# *Tensaurus*: A Versatile Accelerator for Mixed Sparse-Dense Tensor Computations

**Nitish Srivastava**, Hanchen Jin, Shaden Smith[2], Hongbo Rong[3],
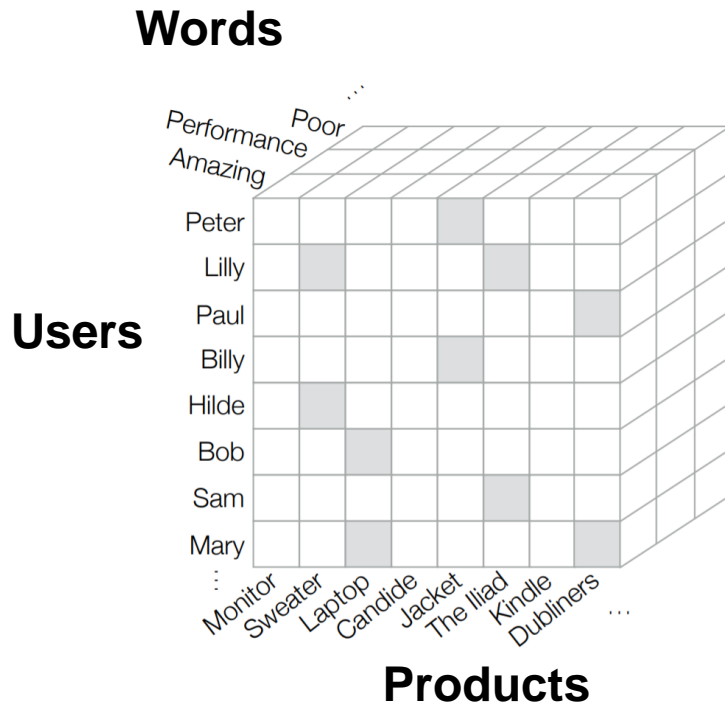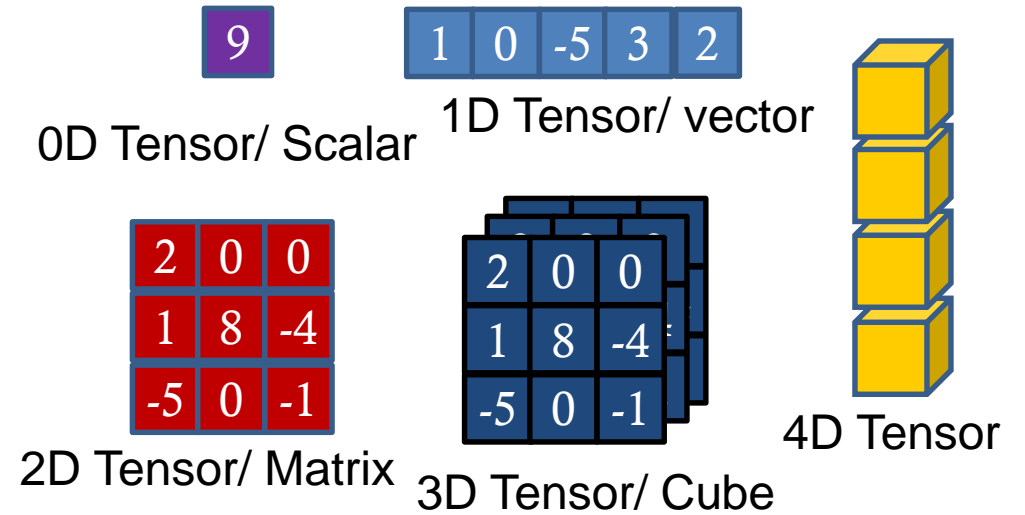David Albonesi, and Zhiru Zhang

Cornell University
[2]Microsoft AI & Research
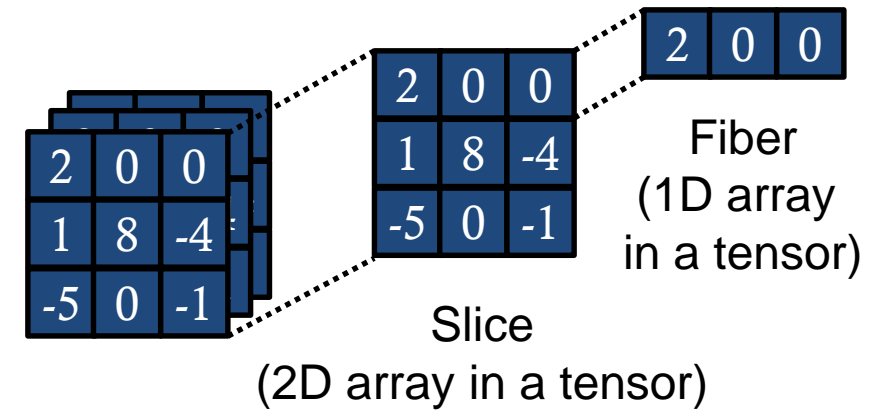[3]Intel Parallel Computing Lab

# What is a Tensor?

- **Tensors** are generalization of matrices to n dimensions
  - Scalar is tensor with 0 dimensions
  - Vector is tensor with 1 dimension
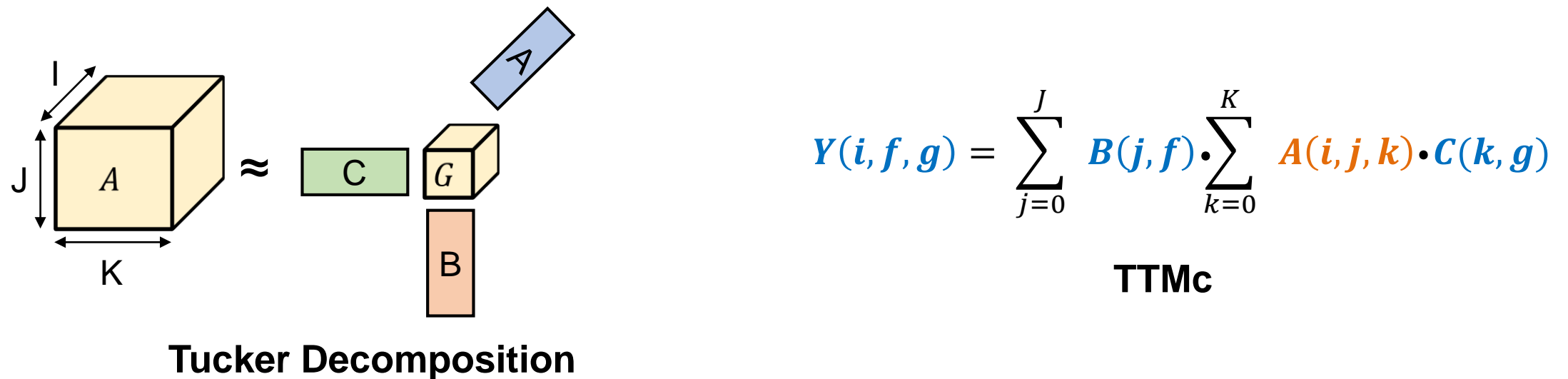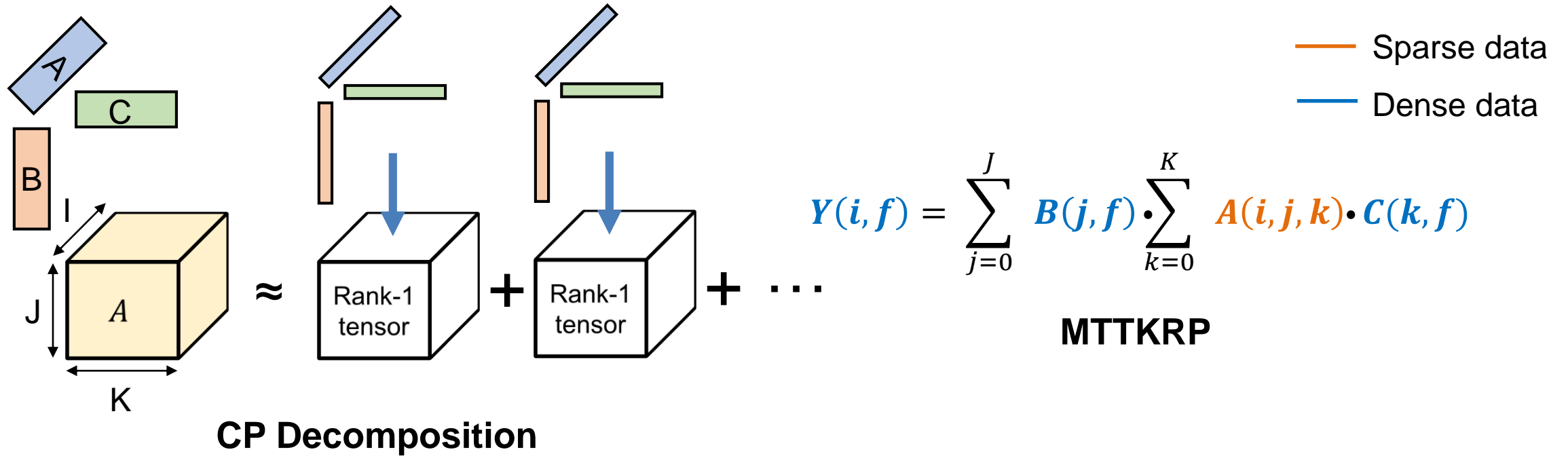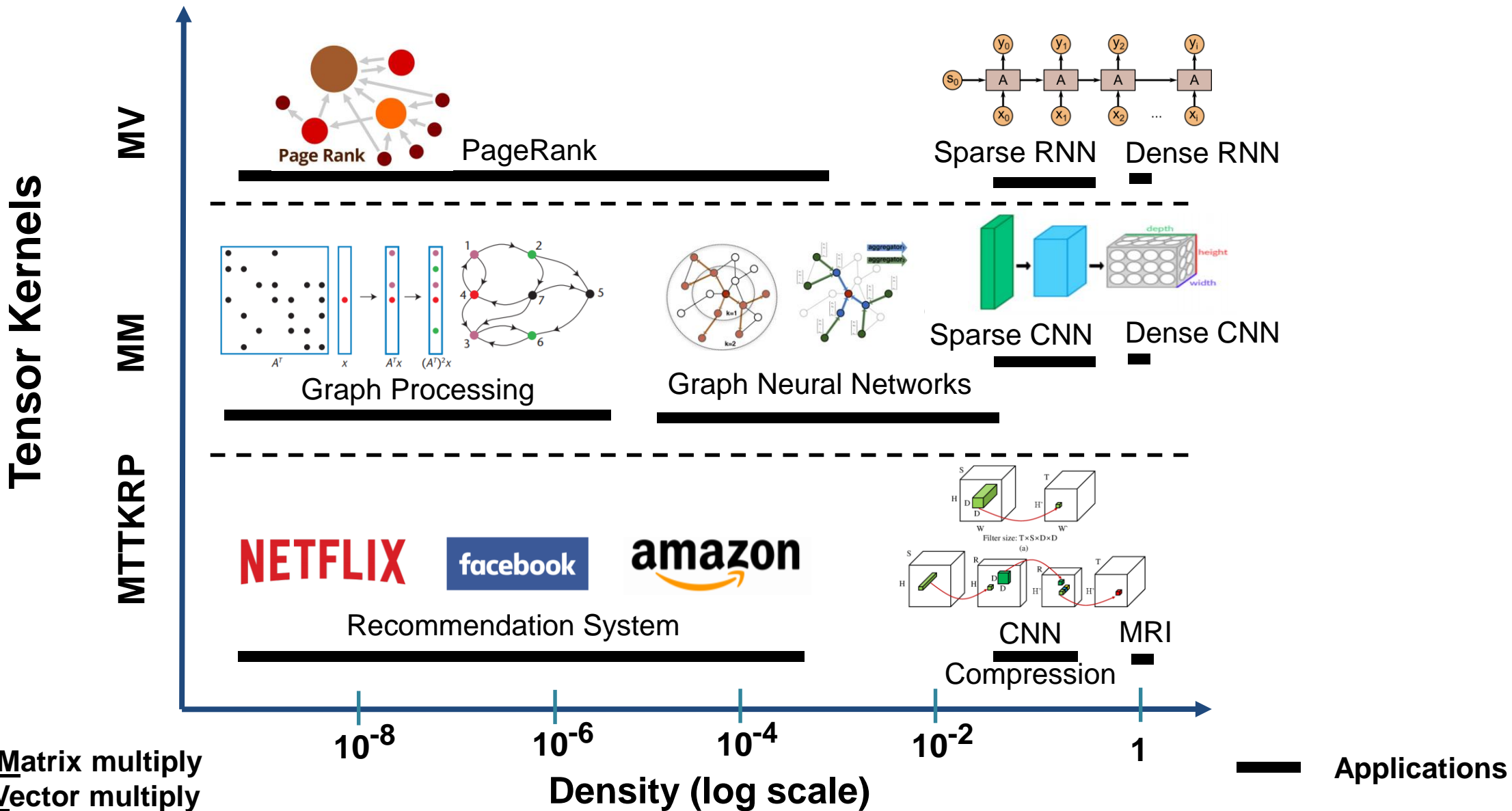  - Matrix is tensor with 2 dimensions, and so on

| 9 |
| --- |

0D Tensor/ Scalar

| 1 | 0 | -5 | 3 | 2 |
| --- | --- | --- | --- | --- |

1D Tensor/ vector

| 2 | 0 | 0 |
| --- | --- | --- |
| 1 | 8 | -4 |
| -5 | 0 | -1 |

2D Tensor/ Matrix

| 2 | 0 | 0 |
| --- | --- | --- |
| 1 | 8 | -4 |
| -5 | 0 | -1 |

3D Tensor/ Cube

4D Tensor

**Words**

**Users**

**Products**

Peter, Lilly, Paul, Billy, Hilde, Bob, Sam, Mary

Performance, Poor, Amazing

Monitor, Sweater, Laptop, Candide, Jacket, The Iliad, Kindle, Dubliners

## amazon
### Product Reviews

- High-dimensional data

- Density: $10^{-7}$ %

- Requires low-dimensional representation for ease of analysis

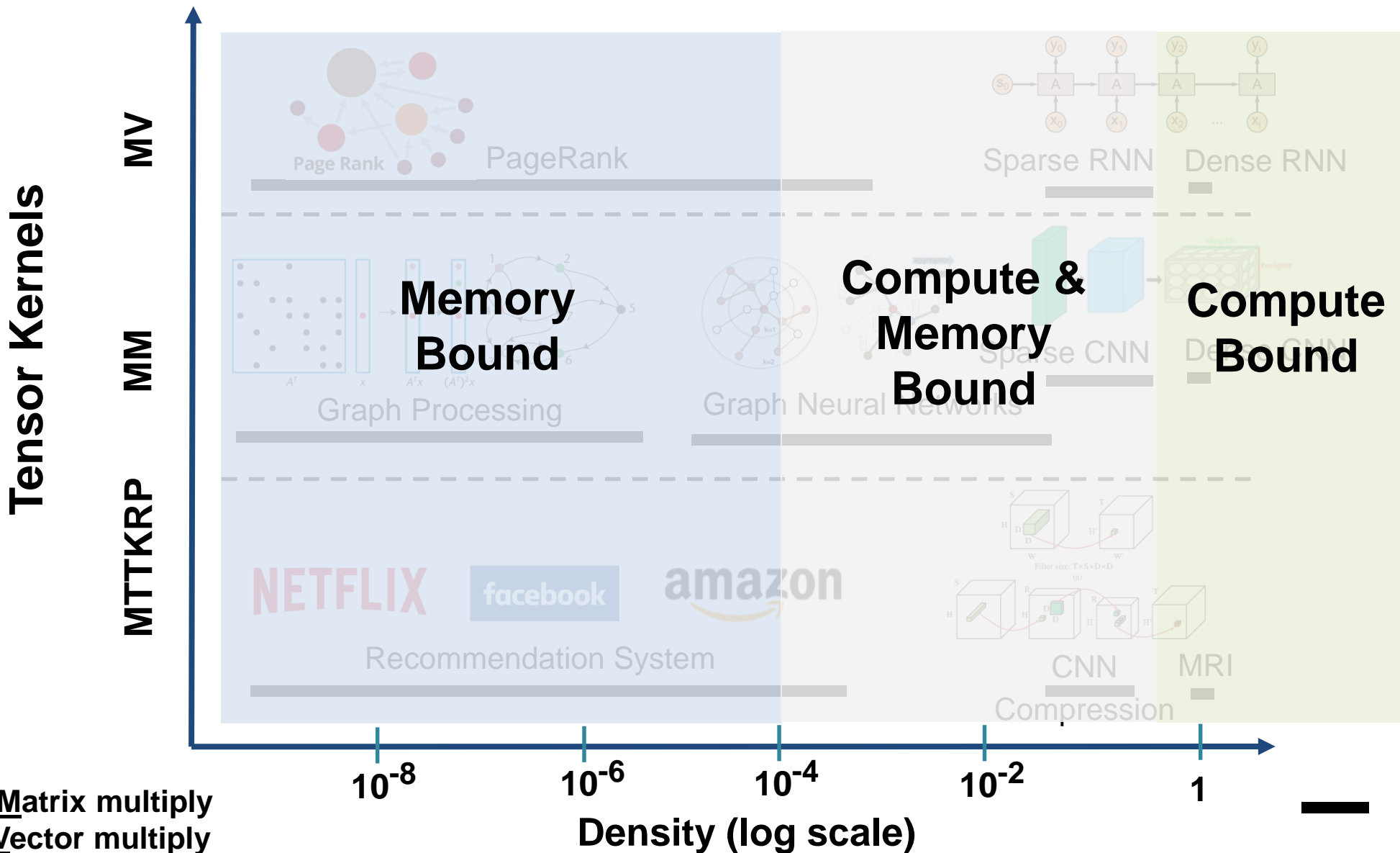| 2 | 0 | 0 |
| --- | --- | --- |
| 1 | 8 | -4 |
| -5 | 0 | -1 |

| 2 | 0 | 0 |
| --- | --- | --- |
| 1 | 8 | -4 |
| -5 | 0 | -1 |

Slice
(2D array in a tensor)

| 2 | 0 | 0 |
| --- | --- | --- |

Fiber
(1D array in a tensor)

# Tensor Decompositions for Low-Dimensional Representation



**CP Decomposition**

$$Y(i,f) = \sum_{j=0}^{J} B(j,f) \cdot \sum_{k=0}^{K} A(i,j,k) \cdot C(k,f)$$

**MTTKRP**

**Tucker Decomposition**

$$Y(i,f,g) = \sum_{j=0}^{J} B(j,f) \cdot \sum_{k=0}^{K} A(i,j,k) \cdot C(k,g)$$
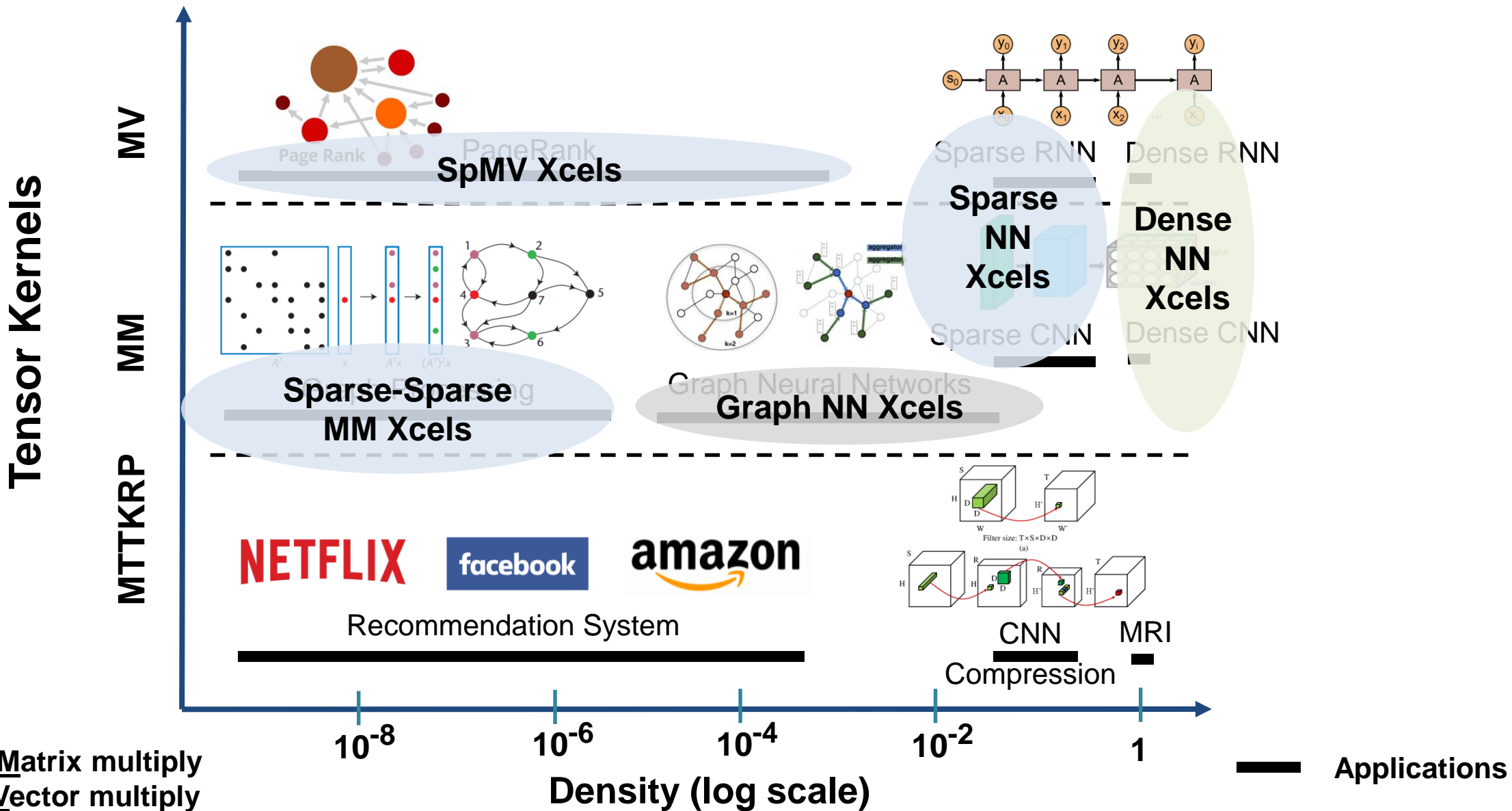
**TTMc**

Sparse data
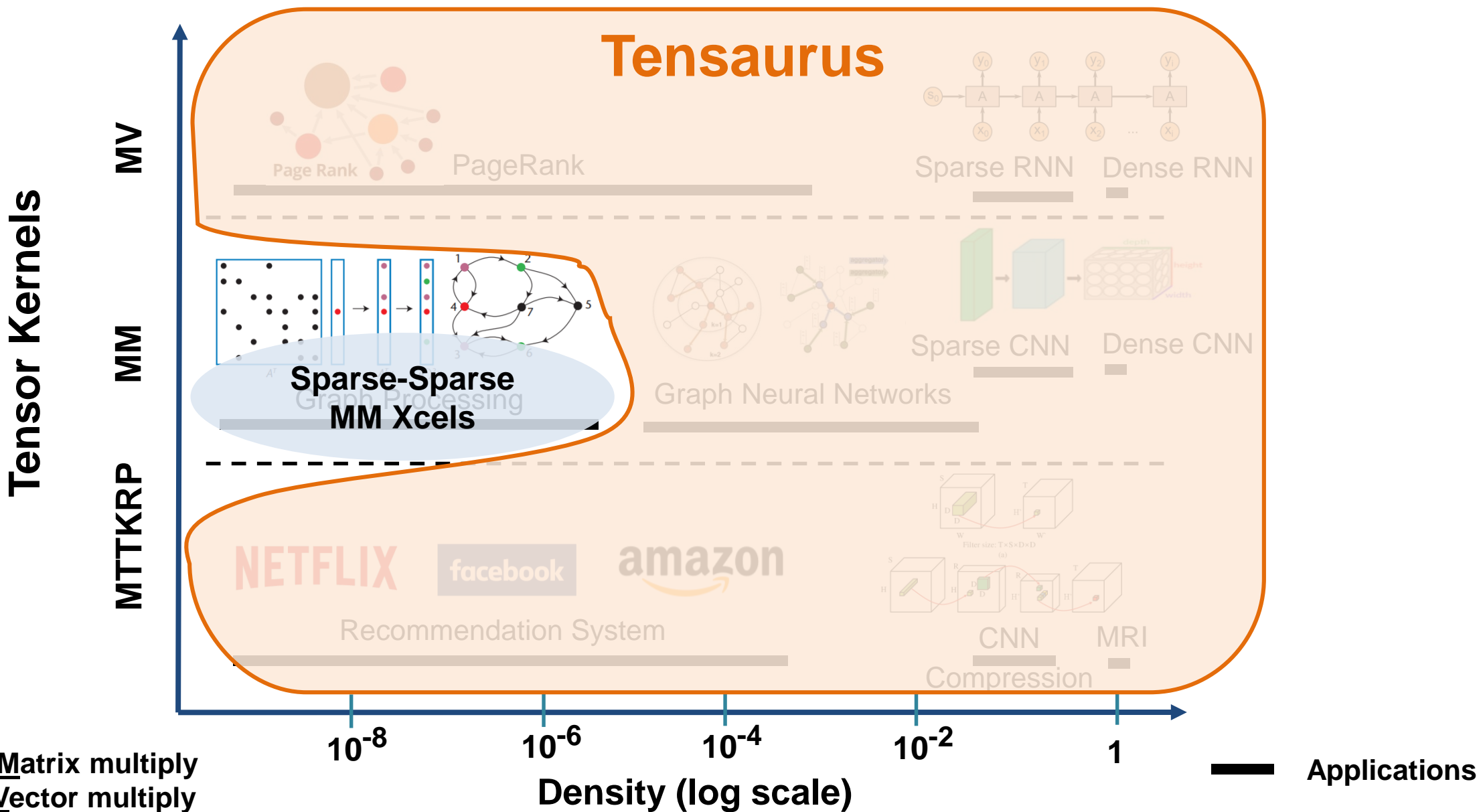Dense data

# Kernel-Sparsity Spectrum of Tensor Applications



MM = **M**atrix-**M**atrix multiply
MV = **M**atrix-**V**ector multiply

3

# Kernel-Sparsity Spectrum of Tensor Applications



**Tensor Kernels**

MV · MM · MTTKRP

PageRank · Page Rank

**Memory Bound**

Graph Processing

**Compute & Memory Bound**

Graph Neural Networks

Sparse CNN

**Compute Bound**

Sparse RNN · Dense RNN

Dense CNN

NETFLIX · facebook · amazon

Recommendation System

CNN · MRI
Compression

Density (log scale): $10^{-8}$ $10^{-6}$ $10^{-4}$ $10^{-2}$ 1

MM = **M**atrix-**M**atrix multiply
MV = **M**atrix-**V**ector multiply

━━ **Applications**

# Kernel-Sparsity Spectrum of Tensor Applications



**Tensor Kernels**

MV

MM

MTTKRP

SpMV Xcels

Sparse RNN    Dense RNN

Sparse NN Xcels

Dense NN Xcels

Sparse-Sparse MM Xcels

Graph NN Xcels

Sparse CNN    Dense CNN

NETFLIX    facebook    amazon

Recommendation System

CNN    MRI

Compression

$10^{-8}$    $10^{-6}$    $10^{-4}$    $10^{-2}$    1

**Density (log scale)**

MM = **M**atrix-**M**atrix multiply
MV = **M**atrix-**V**ector multiply

▬▬ **Applications**

5

# Kernel-Sparsity Spectrum of Tensor Applications



**Tensor Kernels** (y-axis): MV, MM, MTTKRP

**Density (log scale)** (x-axis): $10^{-8}$, $10^{-6}$, $10^{-4}$, $10^{-2}$, $1$

PageRank

Sparse RNN    Dense RNN

Sparse-Sparse MM

Graph Processing

Sparse-Dense and Dense-Dense Tensor Kernels

Graph Neural Networks

Sparse CNN    Dense CNN

Recommendation System

CNN    MRI

Compression

MM = **M**atrix-**M**atrix multiply
MV = **M**atrix-**V**ector multiply

━━ Applications

6

# Kernel-Sparsity Spectrum of Tensor Applications



MM = **M**atrix-**M**atrix multiply
MV = **M**atrix-**V**ector multiply

# Challenges with Sparse-Dense Tensor Acceleration

▸ **Dense Tensor Acceleration**

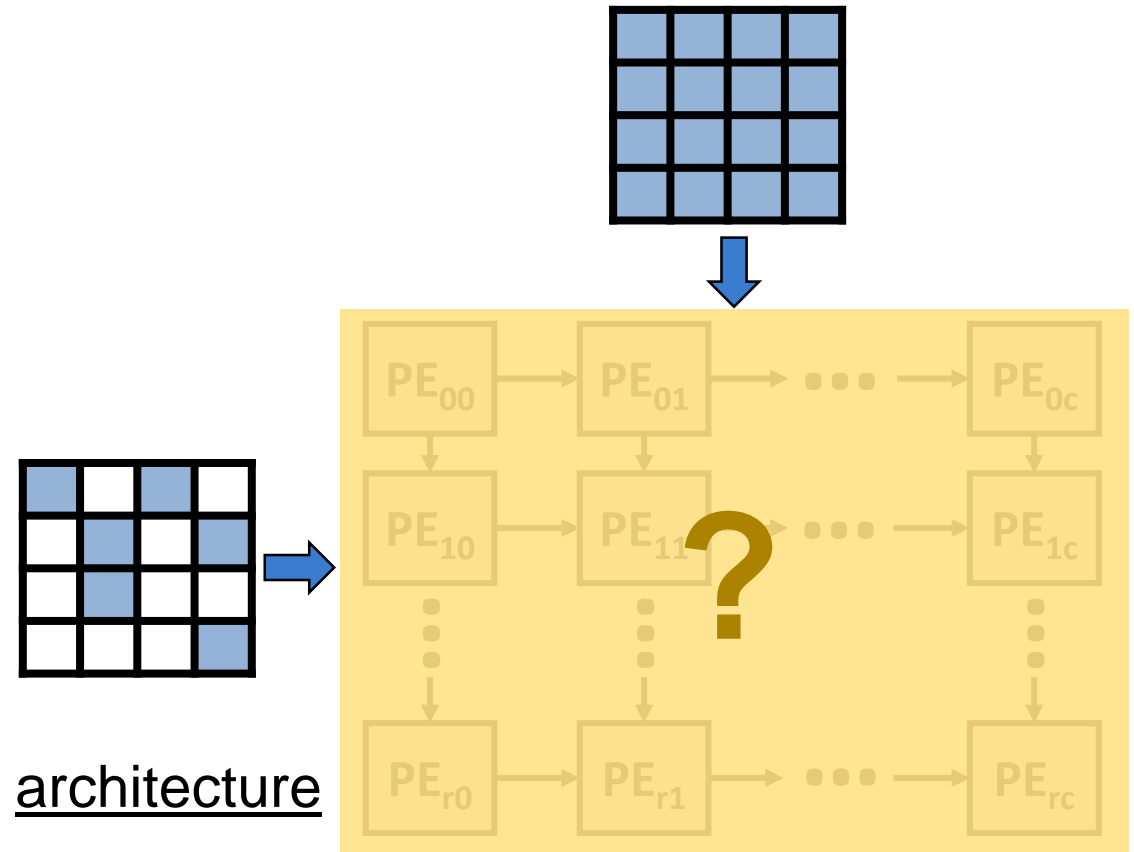    – Systolic arrays provide high utilization of both memory and compute



**Highly efficient**

# Challenges with Sparse-Dense Tensor Acceleration

- **Dense Tensor Acceleration**
  - Systolic arrays provide high utilization of both memory and compute

- **Mixed Sparse-Dense Tensor Acceleration**
  - Memory bound
  - Hard to achieve high compute and bandwidth utilization

- **Goal:** Leverage a dense accelerator to efficiently perform sparse-dense compute

- **Key approach:** Co-design of accelerator <u>architecture</u> and <u>sparse format</u>
  - Low overhead of supporting sparse compute
  - High compute and bandwidth utilization



**How to push sparse data in dense systolic array?**
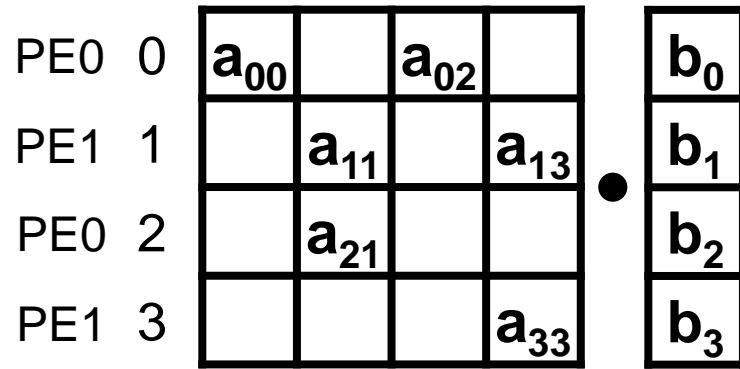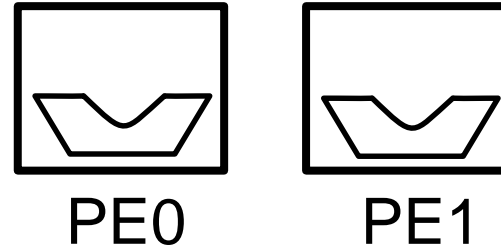
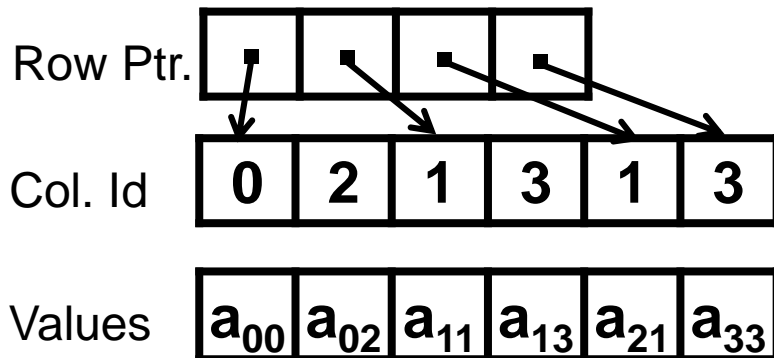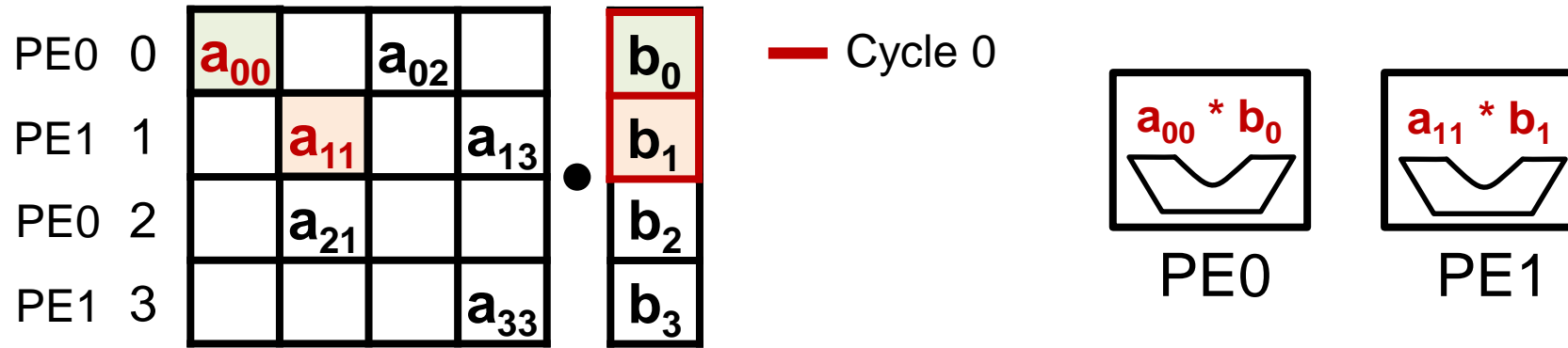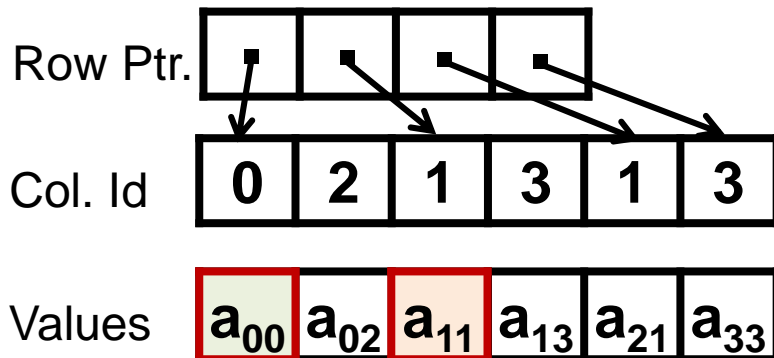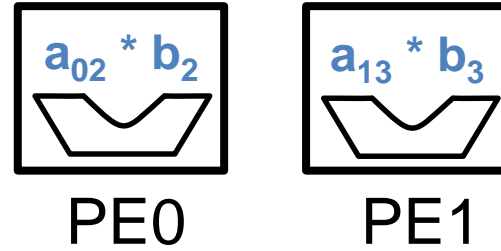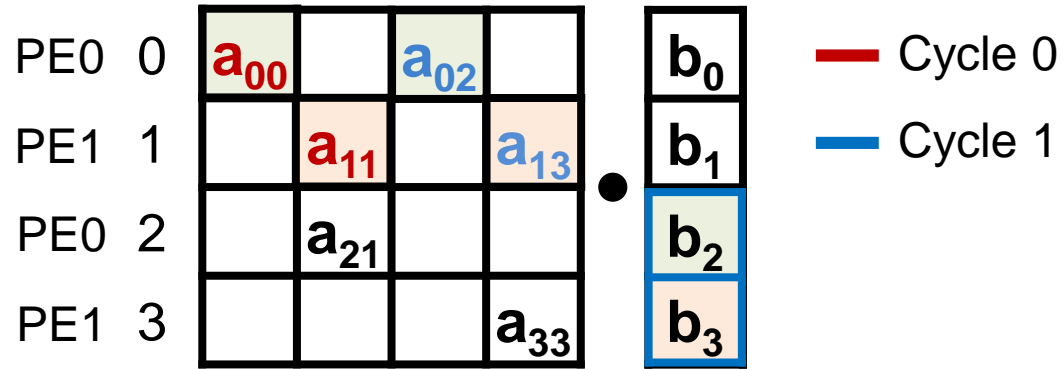# Importance of Accelerator Friendly Formats



SpMV
**(Sparse Matrix dense Vector multiply)**

**Compressed Sparse Row Format
(CSR)**

# Importance of Accelerator Friendly Formats



**SpMV**
**(Sparse Matrix dense Vector multiply)**

**Compressed Sparse Row Format**
**(CSR)**

Row Ptr.

Col. Id

| 0 | 2 | 1 | 3 | 1 | 3 |

Values

| $a_{00}$ | $a_{02}$ | $a_{11}$ | $a_{13}$ | $a_{21}$ | $a_{33}$ |

# Importance of Accelerator Friendly Formats

PE0 0 | $a_{00}$ | | $a_{02}$ | |

PE1 1 | | $a_{11}$ | | $a_{13}$ |

PE0 2 | | $a_{21}$ | | |

PE1 3 | | | | $a_{33}$ |

$b_0$
$b_1$
$b_2$
$b_3$

— Cycle 0

**SpMV**
**(Sparse Matrix dense Vector multiply)**

PE0: $a_{00} * b_0$
PE1: $a_{11} * b_1$

**Compressed Sparse Row Format**
**(CSR)**

Row Ptr.

Col. Id | 0 | 2 | 1 | 3 | 1 | 3 |

Values | $a_{00}$ | $a_{02}$ | $a_{11}$ | $a_{13}$ | $a_{21}$ | $a_{33}$ |

# Importance of Accelerator Friendly Formats

PE0  0   $a_{00}$     $a_{02}$            $b_0$         ▬ Cycle 0
PE1  1          $a_{11}$     $a_{13}$     $b_1$         ▬ Cycle 1
PE0  2          $a_{21}$                  $b_2$
PE1  3                       $a_{33}$     $b_3$

**SpMV**
**(Sparse Matrix dense Vector multiply)**

PE0: $a_{02} * b_2$

PE1: $a_{13} * b_3$

**Compressed Sparse Row Format**
**(CSR)**

Row Ptr.

Col. Id | 0 | 2 | 1 | 3 | 1 | 3 |

Values | $a_{00}$ | $a_{02}$ | $a_{11}$ | $a_{13}$ | $a_{21}$ | $a_{33}$ |

# Importance of Accelerator Friendly Formats

PE0  0  | $a_{00}$ |  | $a_{02}$ |  |
PE1  1  |  | $a_{11}$ |  | $a_{13}$ |
PE0  2  |  | $a_{21}$ |  |  |
PE1  3  |  |  |  | $a_{33}$ |

$\cdot$

$b_0$
$b_1$
$b_2$
$b_3$

— Cycle 0
— Cycle 1
— Cycle 2

**SpMV**
**(Sparse Matrix dense Vector multiply)**

$a_{21} * b_1$

PE0

$a_{33} * b_3$

PE1

**Compressed Sparse Row Format**
**(CSR)**

Row Ptr.

Col. Id | 0 | 2 | 1 | 3 | 1 | 3 |

Values | $a_{00}$ | $a_{02}$ | $a_{11}$ | $a_{13}$ | $a_{21}$ | $a_{33}$ |

# Importance of Accelerator Friendly Formats

PE0 0 | $a_{00}$ | | $a_{02}$ | |
PE1 1 | | $a_{11}$ | | $a_{13}$ |
PE0 2 | | $a_{21}$ | | |
PE1 3 | | | | $a_{33}$ |

$b_0$
$b_1$
$b_2$
$b_3$

— Cycle 0
— Cycle 1
— Cycle 2

PE0: $a_{21} * b_1$

PE1: $a_{33} * b_3$

**SpMV**
**(Sparse Matrix dense Vector multiply)**

**Compressed Sparse Row Format (CSR)**

Row Ptr.

Col. Id | 0 | 2 | 1 | 3 | 1 | 3 |

Values | $a_{00}$ | $a_{02}$ | $a_{11}$ | $a_{13}$ | $a_{21}$ | $a_{33}$ |

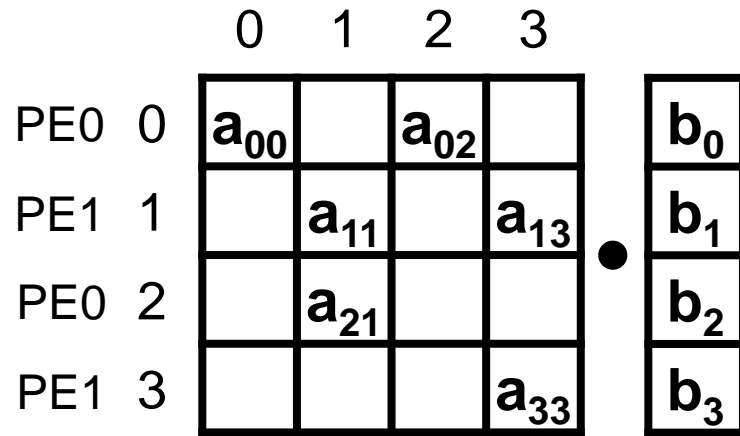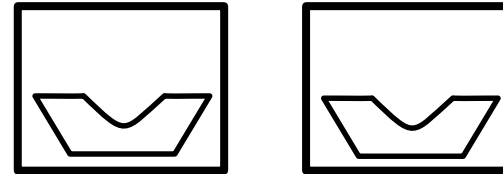**Problems with CSR**

- **Non-streaming & non-vectorized accesses**

- **Indirect memory access**
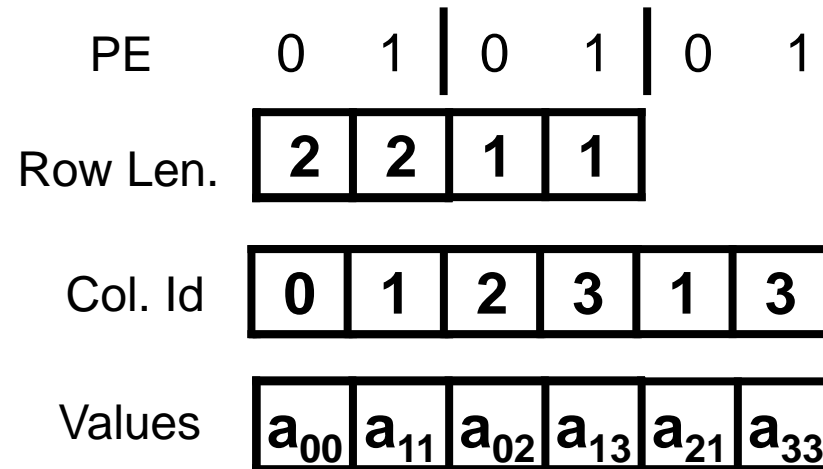
# Importance of Accelerator Friendly Formats



**SpMV**
**(Sparse Matrix dense Vector multiply)**

**Compressed Interleaved Sparse Row (CISR)** [1]

| PE | 0 | 1 | 0 | 1 | 0 | 1 |
|---|---|---|---|---|---|---|
| Row Len. | 2 | 2 | 1 | 1 | | |
| Col. Id | 0 | 1 | 2 | 3 | 1 | 3 |
| Values | $a_{00}$ | $a_{11}$ | $a_{02}$ | $a_{13}$ | $a_{21}$ | $a_{33}$ |

[1] Fowers, et al. A high memory bandwidth FPGA accelerator for sparse matrix-vector multiplication, Int'l Symp. On *Field-Programmable Custom Computing Machines (FCCM), 2014*

# Importance of Accelerator Friendly Formats



**SpMV**
**(<u>S</u>parse <u>M</u>atrix dense <u>V</u>ector multiply)**

**Compressed Interleaved Sparse Row (CISR)** [1]

# Importance of Accelerator Friendly Formats



SpMV
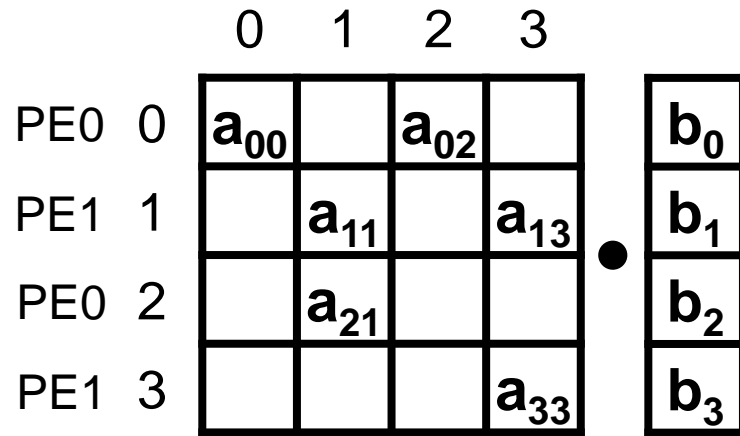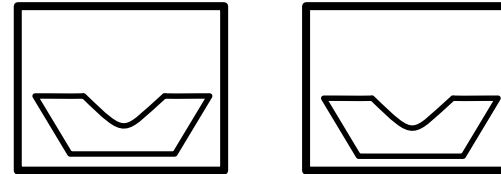**(Sparse Matrix dense Vector multiply)**

Compressed Interleaved Sparse Row
**(CISR)** [1]

[1] Fowers, et al. A high memory bandwidth FPGA accelerator for sparse matrix-vector multiplication, Int'l Symp. On *Field-Programmable Custom Computing Machines (FCCM), 2014*

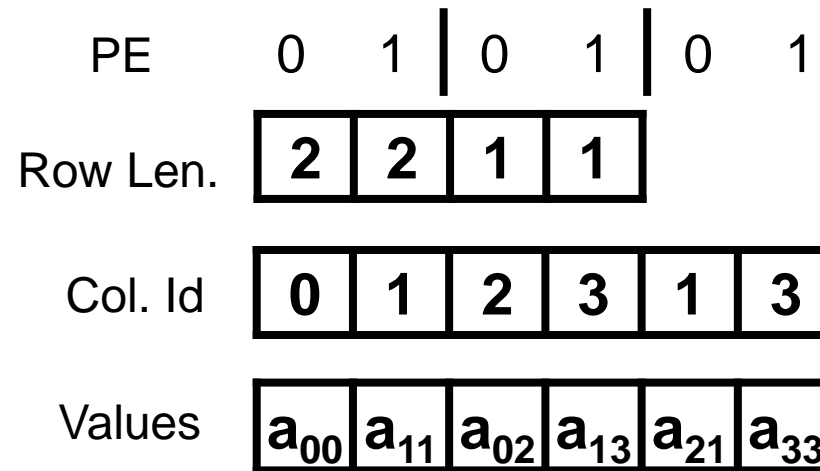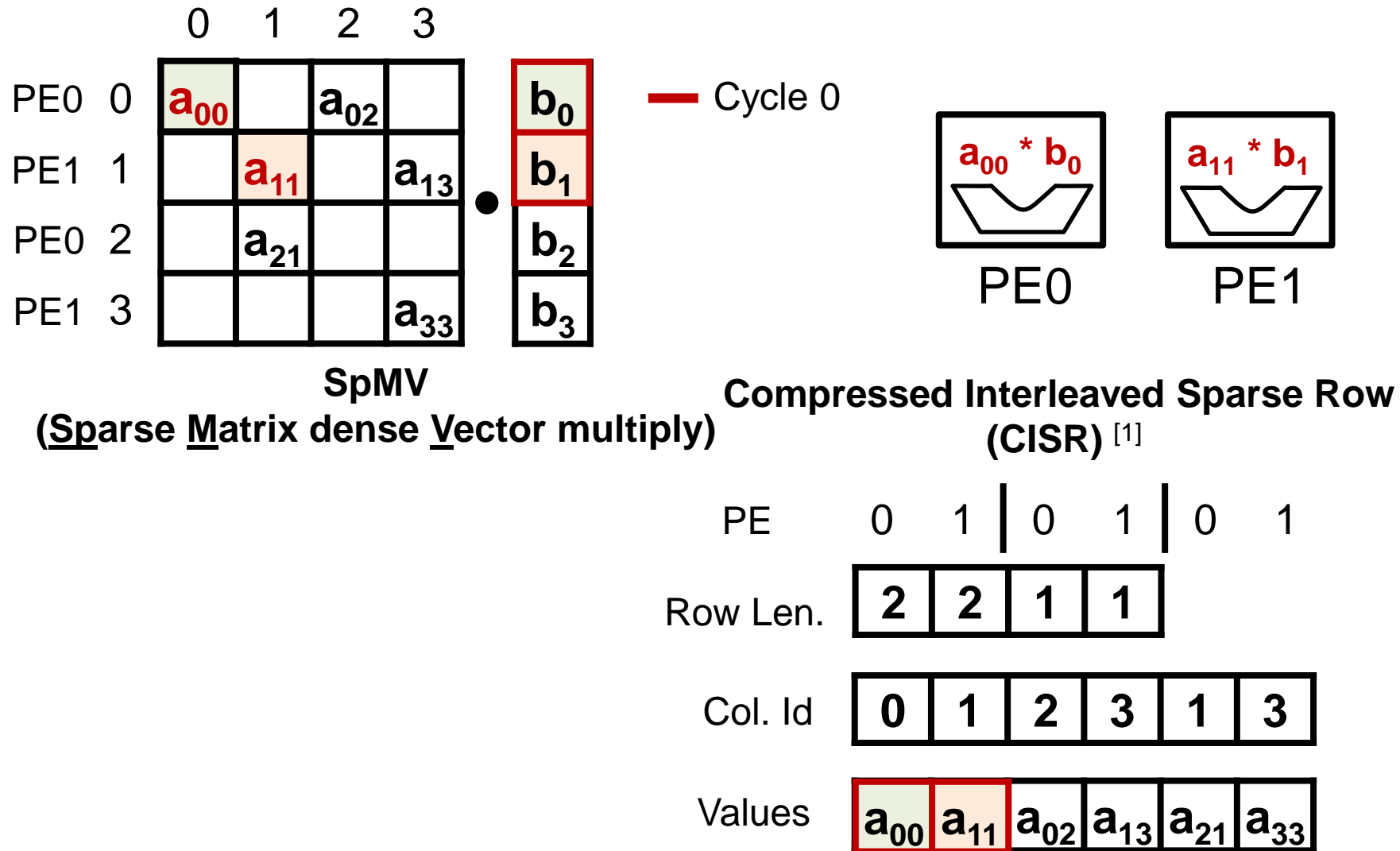# Importance of Accelerator Friendly Formats



SpMV
**(<u>S</u>parse <u>M</u>atrix dense <u>V</u>ector multiply)**

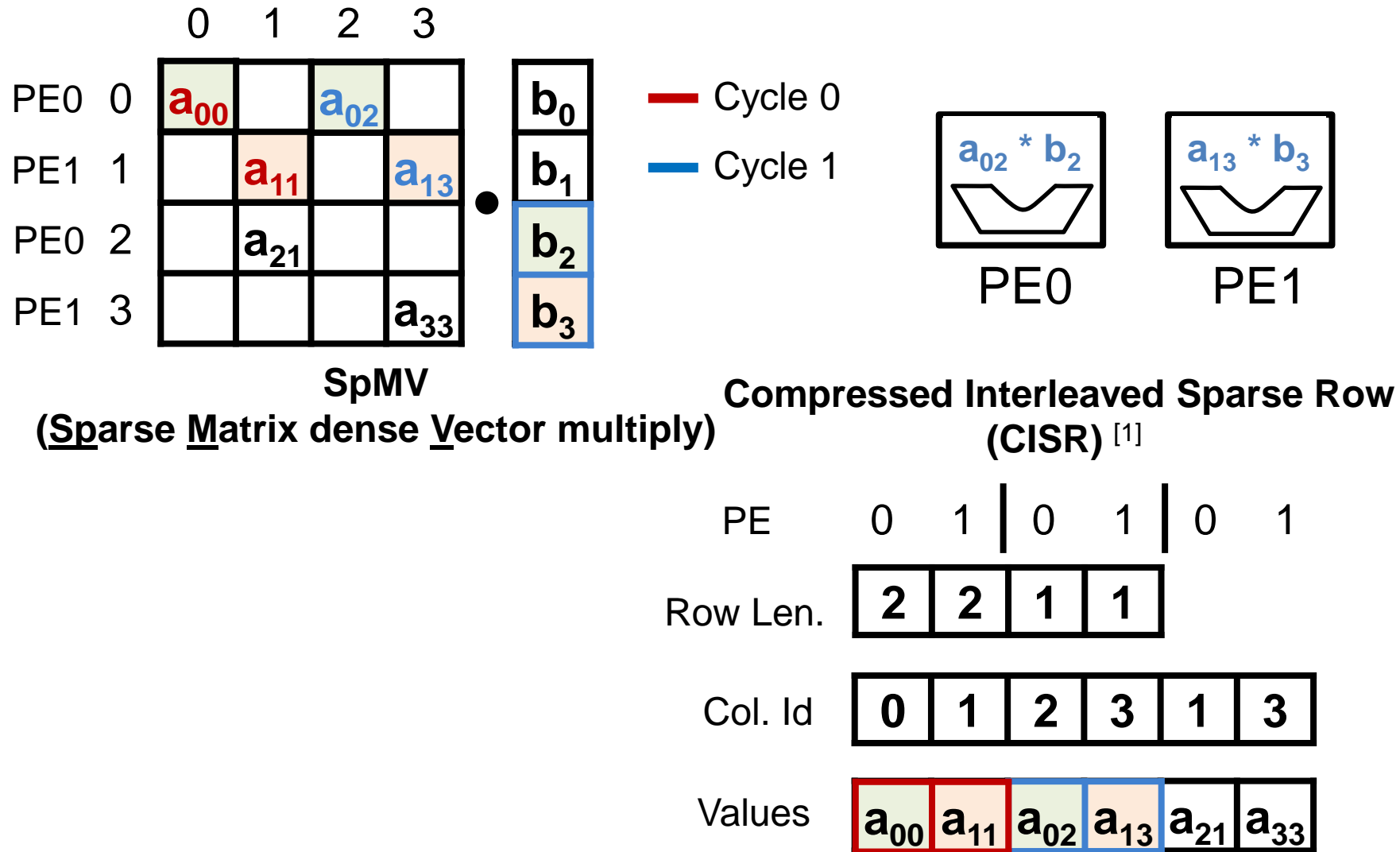**Compressed Interleaved Sparse Row (CISR)** [1]

[1] Fowers, et al. A high memory bandwidth FPGA accelerator for sparse matrix-vector multiplication, Int'l Symp. On *Field-Programmable Custom Computing Machines (FCCM), 2014*

# Importance of Accelerator Friendly Formats



SpMV
**(Sparse Matrix dense Vector multiply)**

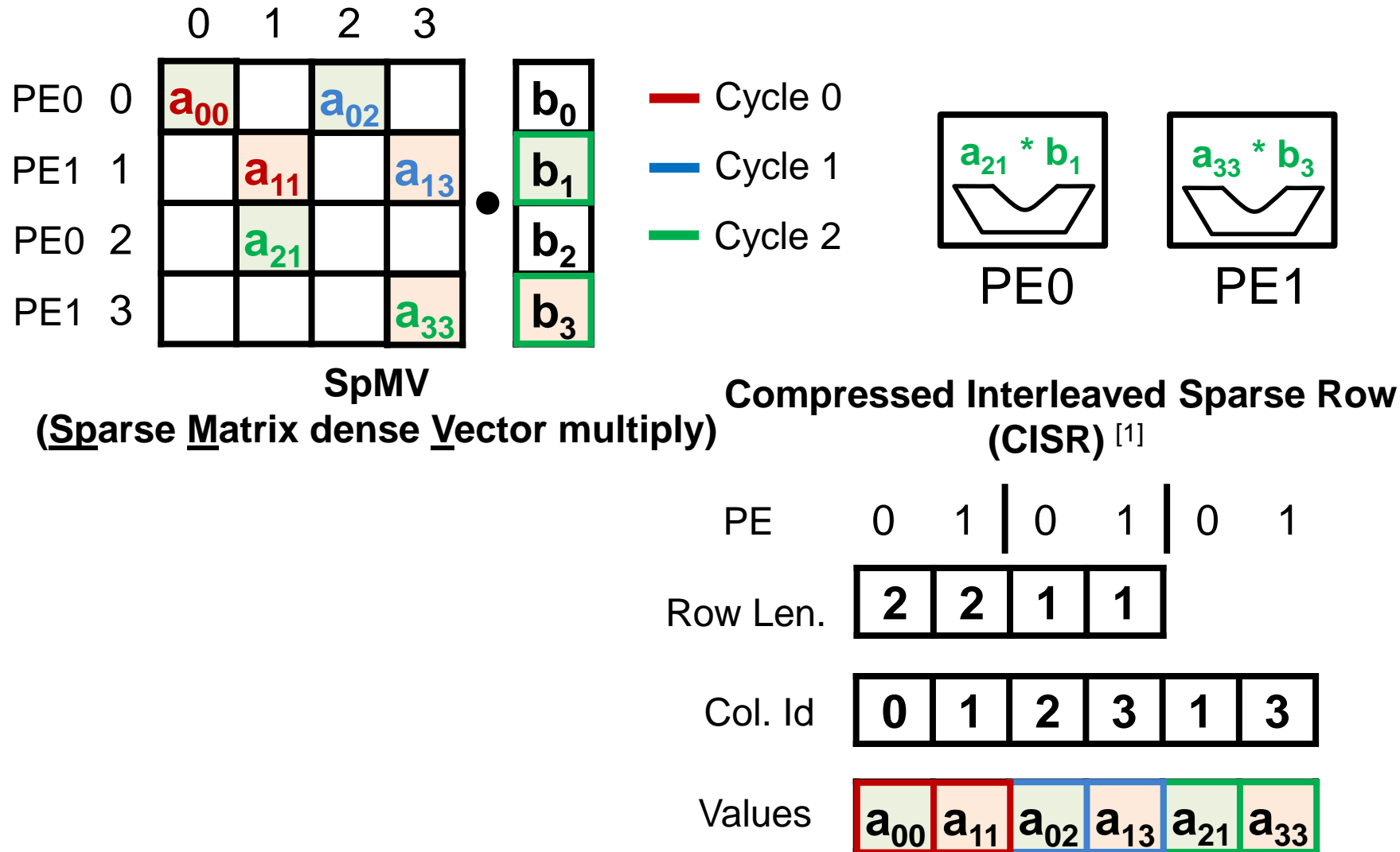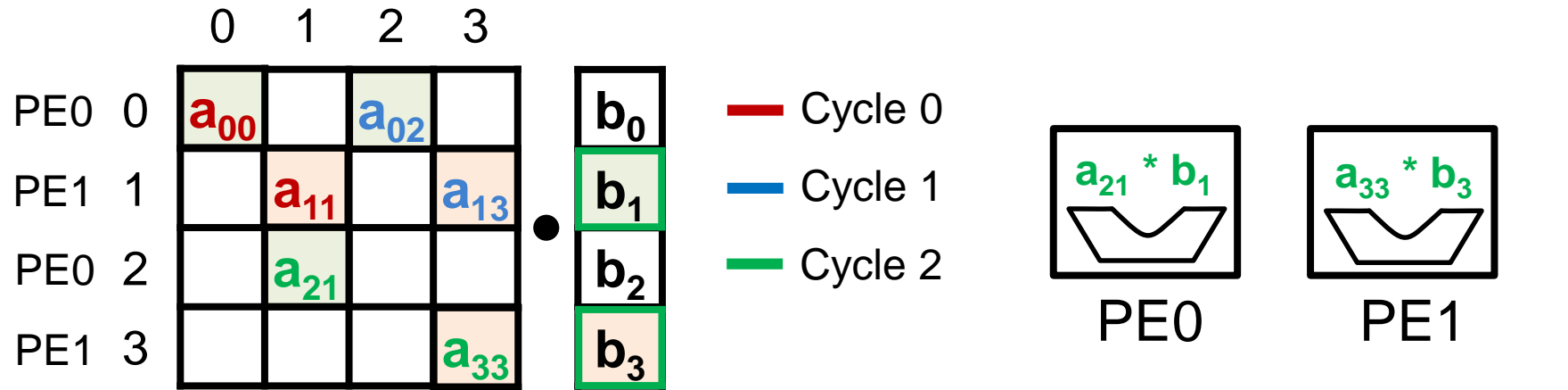Compressed Interleaved Sparse Row **(CISR)** [1]

[1] Fowers, et al. A high memory bandwidth FPGA accelerator for sparse matrix-vector multiplication, Int'l Symp. On *Field-Programmable Custom Computing Machines (FCCM), 2014*

# Importance of Accelerator Friendly Formats

Max: **16GB/s** (DDR3)

PE0 0   $a_{00}$   $a_{02}$

PE1 1   $a_{11}$   $a_{13}$

PE0 2   $a_{21}$

PE1 3   $a_{33}$

$b_0$   $b_1$   $b_2$   $b_3$

— Cycle 0
— Cycle 1
— Cycle 2

$a_{21} * b_1$    $a_{33} * b_3$

PE0    PE1

**SpMV**
**(Sparse Matrix dense Vector multiply)**

**Compressed Interleaved Sparse Row (CISR)** [1]

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

Row Len.

| 2 | 2 | 1 | 1 |
|---|---|---|---|

Col. Id

| 0 | 1 | 2 | 3 | 1 | 3 |
|---|---|---|---|---|---|

Values

| $a_{00}$ | $a_{11}$ | $a_{02}$ | $a_{13}$ | $a_{21}$ | $a_{33}$ |
|---|---|---|---|---|---|

**1.8GB/s** : 8 PEs

**CSR**

**1.6GB/s** : 2 PEs

**Utilized Bandwidth vs. # PEs**

[1] Fowers, et al. A high memory bandwidth FPGA accelerator for sparse matrix-vector multiplication, Int'l Symp. On *Field-Programmable Custom Computing Machines (FCCM)*, 2014
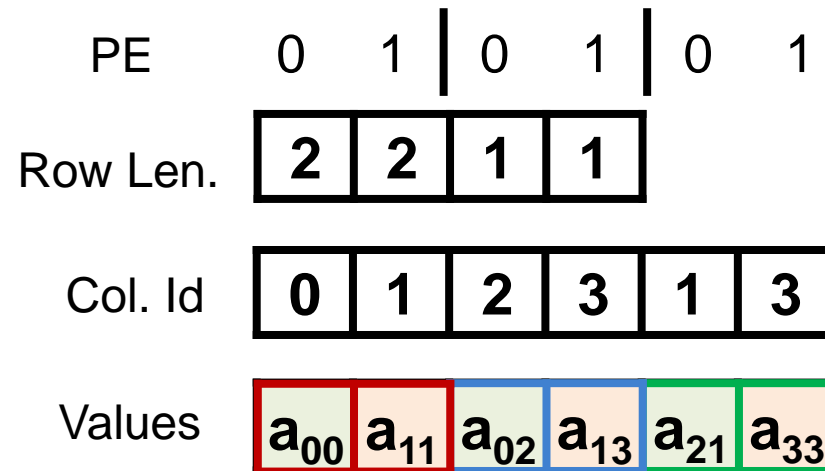
17

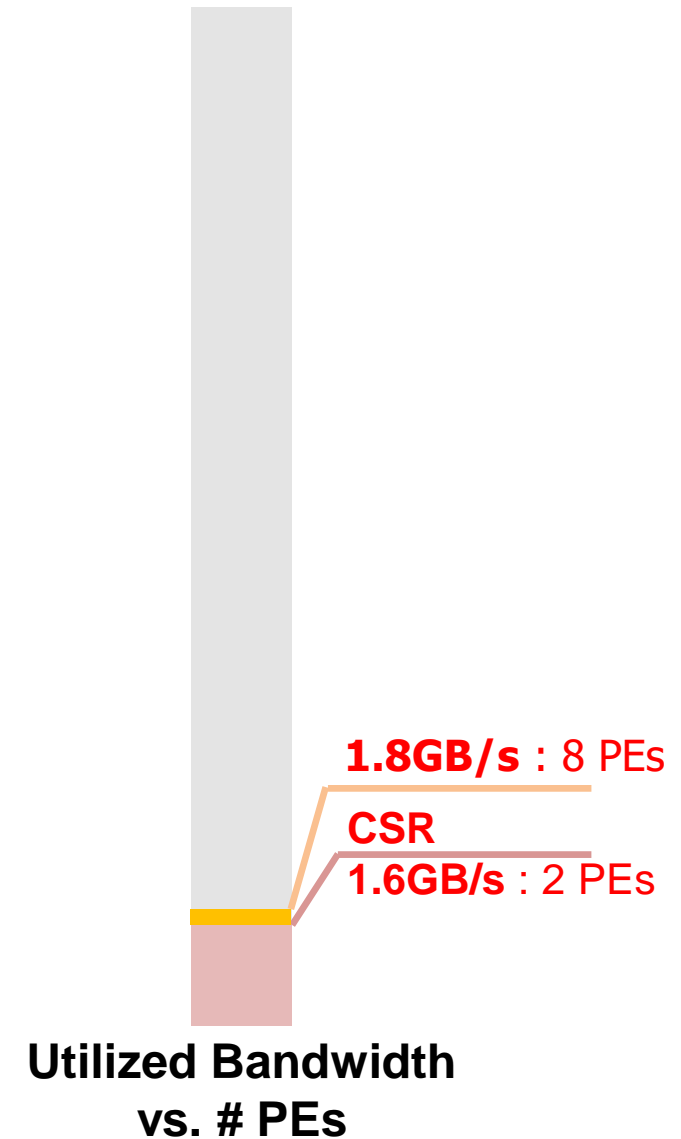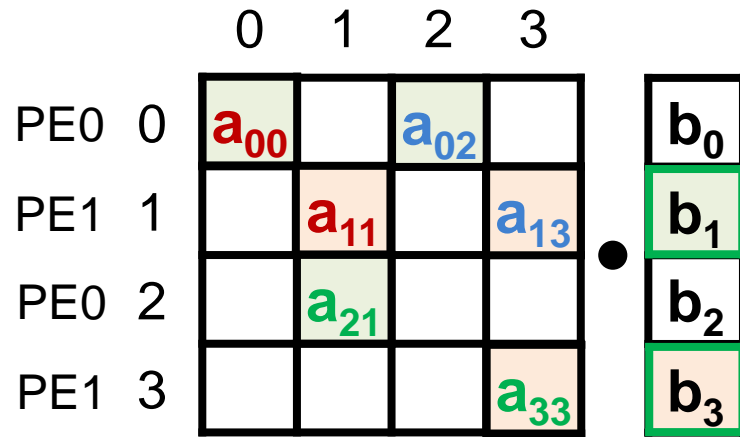# Importance of Accelerator Friendly Formats



**Max: 16GB/s (DDR3)**

**CISR:**
**11.2GB/s : 8 PEs**

**CISR:**
**6.1GB/s : 4 PEs**

**1.8GB/s : 8 PEs**

**CSR**
**1.6GB/s : 2 PEs**

6x higher bandwidth utilization

**Utilized Bandwidth vs. # PEs**

| | | | | Cycle 0 |
| | | | | Cycle 1 |
| | | | | Cycle 2 |

**SpMV**
**(Sparse Matrix dense Vector multiply)**

$a_{21} * b_1$  PE0

$a_{33} * b_3$  PE1

**Compressed Interleaved Sparse Row (CISR)** [1]

| PE | 0 | 1 | 0 | 1 | 0 | 1 |

Row Len. | 2 | 2 | 1 | 1 |

Col. Id | 0 | 1 | 2 | 3 | 1 | 3 |

Values | $a_{00}$ | $a_{11}$ | $a_{02}$ | $a_{13}$ | $a_{21}$ | $a_{33}$ |

[1] Fowers, et al. A high memory bandwidth FPGA accelerator for sparse matrix-vector multiplication, Int'l Symp. On *Field-Programmable Custom Computing Machines (FCCM), 2014*
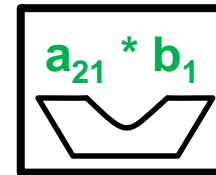
# Importance of Accelerator Friendly Formats

## SpMV


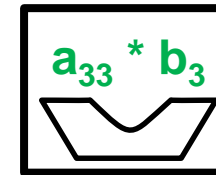
**SpMV**
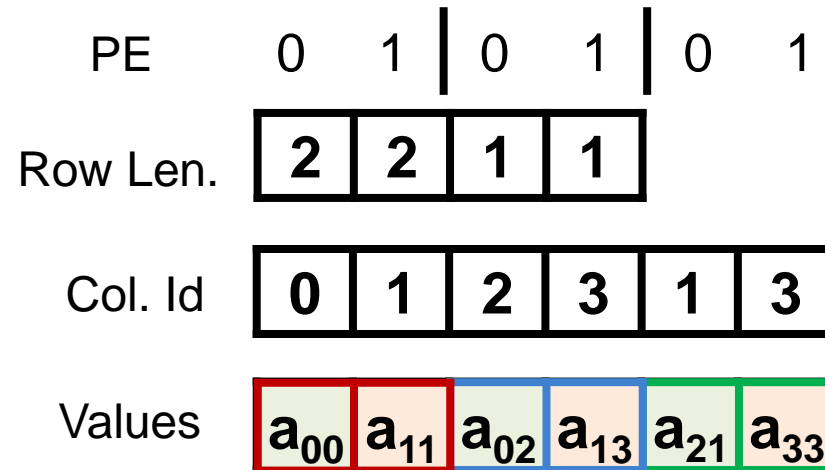**(Sparse Matrix dense Vector multiply)**

| Cycle 0 | (dark red) |
| Cycle 1 | (blue) |
| Cycle 2 | (green) |

$a_{21} * b_1$ — PE0

$a_{33} * b_3$ — PE1

## Compressed Interleaved Sparse Row (CISR) [1]

| PE | 0 | 1 | 0 | 1 | 0 | 1 |

Row Len. | 2 | 2 | 1 | 1 |

Col. Id | 0 | 1 | 2 | 3 | 1 | 3 |

Values | $a_{00}$ | $a_{11}$ | $a_{02}$ | $a_{13}$ | $a_{21}$ | $a_{33}$ |

### Benefits of CISR

- **Streaming & vectorized accesses**

## Utilized Bandwidth vs. # PEs

Max: **16GB/s** (DDR3)

6x higher bandwidth utilization

**CISR: 11.2GB/s** : 8 PEs

**CISR: 6.1GB/s** : 4 PEs

**1.8GB/s** : 8 PEs

**CSR 1.6GB/s** : 2 PEs

[1] Fowers, et al. A high memory bandwidth FPGA accelerator for sparse matrix-vector multiplication, Int'l Symp. On *Field-Programmable Custom Computing Machines (FCCM)*, 2014

# Extending CISR to CISS for Tensors



**CISR+**

Extending CISR+ to Tensors

slice: two-dimensional arrays in a tensor

| PE | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|
| **Row Id**/Col. Id | **0** | **1** | 0 | 1 | 2 | 3 | **2** | **3** | 1 | 3 |
| Values | **0** | **0** | $a_{00}$ | $a_{11}$ | $a_{02}$ | $a_{13}$ | **0** | **0** | $a_{21}$ | $a_{33}$ |

**CISR+ avoids row-decoding**

**Adds extra zeros to make design more scalable ⇒ HW/SW co-design**

**Compressed Interleaved Sparse Slice (CISS)**

| PE | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|
| $k$ | x | x | 0 | 1 | 1 | x | x | 0 | 1 | 1 |
| $i$/$j$ | **0** | **1** | 0 | 1 | 1 | **2** | **3** | 0 | 0 | 0 |
| Values | **0** | **0** | $a_{000}$ | $a_{111}$ | $a_{011}$ | **0** | **0** | $a_{200}$ | $a_{301}$ | $a_{201}$ |

# Computation Pattern for Tensor Kernels

## MTTKRP

$$\forall i \quad Y_{(i,:)} = \sum_j \underbrace{B(j,:)}_{\text{fiber}} \circ \sum_k \underbrace{A(i,j,k)}_{\text{scalar}} \cdot \underbrace{C(k,:)}_{\text{fiber}}$$

**fiber**    **fiber**   **op**    **scalar**    **fiber**

## Scalar-Fiber product followed by Fiber-Fiber product (SF³) Pattern

$$fibers_{out} = \sum_{D_1} fiber_1 \; op \sum_{D_0} (scalar \cdot fiber_0)$$

**For MTTKRP**

$$[0, J) \qquad\qquad [0, K) \quad \textbf{dense}$$

**or**                **or**

$$\{j \mid \exists k \; st. \; A(i,j,k) \neq 0\} \qquad \{k \mid A(i,j,k) \neq 0\} \quad \textbf{sparse}$$
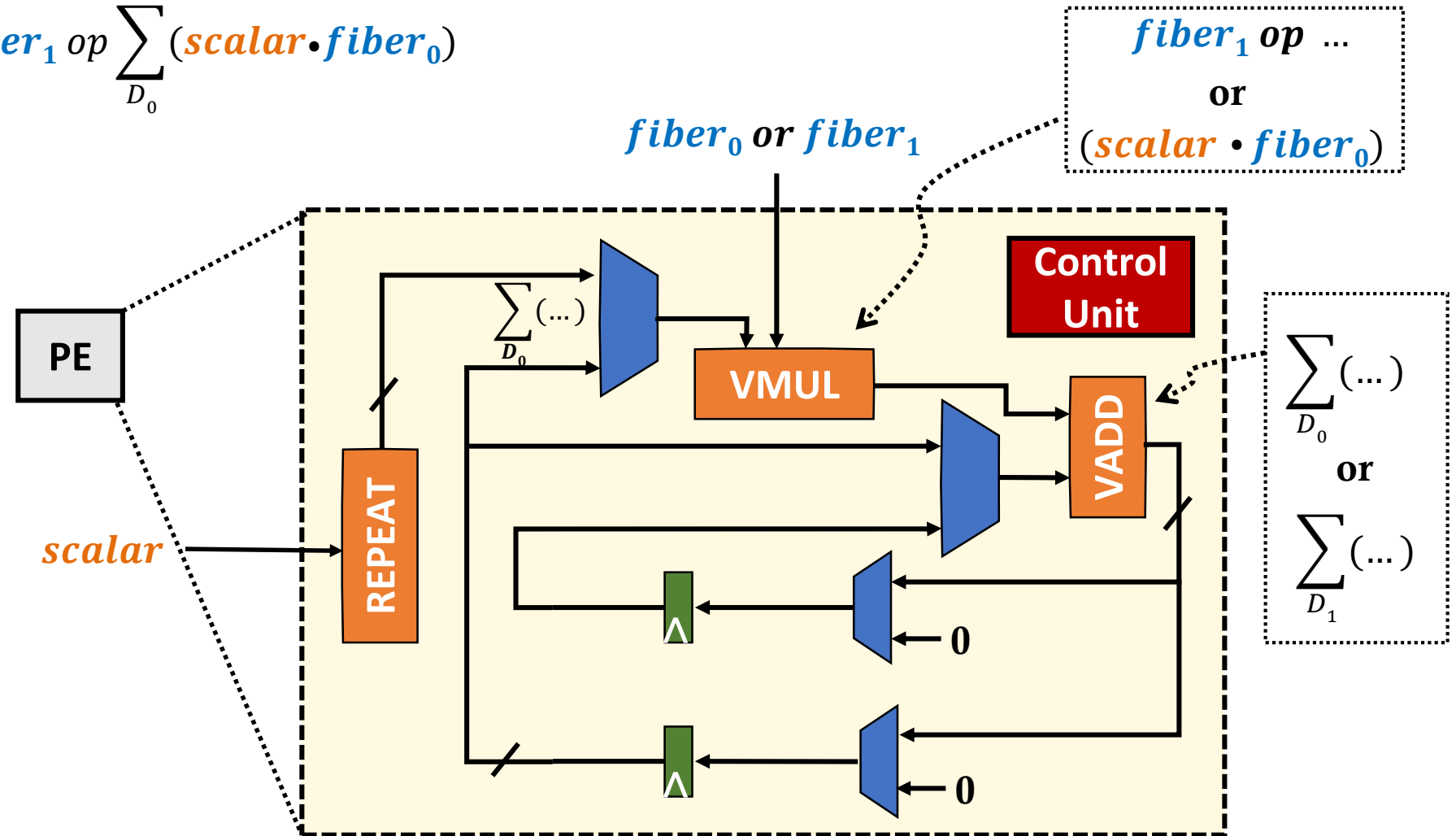
**Similarly**

**TTMc** $\quad \forall i \quad Y(i,:) = \sum_j B(j,:) \otimes \sum_k (A(i,j,k) \cdot C(k,:))$

**MM** $\quad\quad \forall i \quad Y(i,:) = \sum_\phi null \quad op \sum_j (A(i,j) \cdot B(j,:))$

**MV** $\quad\quad \forall i \quad Y(i,:) = \sum_\phi null \quad op \sum_j (A(i,j) \cdot b(j,:))$

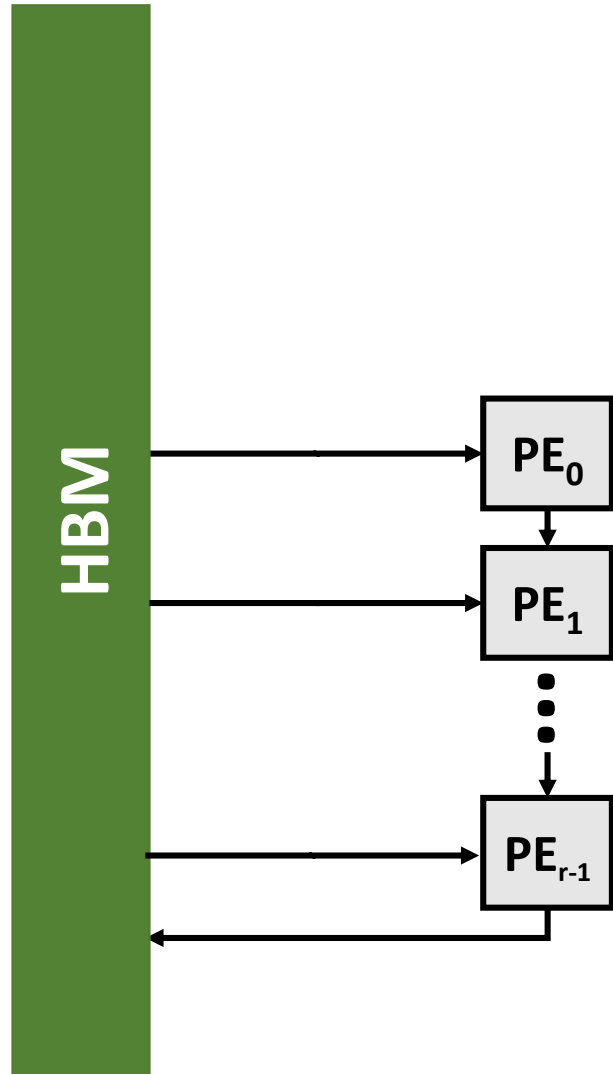SF³ compute pattern can express all the common dense and mixed sparse-dense tensor kernels

# PE for SF³ Compute Pattern

$$fibers_{out} = \sum_{D_1} fiber_1 \; op \sum_{D_0} (scalar \bullet fiber_0)$$

**Different output fibers can be computed in parallel**

$$fiber_1 \; op \; ...$$
**or**
$$(scalar \bullet fiber_0)$$

$fiber_0$ **or** $fiber_1$



$$\sum_{D_0}(...)$$
**or**
$$\sum_{D_1}(...)$$

# Vertical Scaling Using Coarse-Grained Parallelism



HBM

PE$_0$

PE$_1$

PE$_{r-1}$

PE$_0$ :    $fiber_{out}[0]$    $= \sum\limits_{D'_1} fiber_1\ op \sum\limits_{D'_0}(scalar \cdot fiber_0)$

PE$_1$ :    $fiber_{out}[1]$    $= \sum\limits_{D''_1} fiber_1\ op \sum\limits_{D''_0}(scalar \cdot fiber_0)$

PE$_{r-1}$ :  $fiber_{out}[r-1] = \sum\limits_{D'''_1} fiber_1\ op \sum\limits_{D'''_0}(scalar \cdot fiber_0)$

**Long vector SIMD compute**

# Horizontal Scaling Using SIMD-Vector Parallelism

# Tensaurus Architecture



MLU: Matrix Load Unit
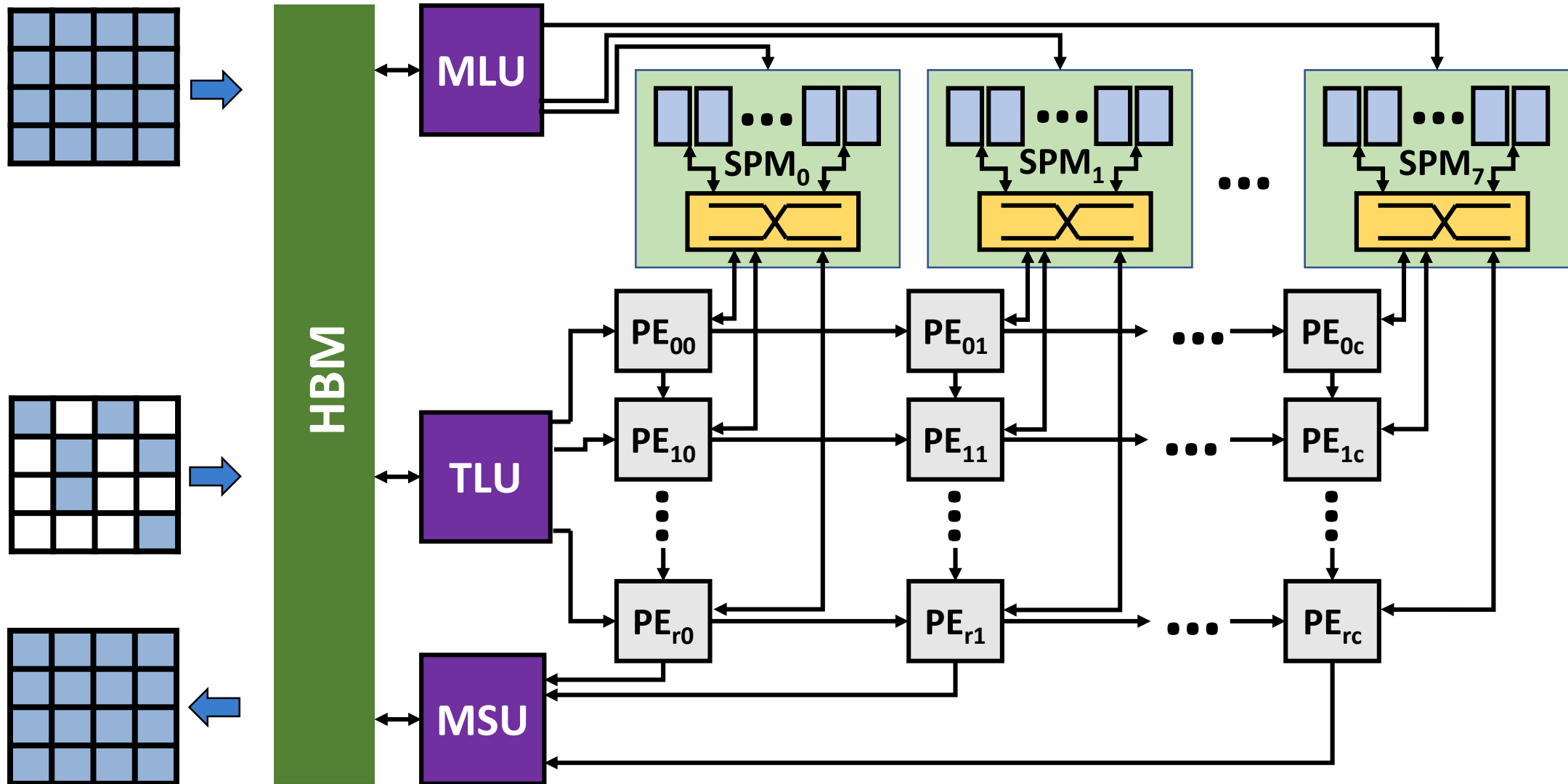TLU: Tensor Load Unit
MSU: Matrix Store Unit

Double Buffer

Crossbar

# Tensaurus Architecture

**Accelerator for both dense and sparse-dense!!**

# Evaluation Methodology

- **Cycle-level simulation in gem5**
  - 8 x 8 PE array, VLEN = 8
  - 8 16KB RAMs per SPM
  - HBM: 8 128-bit physical channels (128 GB/s peak bandwidth)

- **RTL Modeling of a PE using PyMTL**
  - 28 nm (Synopsys & Cadence Tools)

- **Baselines**
  - CPU: Intel(R) Xeon(R) CPU E7-8867
    - SparseBLAS and SPLATT
  - GPU: Titan XP
    - CuSparse, PaRTI
  - Sparse NN Accelerator:
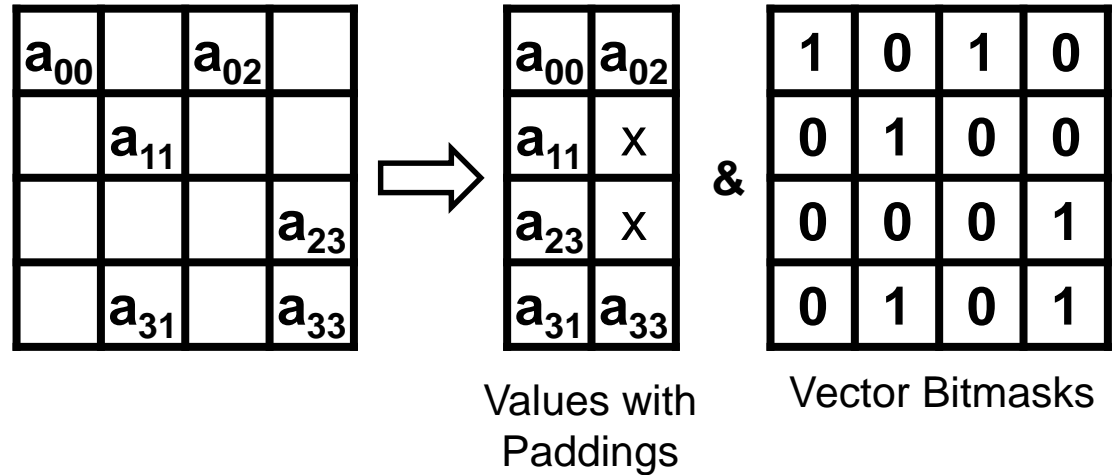    - Cambricon-X [1]

- **Datasets**
  - FROSTT Tensors, Florida Sparse Matrices, AlexNet, VGG-16

**Area and Power Breakdown**

| Component | Area($mm^2$) | % | Power (mW) | % |
|---|---|---|---|---|
| PE | 0.625 | 27.2 % | 402.30 | 40.9 % |
| Xbar | 0.066 | 2.8 % | 24.27 | 2.5% |
| SPM | 0.832 | 36.2 % | 296.05 | 30.1 % |
| MSU | 0.759 | 33.0 % | 247.03 | 25.2 % |
| TLU | 0.009 | 0.4 % | 6.28 | 0.6% |
| MLU | 0.009 | 0.4 % | 6.28 | 0.6 % |
| Total | 2.3 | 100 % | 982.21 | 100 % |

[1] Zhang, Shijin, et al. "Cambricon-x: An accelerator for sparse neural networks.", *Int'l Symp. on Microarchitecture (MICRO)*, 2016.

# Cambricon-X Baseline
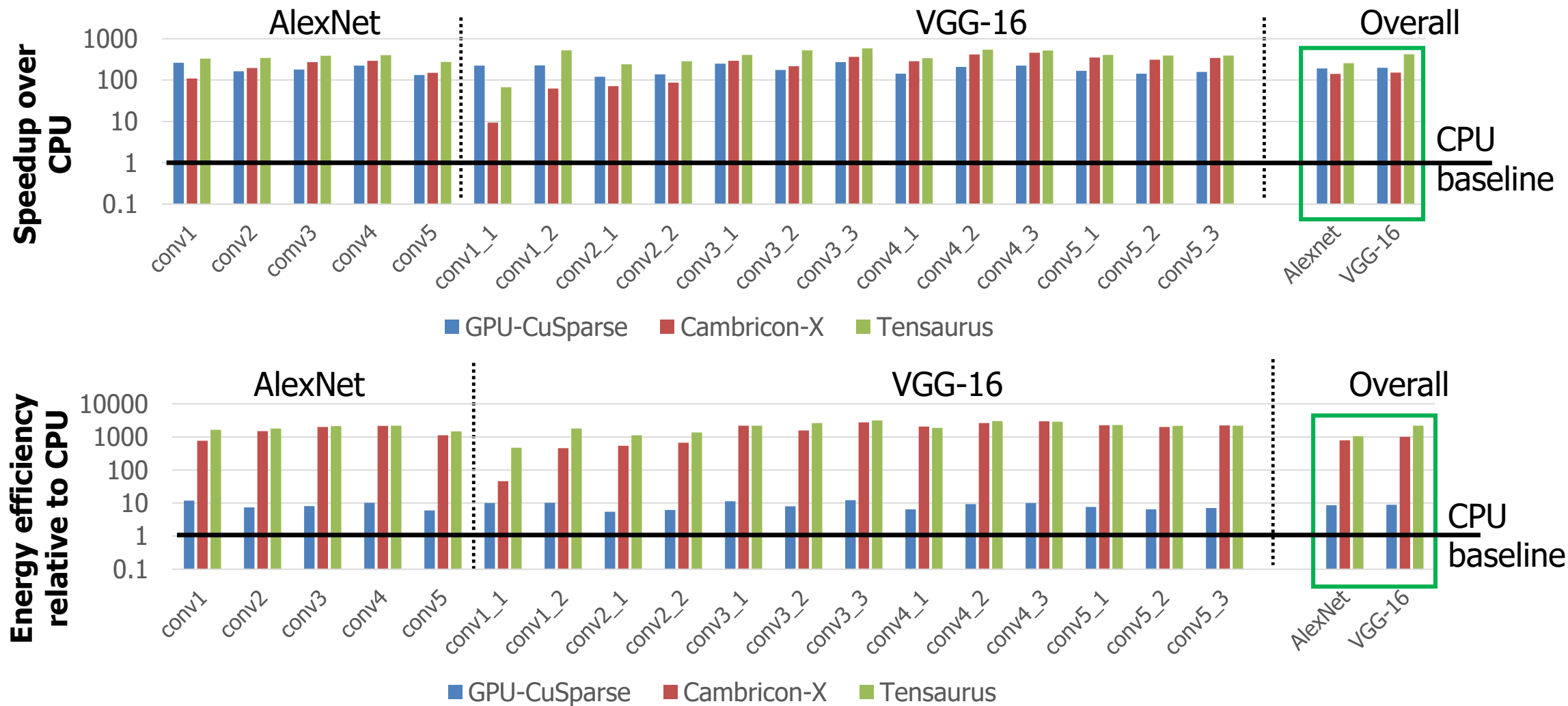
▸ Cambricon-X uses a CSR-variant
  – Pads empty entries with padding (x)
  – Uses vector bit-masks to indicate non-zero positions
  – Specialized for CNNs with low sparsity



Values with Paddings    Vector Bitmasks

▸ CSR results in load-imbalanced schedule
  – Cambricon-X has synchronization boundaries across rows



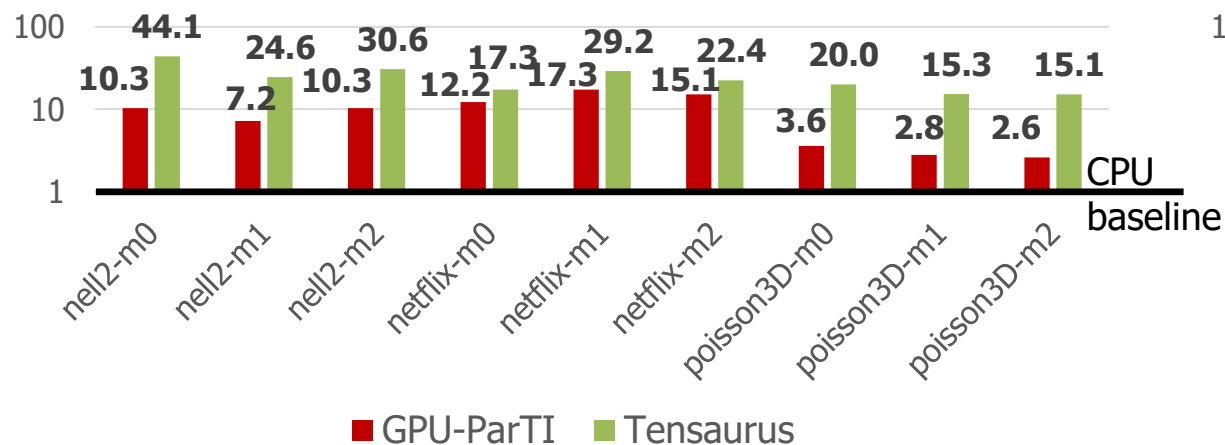Load imbalance due to synchronization at row boundaries
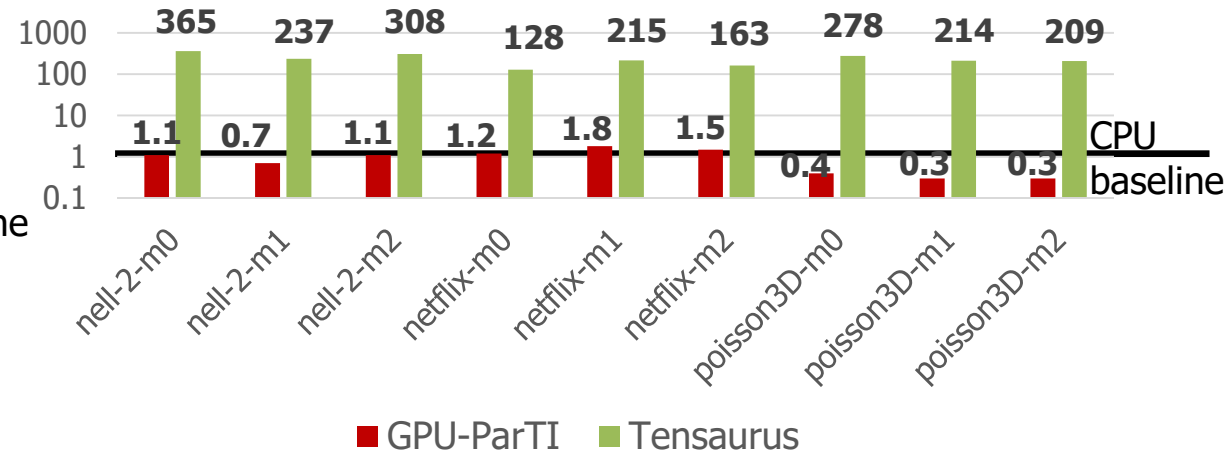
# Results on Sparse Neural Nets



Overall Tensaurus is 1.9x faster and 1.7x more energy-efficient than Cambricon-X even for Sparse Neural Nets

# Results on Sparse Tensor Decomposition



**Speedup for MTTKRP**

- GPU-ParTI
- Tensaurus

nell2-m0: 10.3, 44.1
nell2-m1: 7.2, 24.6
nell2-m2: 10.3, 30.6
netflix-m0: 12.2, 17.3
netflix-m1: 17.3, 29.2
netflix-m2: 15.1, 22.4
poisson3D-m0: 3.6, 20.0
poisson3D-m1: 2.8, 15.3
poisson3D-m2: 2.6, 15.1

**Performance/Watt for MTTKRP**

- GPU-ParTI
- Tensaurus

nell-2-m0: 1.1, 365
nell-2-m1: 0.7, 237
nell-2-m2: 1.1, 308
netflix-m0: 1.2, 128
netflix-m1: 1.8, 215
netflix-m2: 1.5, 163
poisson3D-m0: 0.4, 278
poisson3D-m1: 0.3, 214
poisson3D-m2: 0.3, 209

**Tensaurus is 22.9x & 3.1x faster, and 220x & 290x more energy-efficient than CPU & GPU for MTTKRP**

# Concluding Remarks

▸ **Tensaurus: A versatile accelerator for sparse-dense tensor acceleration**

- First accelerator for sparse tensor decompositions (MTTKRP, TTMc)

- **Versatile:** NOT limited to tensor decompositions. Also efficient for <u>sparse-dense matrix computations</u>

- **Adaptable:**

  - Also <u>accelerates dense kernels</u>

  - Easily <u>adapts to different levels of sparsity</u> found in various domains


▸ **Key Approach:** Co-design <u>sparse format</u> and <u>architecture</u>


▸ **Key Results:**

- High bandwidth utilization (> 70% of peak bandwidth)

- High speedup and energy efficiency compared to CPU, GPU and Cambricon-X

# Thank you! Questions?

## *Tensaurus*: A Versatile Accelerator for Mixed Sparse-Dense Tensor Computations

**Nitish Srivastava**, Hanchen Jin, Shaden Smith[2], Hongbo Rong[3],

David Albonesi, and Zhiru Zhang

Cornell University

[2]Microsoft AI & Research

[3]Intel Parallel Computing Lab