

Software and Hardware for Sparse ML

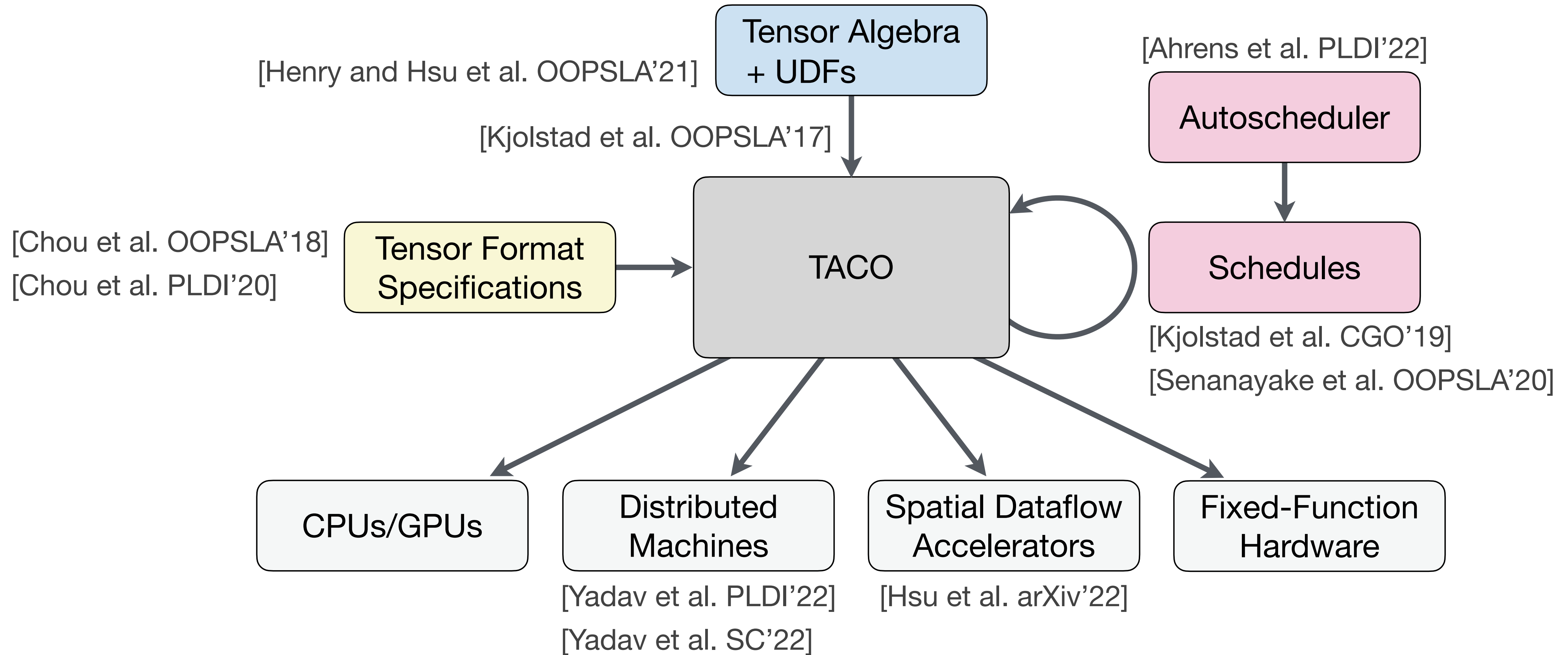
Fredrik Kjolstad

Assistant Professor, Department of Computer Science

Stanford University



Background: Sparse Tensor Algebra Compilation



Other systems:

COMET [Mutlu et al. LCPC'20]

SPF [Zhao et al. arXiv'22]

MLIR SparseTensor Dialect [Bik et al. TACO'22]

SparseTIR [Ye et al. arXiv'22]

Overview

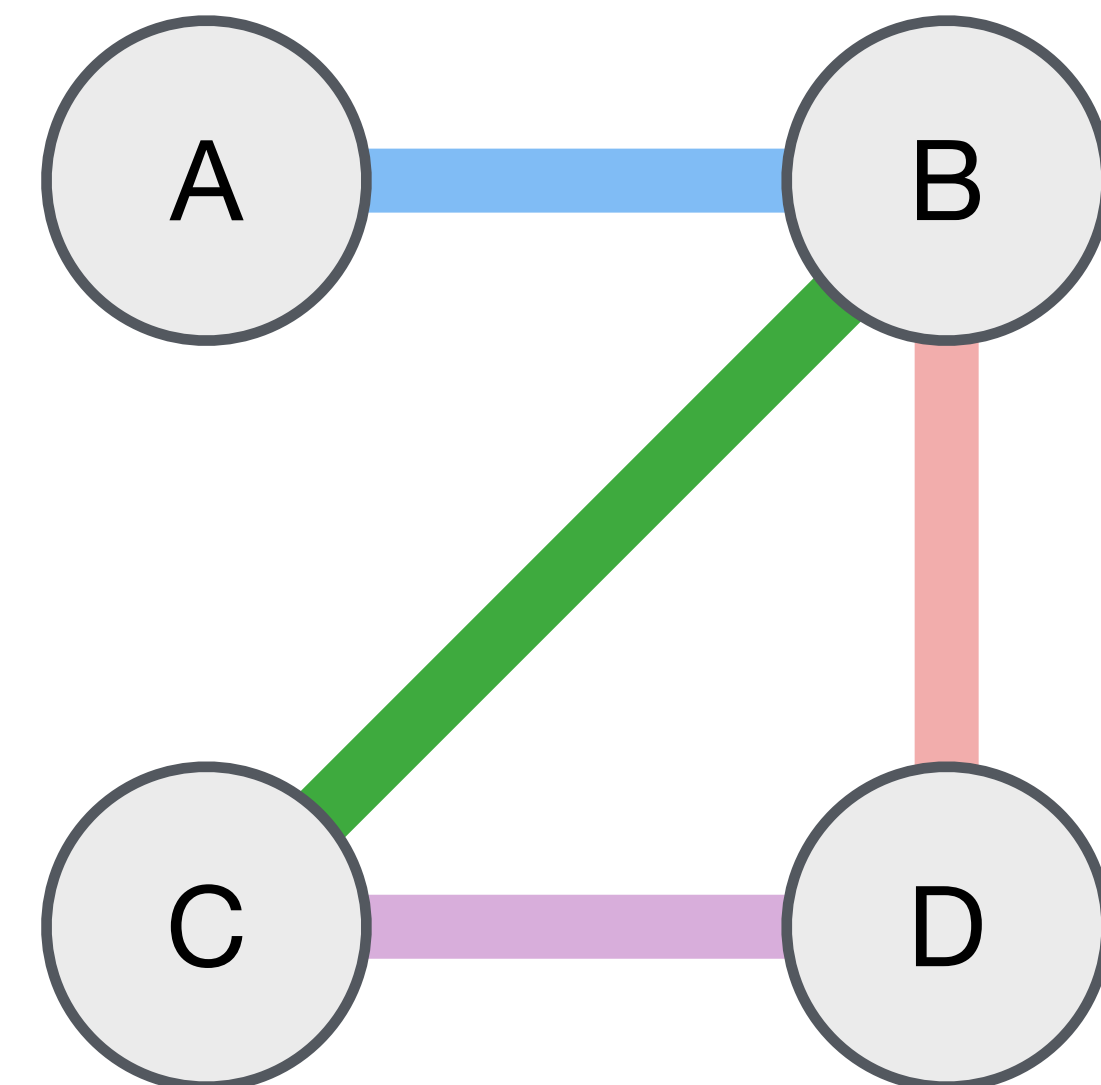
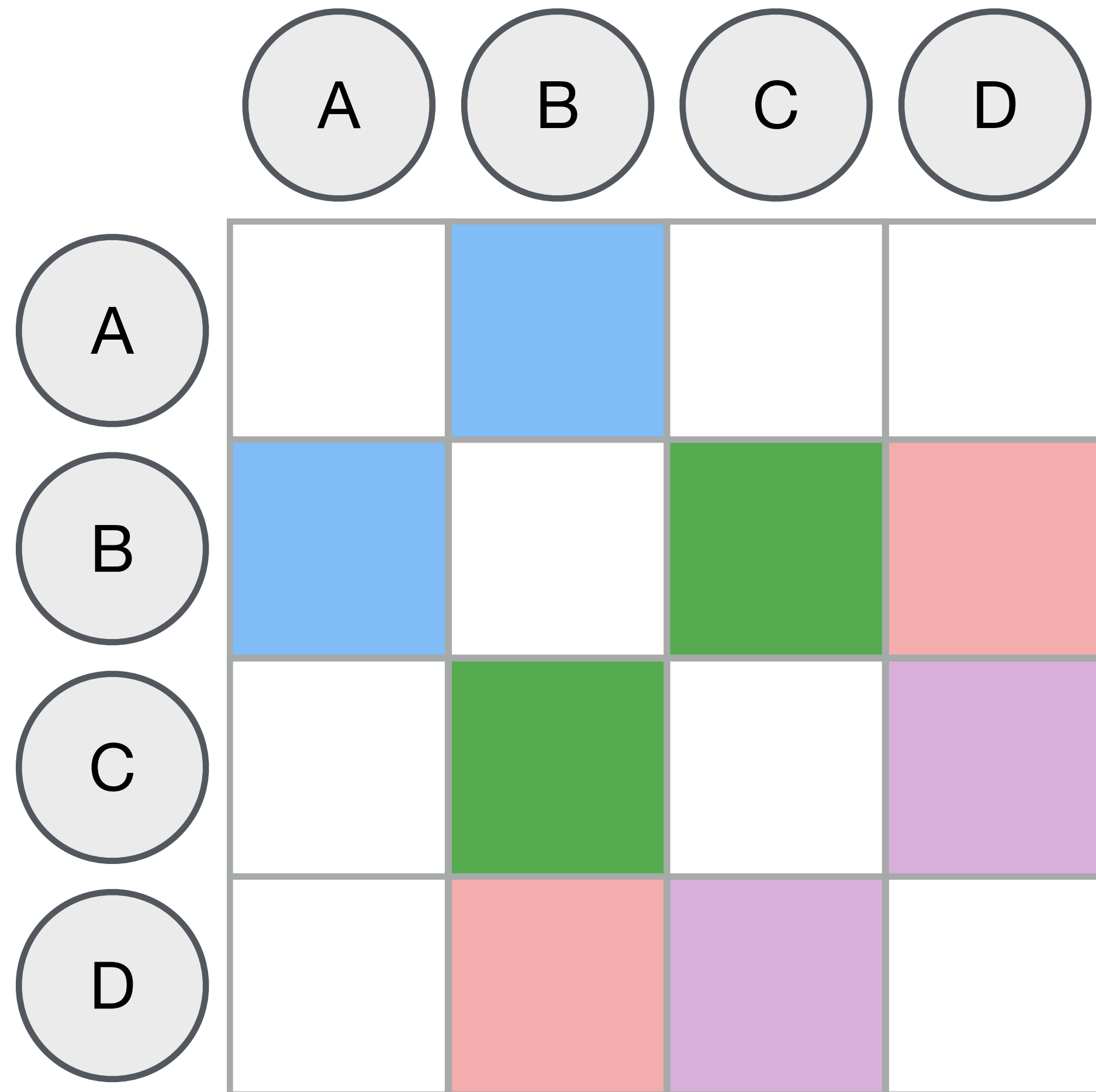
My view of sparsity

Why sparsity requires compilers and general hardware

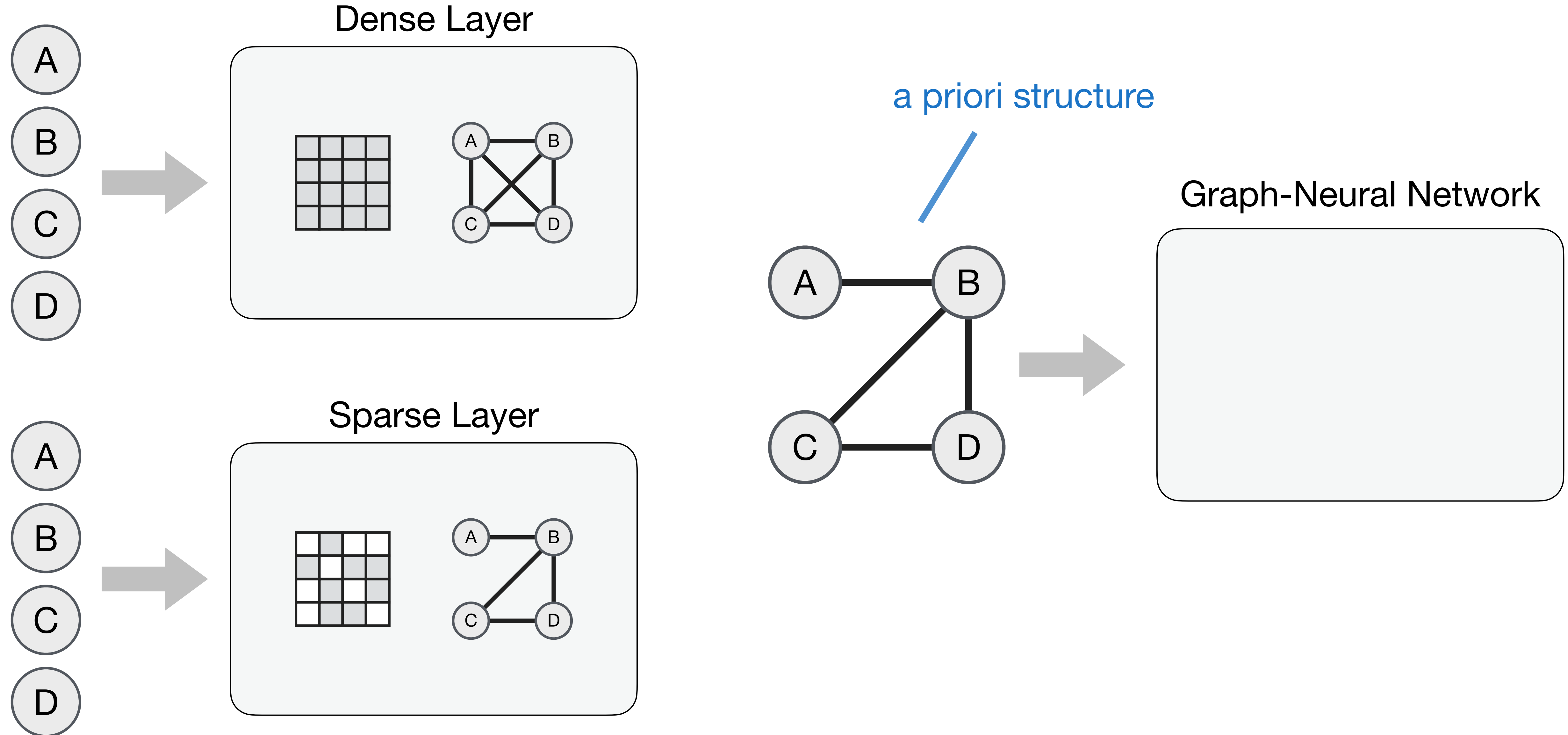
Thesis

Unlike dense neural networks that can be reduced to GEMM,
it will not be possible to reduce sparse neural networks
to one optimized function

Sparsity as system connectivity

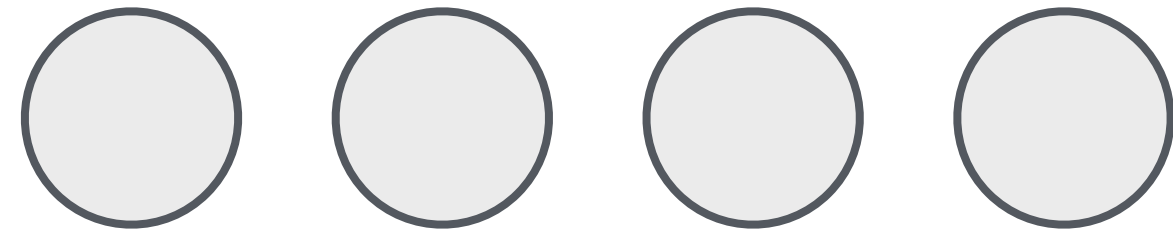


Two types of sparsity: learned and a priori



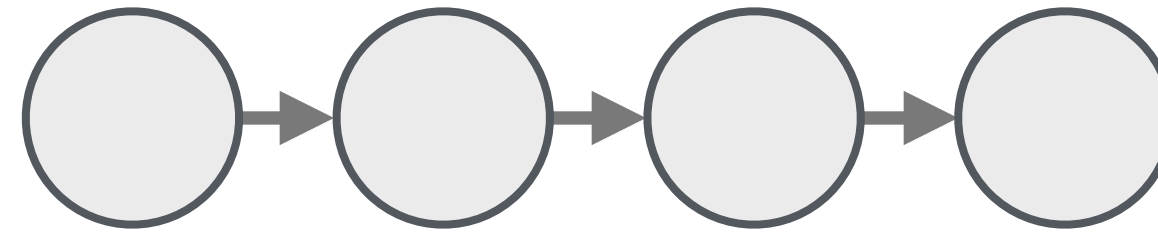
Structure of input data

Sets



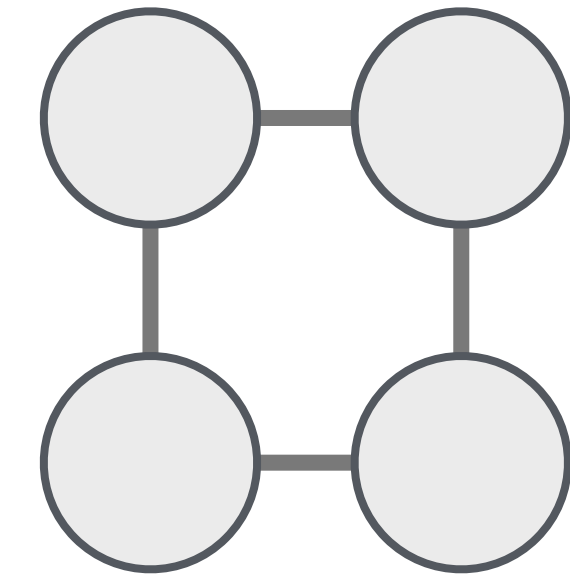
- No a priori connections
- Postulate fully connected layer
- Then may try to learn sparsity

Sequences



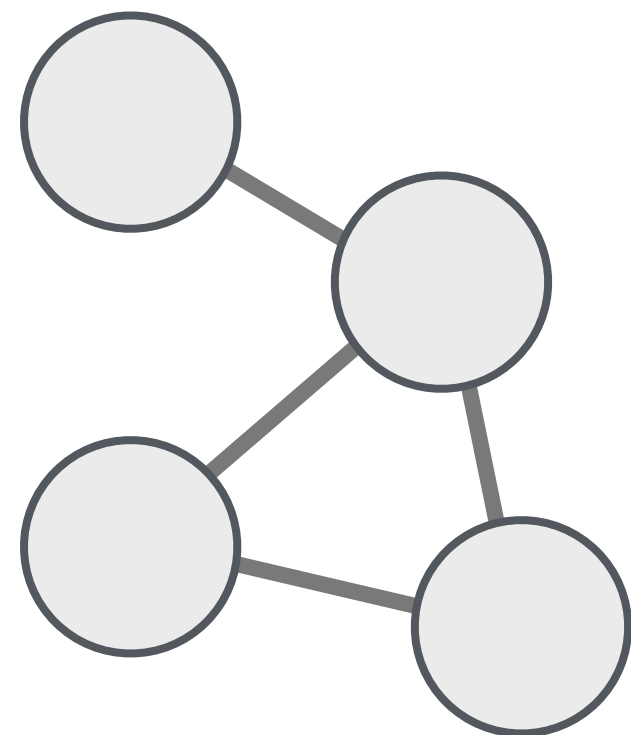
- Triangular matrices in transformers
- Recurrences in RNNs

Grids



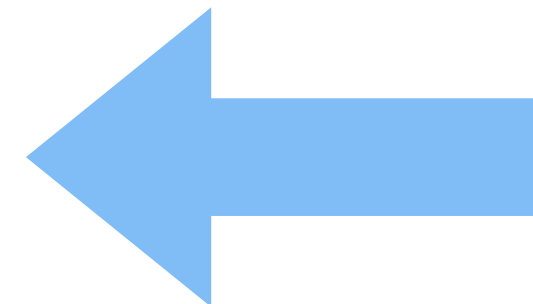
- Pixel locality in CNNs

Graphs



Relational Data

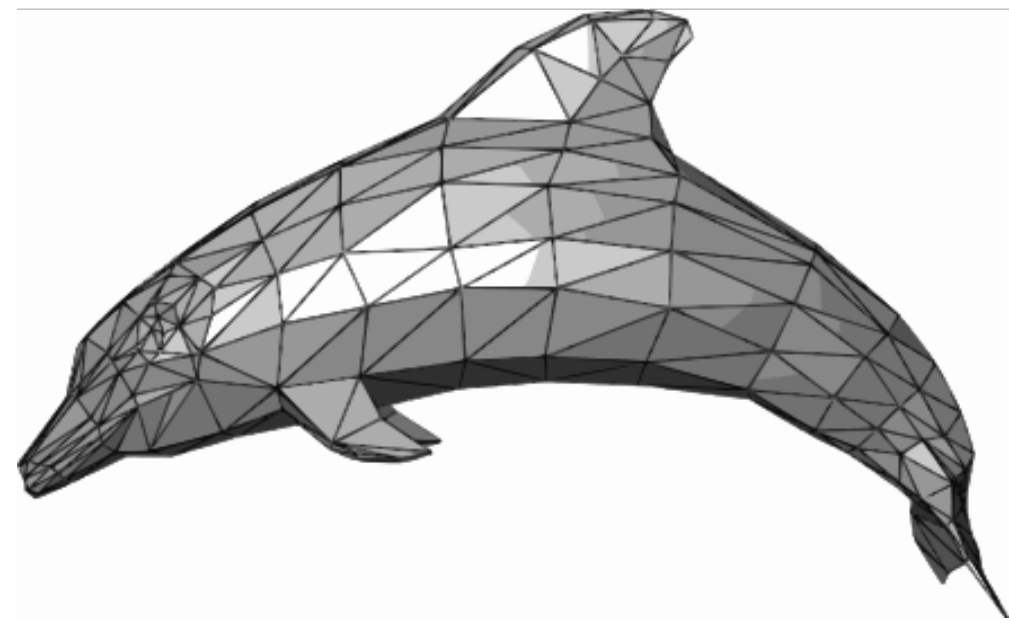
John	Jill
Jill	Kim
Kim	Mary
Mary	Kim



How sparse is graph/relational data? Often asymptotically sparse.

Assume an average degree of 150 (e.g., 150 friends)

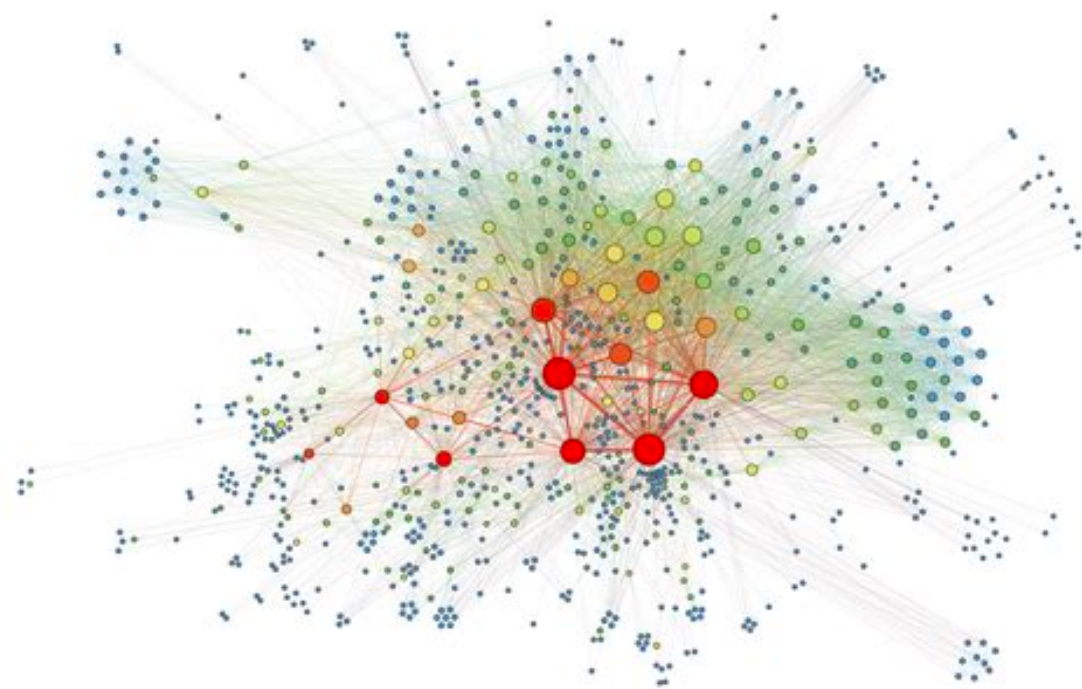
Conditioned Meshes



Each matrix row then has 150 nonzeros

$$\text{At 10,000 rows: } \frac{150 \cdot 10,000}{10,000^2} = 1.5 \% \text{ nonzeros}$$

Power-law graphs



$$\text{At 100,000 rows: } \frac{150 \cdot 100,000}{100,000^2} = 0.15 \% \text{ nonzeros}$$

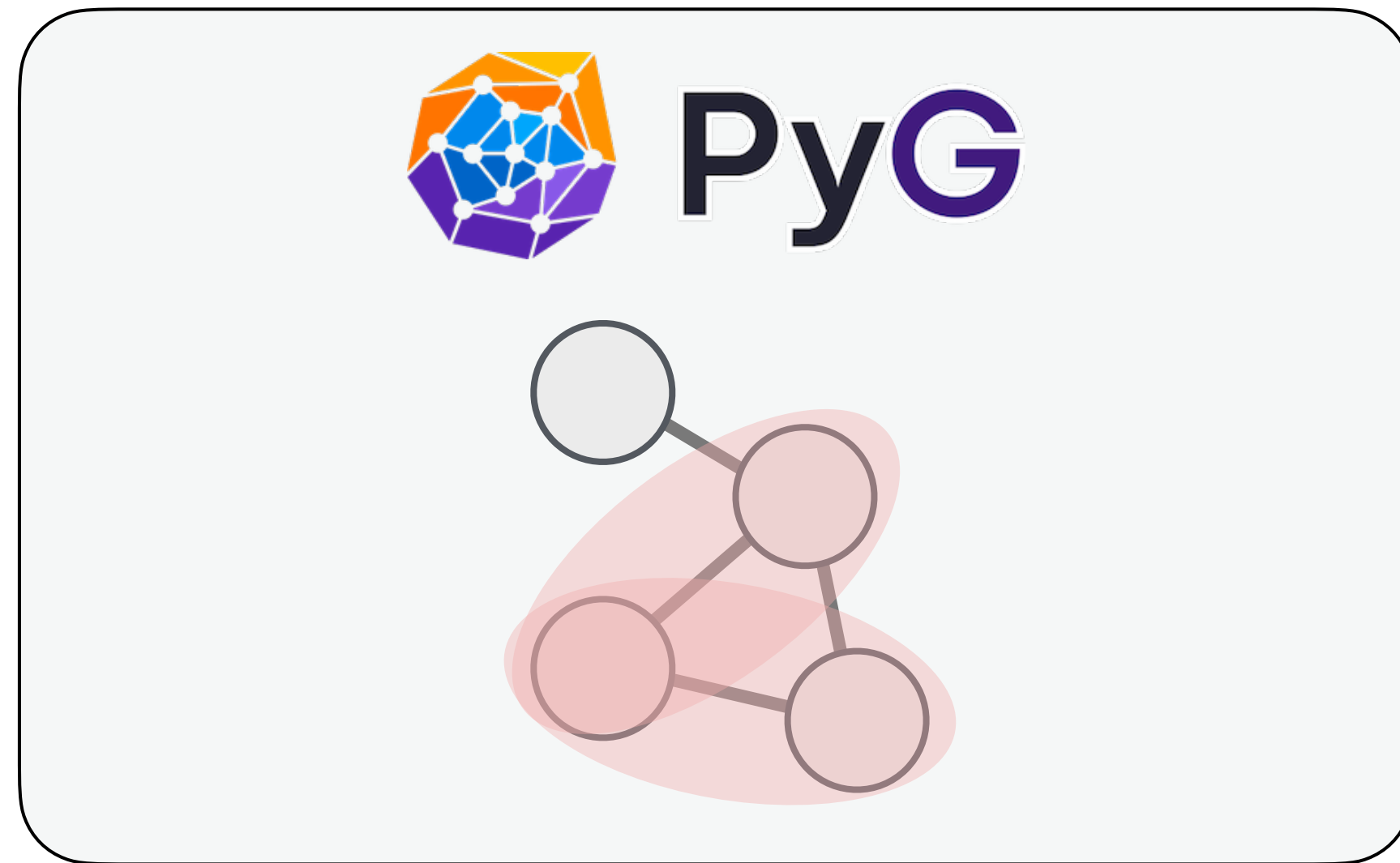
Matrix components: $O(n^2)$

Nonzeros: $O(n)$

Fraction of nonzeros: $O(1/n)$

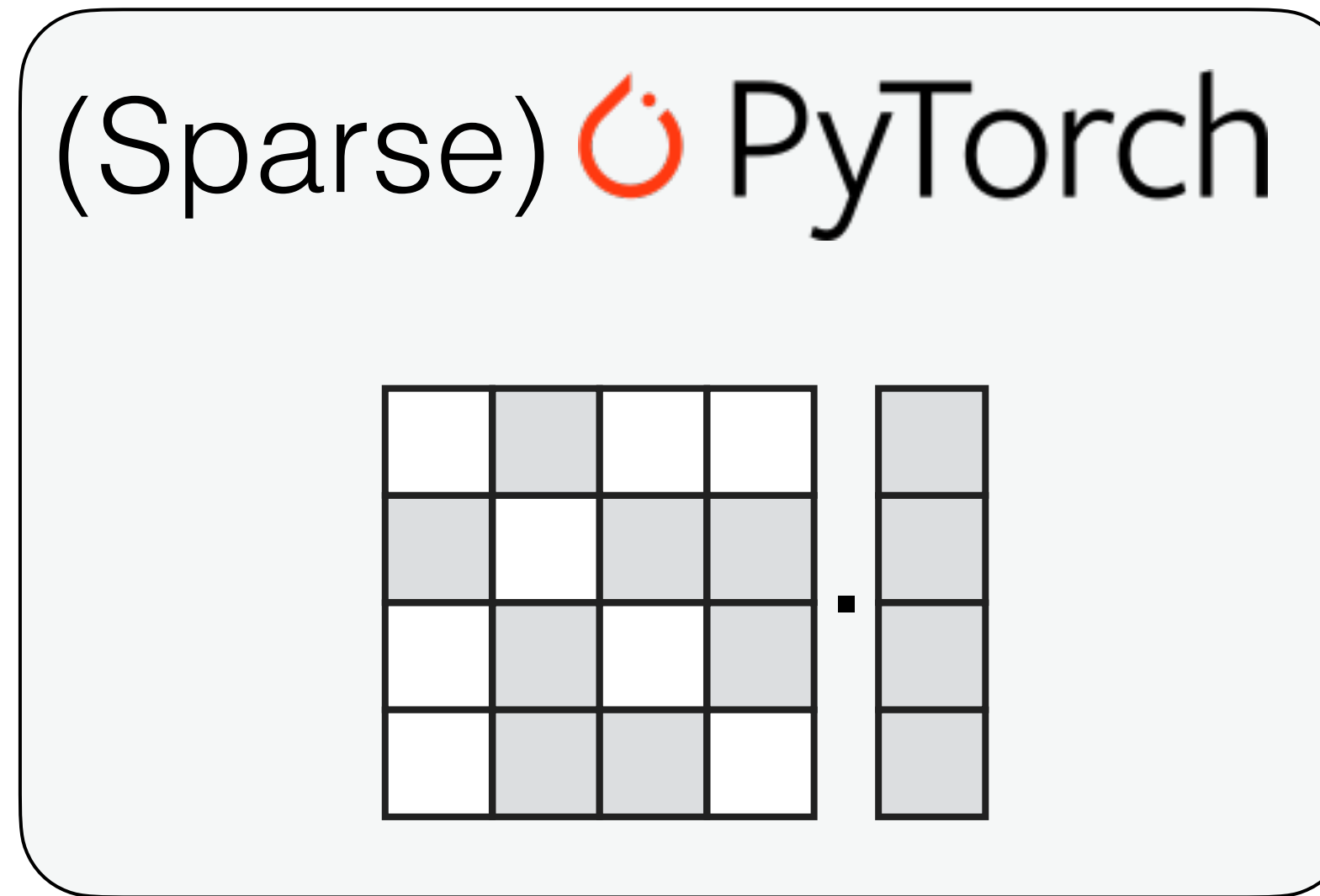
Should we build graph frameworks or sparse tensor frameworks?

Graph Framework



- Lower level
- Can hand-implement fusion

Sparse Tensor Framework



- Higher-level natural notation
 - Same notation as in papers
 - Composes with dense/CNNs
 - Easy to compose multiple graphs
- Compiler can
 - Fuse computation
 - Reorder and tile
 - Port across machines

How general should tensor frameworks be?

Kernel library

- Fixed number of hand-optimized operations
- Fixes tensor formats

Full tensor support

- General tensor algebra
- User-defined functions
- Tensor reshapes and composition
- Portable across tensor formats

```
class SAGEConv ( in_channels: Union[int, Tuple[int, int]], out_channels: int, aggr: Optional[Union[str, List[str], Aggregation]] = 'mean', normalize: bool = False, root_weight: bool = True, project: bool = False, bias: bool = True, **kwargs ) [source]
```

```
class ChebConv ( in_channels: int, out_channels: int, K: int, normalization: Optional[str] = 'sym', bias: bool = True, **kwargs ) [source]
```

The chebyshev spectral graph convolutional operator from the “Convolutional Neural Networks on Graphs with Fast Localized Spectral Filtering” paper

$$\mathbf{X}' = \sum_{k=1}^K \mathbf{Z}^{(k)} \cdot \Theta^{(k)}$$

where $\mathbf{Z}^{(k)}$ is computed recursively by

$$\begin{aligned} \mathbf{Z}^{(1)} &= \mathbf{X} \\ \mathbf{Z}^{(2)} &= \hat{\mathbf{L}} \cdot \mathbf{X} \\ \mathbf{Z}^{(k)} &= 2 \cdot \hat{\mathbf{L}} \cdot \mathbf{Z}^{(k-1)} - \mathbf{Z}^{(k-2)} \end{aligned}$$

and $\hat{\mathbf{L}}$ denotes the scaled and normalized Laplacian $\frac{2\mathbf{L}}{\lambda_{\max}} - \mathbf{I}$.

ator from the “Inductive Representation Learning on Large Graphs” paper

$$\mathbf{x}'_i = \mathbf{W}_1 \mathbf{x}_i + \mathbf{W}_2 \cdot \text{mean}_{j \in \mathcal{N}(i)} \mathbf{x}_j$$

in \mathbf{x}_j will first get projected via

$$\mathbf{x}_j \leftarrow \sigma(\mathbf{W}_3 \mathbf{x}_j + \mathbf{b})$$

```
class GraphConv ( in_channels: Union[int, Tuple[int, int]], out_channels: int, aggr: str = 'add', bias: bool = True, **kwargs ) [source]
```

The graph neural network operator from the “Weisfeiler and Leman Go Neural: Higher-order Graph Neural Networks” paper

$$\mathbf{x}'_i = \mathbf{W}_1 \mathbf{x}_i + \mathbf{W}_2 \sum_{j \in \mathcal{N}(i)} e_{j,i} \cdot \mathbf{x}_j$$

edge weight from source node j to target node i (default: 1)

```
class GCNConv ( in_channels: int, out_channels: int, improved: bool = False, cached: bool = False, add_self_loops: bool = True, normalize: bool = True, bias: bool = True, **kwargs ) [source]
```

The graph convolutional operator from the “Semi-supervised Classification with Graph Convolutional Networks” paper

$$\mathbf{X}' = \hat{\mathbf{D}}^{-1/2} \hat{\mathbf{A}} \hat{\mathbf{D}}^{-1/2} \mathbf{X} \Theta,$$

where $\hat{\mathbf{A}} = \mathbf{A} + \mathbf{I}$ denotes the adjacency matrix with inserted self-loops and $\hat{D}_{ii} = \sum_{j=0} \hat{A}_{ij}$

How to implement a sparse PyTorch in software and hardware

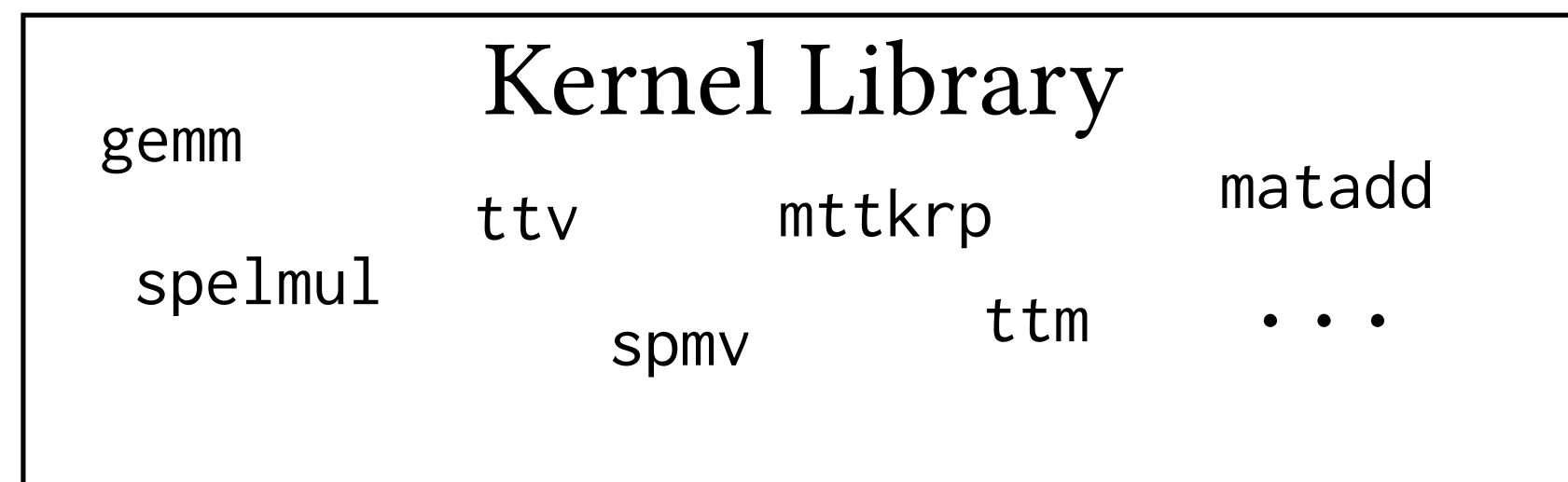
Factorization

$$A = B \odot (CD)$$

↓ factorize

Matrix T = gemm(C,D);

Matrix A = spelmul(B,T);



Compilation

$$A = B \odot (CD)$$

↓ compile

```
int pA2 = 0;
for (int pB1 = B1_pos[0];
     pB1 < B1_pos[1]; pB1++) {
    int i = B1_crd[pB1];
    for (int pB2 = B2_pos[pB1];
         pB2 < B2_pos[pB1+1]; pB2++) {
        int j = B2_crd[pB2];
        double t = 0.0;
        for (int k = 0; k < 0; k++) {
            int pC2 = i * 0 + k;
            int pD2 = k * N + j;
            t += C[pC2] * D[pD2];
        }
        A[pA2++] = B[pB2] * t;
    }
}
```

Factorization in **dense** tensor algebra

$$a = \sum_{ijklmnop} M_{ij} P_{jk} M_{kl} P_{lm} \overline{M_{nm}} \overline{P_{no}} \overline{M_{po}} \overline{P_{ip}}$$



Transposes → GEMM → Transposes → GEMM → ...

- Works pretty well for dense (at least on shared memory machines)
- The cost of transpose is modest
- The benefit of handwritten GEMM is large

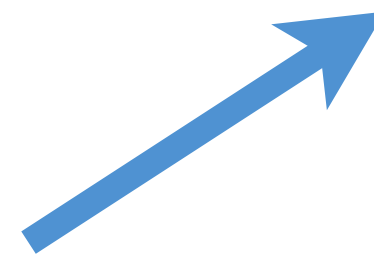
Unlike dense neural networks that can be reduced to GEMM,
it will not be possible to reduce sparse neural networks
to one optimized function

Factorization in **sparse** tensor algebra

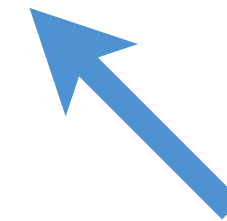
$$a = \sum_{ijklmnop} M_{ij} P_{jk} M_{kl} P_{lm} \overline{M_{nm}} \overline{P_{no}} \overline{M_{po}} \overline{P_{ip}}$$



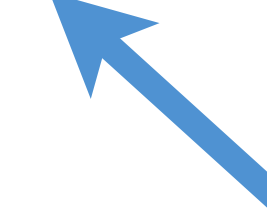
Transposes → SpGEMM → Transposes → SpGEMM → ...



Requires sorting and
potentially data structure
conversion



Compilers better able to
produce competitive code



Need to flatten data structures
(e.g., COO triplets to pairs)

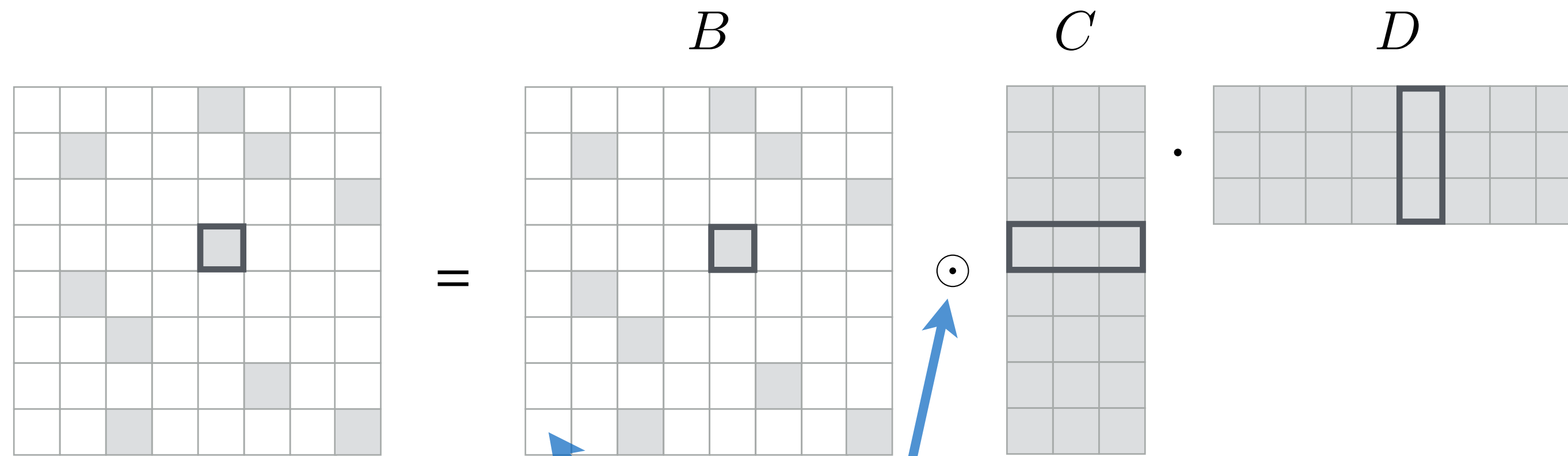
No fusion across operation



Potential asymptotic complexity slowdown!

Factorization destroys fusion

Sampled Dense-Dense Matrix Multiplication (SDDMM)



$O(IJK)$

This dot product need not be computed

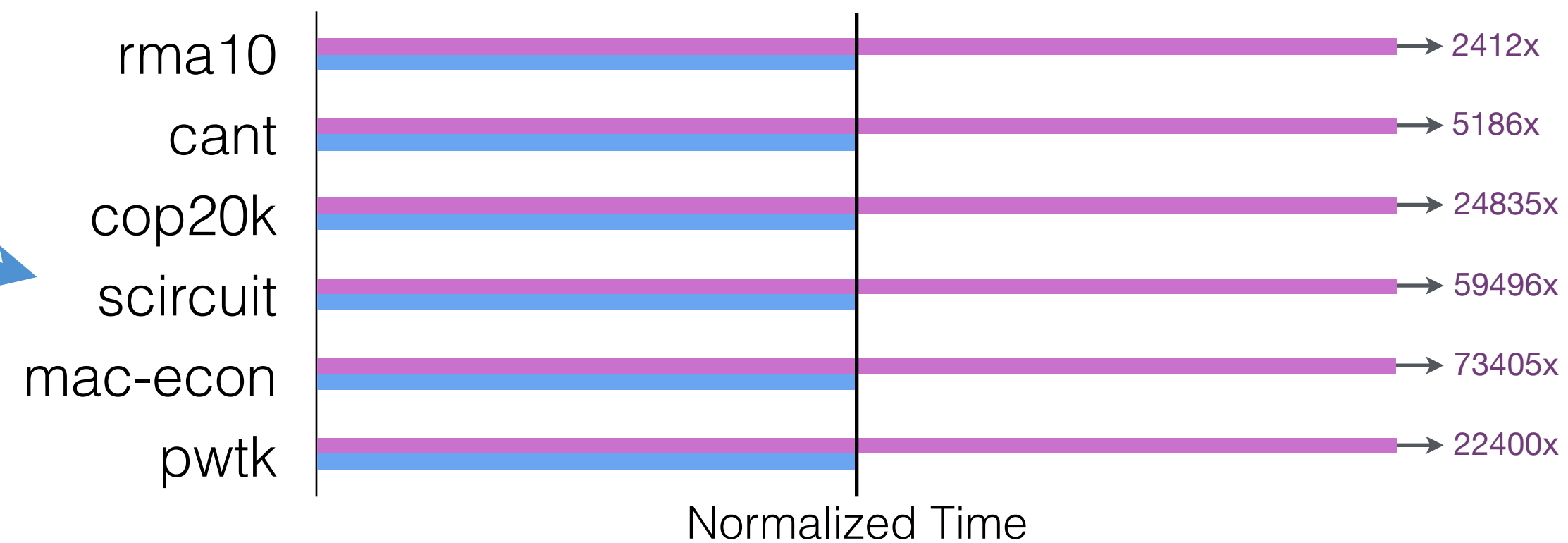
64 inner product
10 inner product

$O(\text{NNZ}_B \cdot K)$

Separate Operations

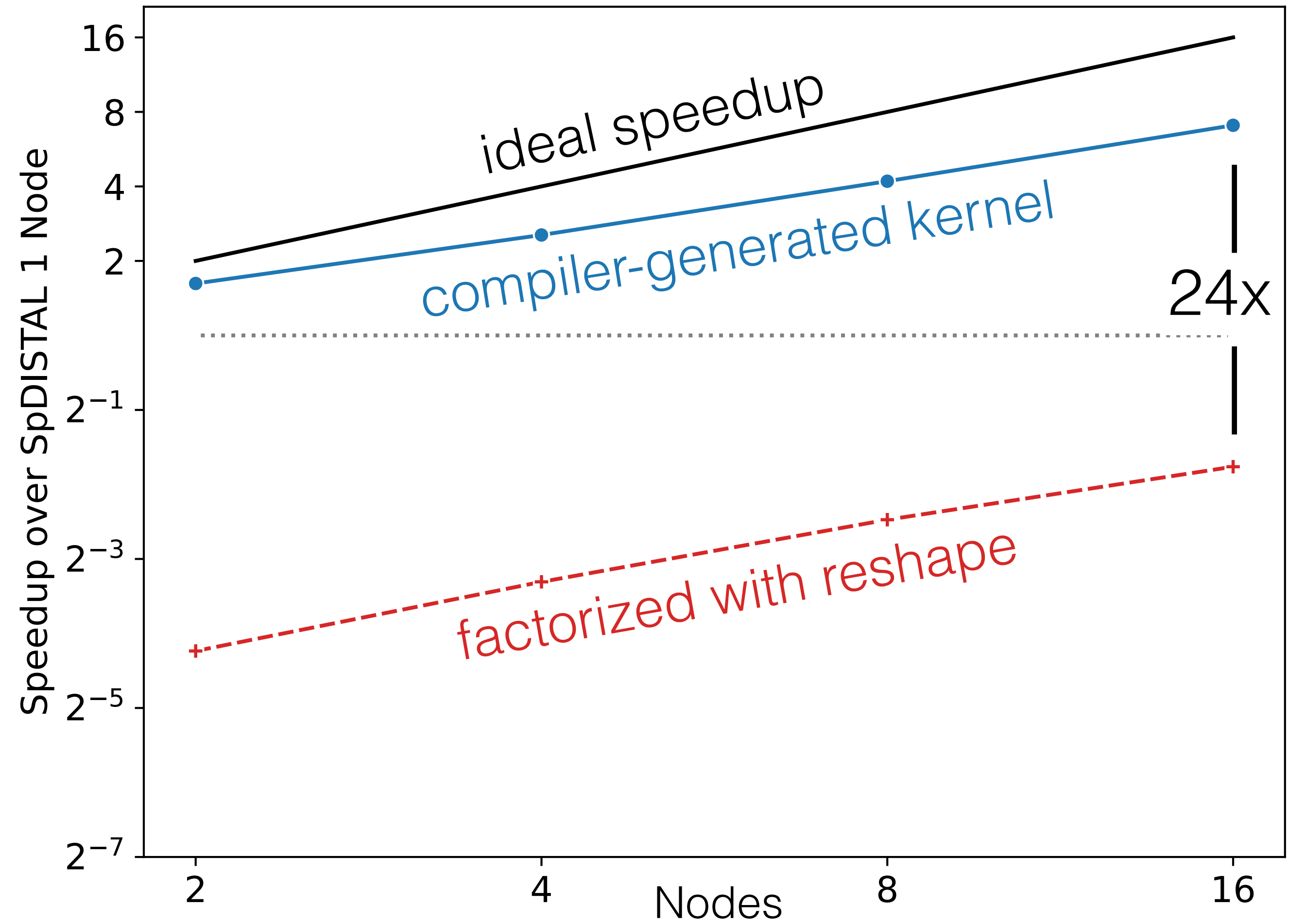
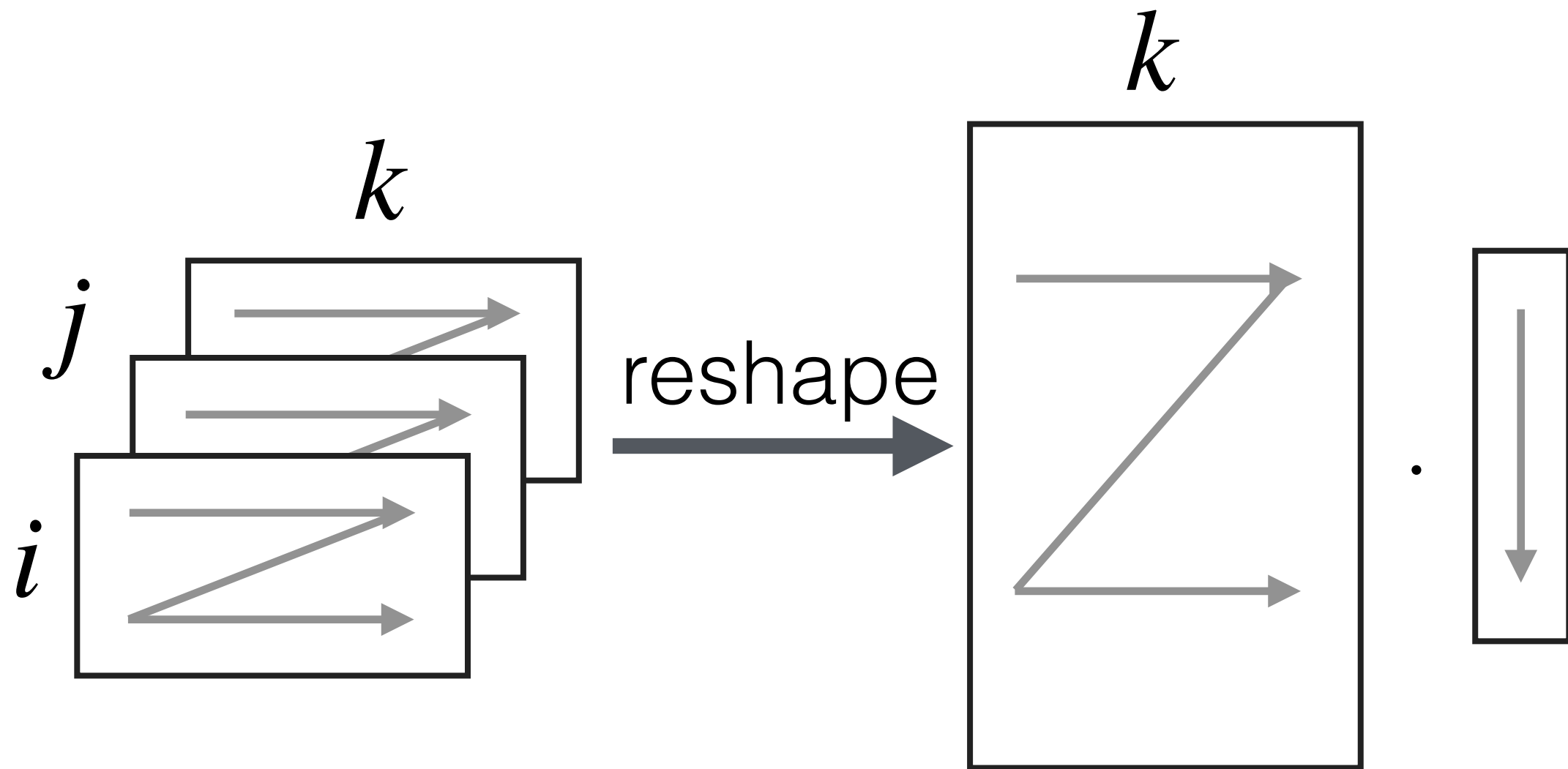
Fused Operation

Sparse matrices



Factorization forces data movement

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



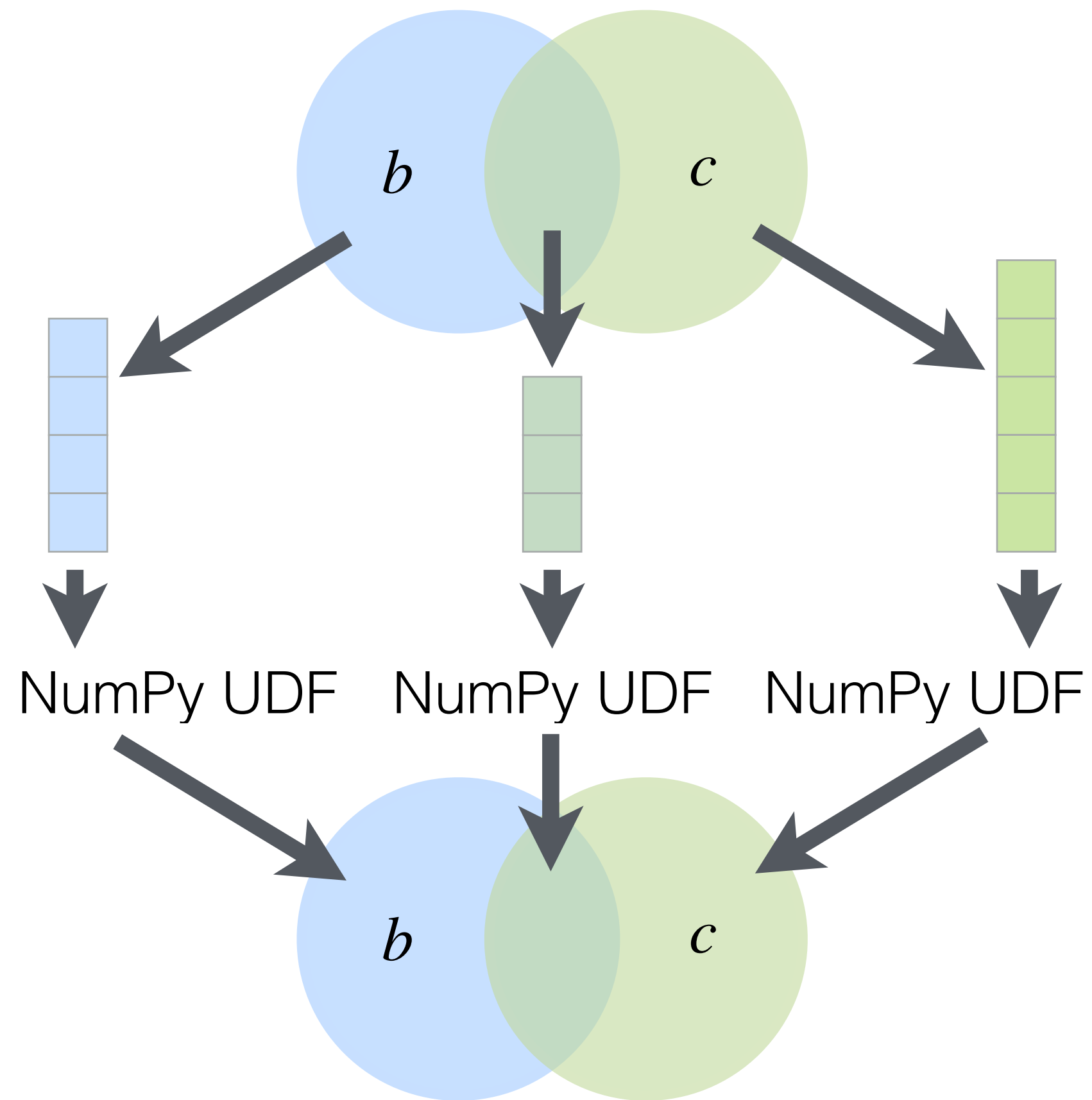
DISTAL

[Yadav et al. PLDI'22]

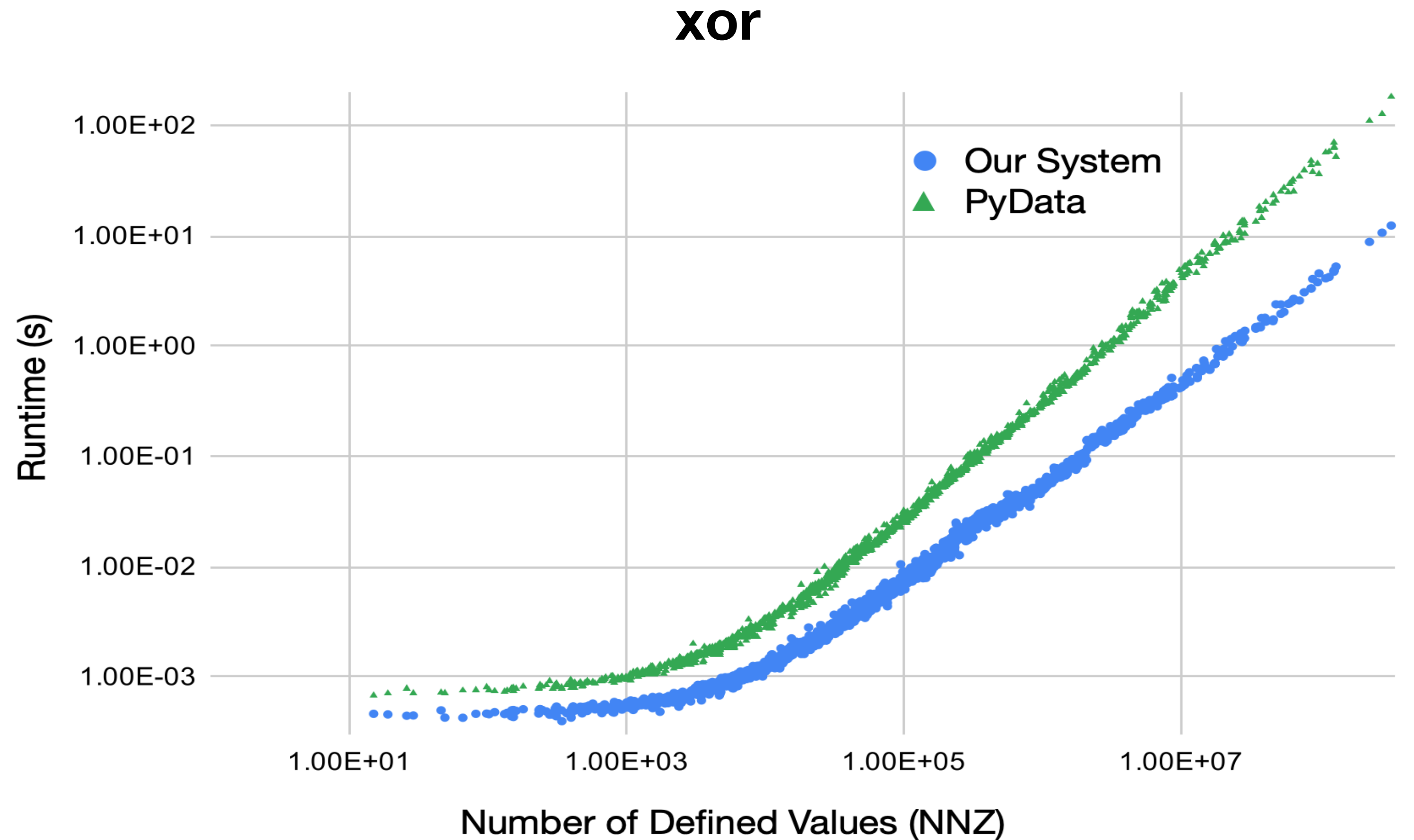
[Yadav et al. SC'22]



Factorization prevents efficient user-defined function support



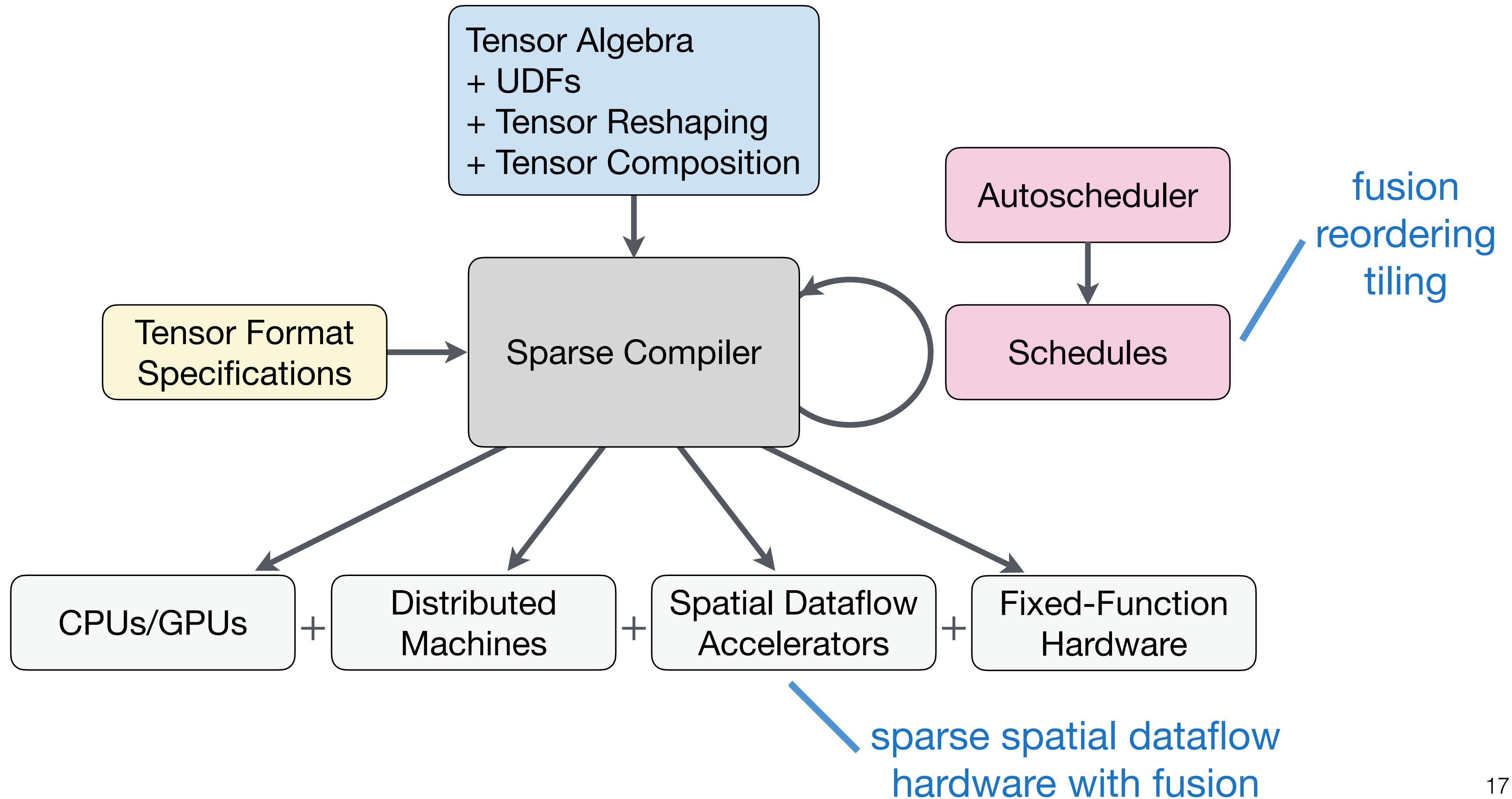
Average: 7.5x cost factorization



[Henry and Hsu et al. OOPSLA'21]



Compiler design for general sparse tensor operations



Hardware design for general sparse tensor operations

Sparse tensor algebra accelerators must support:

1. **Generality:** arbitrary tensor algebra operations
2. **Data Structures:** dense and sparse data structures
3. **Fusion:** Fusion across operations
4. **Reordering:** Changing the order they process tensor dimensions

The Sparse Abstract Machine

- Abstract spatial dataflow machine architecture
- Supports all four properties (generality, data structures, fusion, and reordering)
- Also supports tiling, parallelization, vectorization, and bitvector wire protocols
- Implemented in a simulator and first prototype taped out next month
- Straightforward to compile tensor algebra to the sparse abstract machine

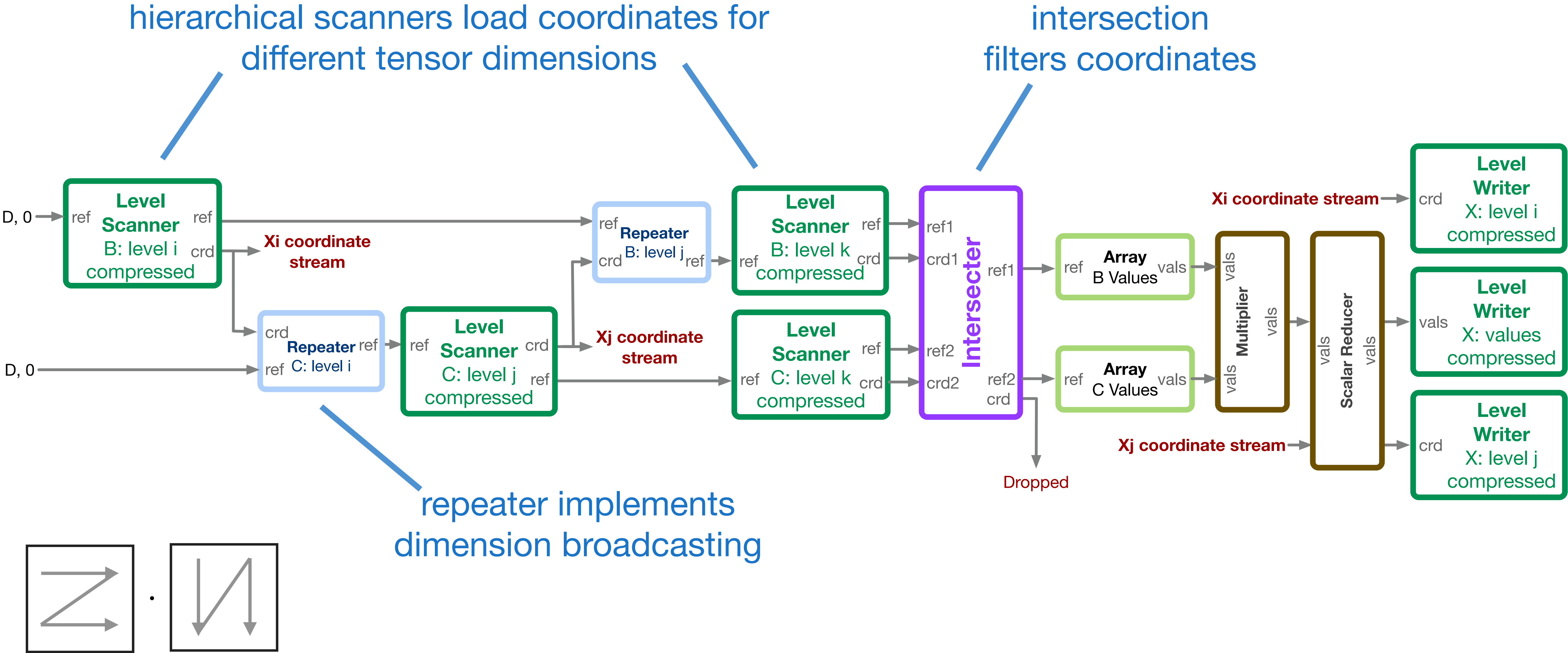


Olivia Hsu Maxwell Strange Jaeyeon Won Ritvik Sharma Kunle Olukotun Joel Emer Mark Horowitz Fred Kjolstad



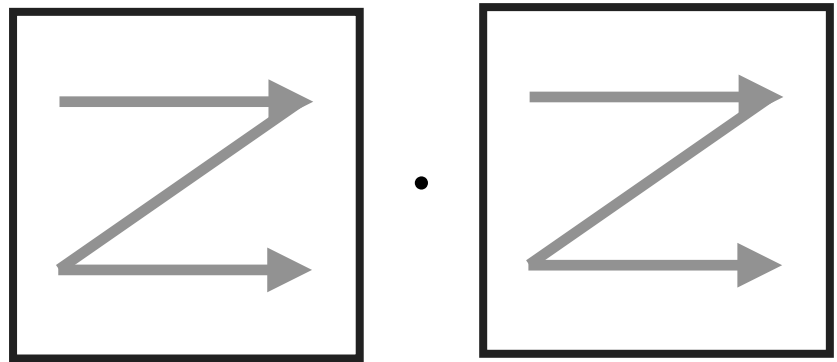
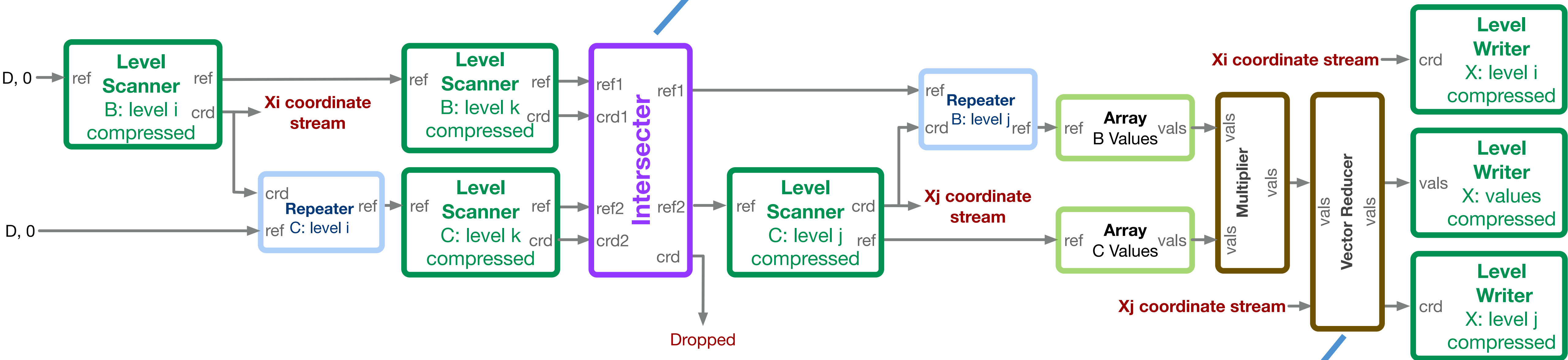
arxiv.org/abs/2208.14610

Inner-product sparse matrix multiplication



Gustafson sparse matrix multiplication

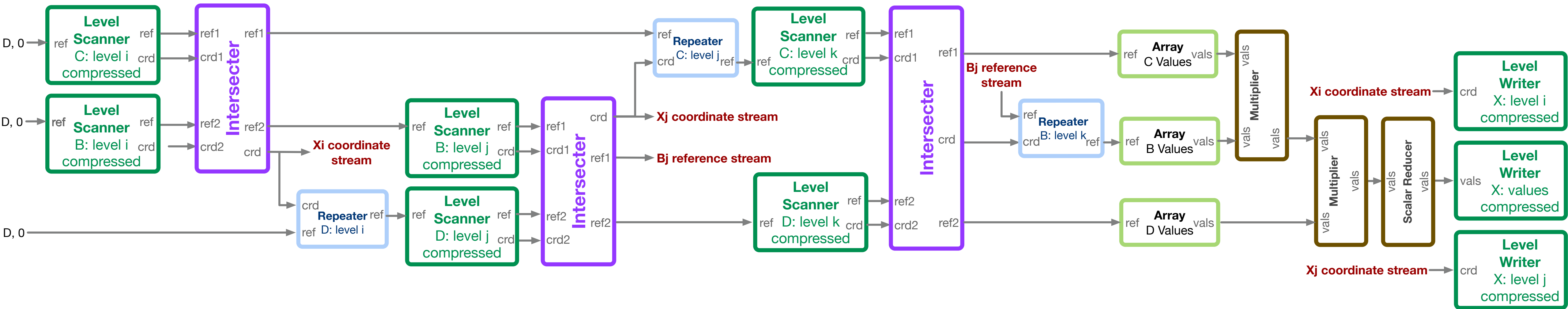
asymptotically less work,
because intersection occurs earlier



asymptotically more temporary memory,
because it must reduce whole rows

Fused SDDMM

$$A_{ij} = B \odot (CD)$$



$$O(\text{NNZ}_B \cdot K)$$

Conclusion and references

Unlike dense neural networks that can be reduced to GEMM, it will not be possible to reduce sparse neural networks to one optimized function

Compilation Approach

[Kjolstad et al. OOPSLA'17]

[Kjolstad et al. MIT'20]

Scheduling Language

[Kjolstad et al. CGO'19]

[Senanayake et al. OOPSLA'20]

Distributed Compilation

[Yadav et al. PLDI'22]

[Yadav et al. SC'22]

Sparse Abstract Machine

[Hsu et al. arXiv'22]

Format Abstractions

[Chou et al. OOPSLA'18]

[Chou et al. PLDI'20]

Autoscheduling

[Ahrens et al. PLDI'22]

User-Defined Functions

[Henry and Hsu et al. OOPSLA'21]

Verification

[Kovach and Kjolstad arXiv'22]