: COSMA ● | CTF ● | ScaLAPACK ● | Our Cannon's ▲ | Our SUMMA ■ | Our PUMMA ◆ | Our Solomonik's ★ | Our Johnson's ✕ | Our COSMA ⬠ | Peak Utilization ▬

Y-axis: GFLOP/s

X-axis: Nodes (CPU Cores)
1 (40), 2 (80), 4 (160), 8 (320), 16 (640), 32 (1280), 64 (2560), 128 (5120), 256 (10240)