# Onyx: A Programmable Accelerator for Sparse Tensor Algebra

Kalhan Koul[1], Maxwell Strange[1], Jackson Melchert[1], Alex Carsello[1], Yuchen Mei[1], Olivia Hsu[1], Taeyoung Kong[1], Po-Han Chen[1], Huifeng Ke[1], Keyi Zhang[1], Qiaoyi Liu[1], Gedeon Nyengele[1], Akhilesh Balasingam[1], Jayashree Adivarahan[1], Ritvik Sharma[1], Zhouhua Xie[1], Christopher Torng[2], Joel Emer[3], Fredrik Kjolstad[1], Mark Horowitz[1], Priyanka Raina[1]

[1]Stanford University, CA, USA; [2]University of Southern California, CA, USA; [3]MIT, MA, USA
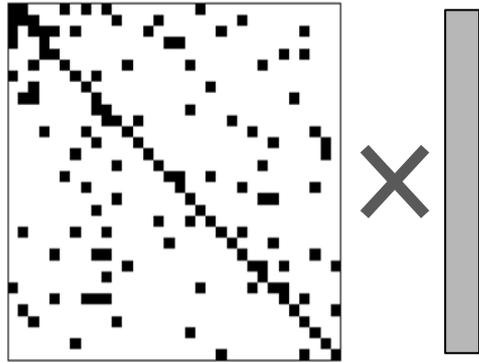
# Sparse Applications

- Applications ranging from scientific computing to machine learning can have **extremely** sparse inputs

# Sparse Applications

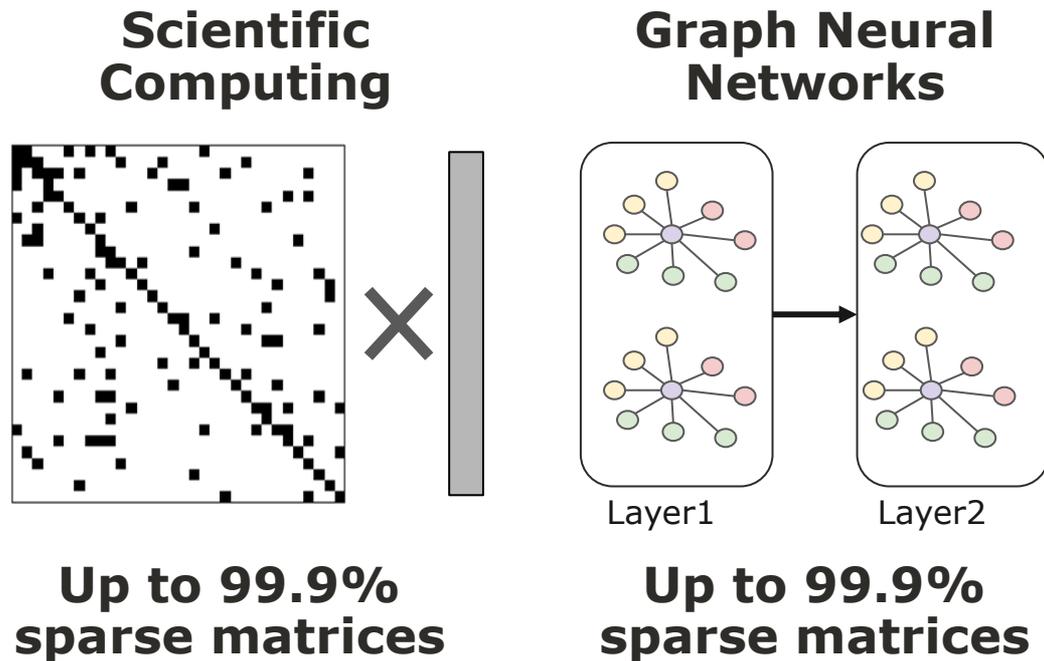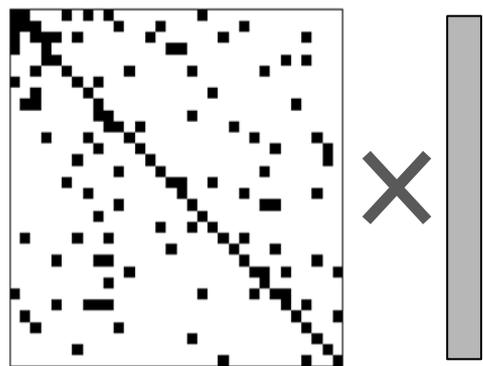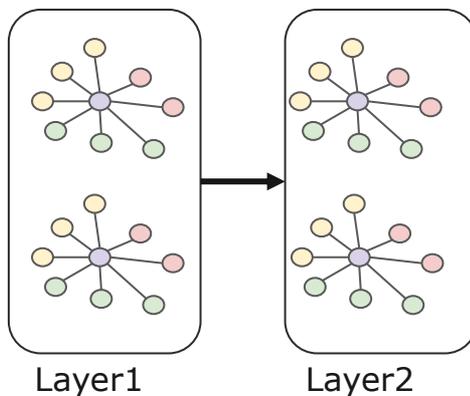- Applications ranging from scientific computing to machine learning can have **extremely** sparse inputs

**Scientific Computing**



**Up to 99.9% sparse matrices**

# Sparse Applications

- Applications ranging from scientific computing to machine learning can have **extremely** sparse inputs



**Scientific Computing**

**Up to 99.9% sparse matrices**

**Graph Neural Networks**

Layer1     Layer2

**Up to 99.9% sparse matrices**

# Sparse Applications

- Applications ranging from scientific computing to machine learning can have **extremely** sparse inputs
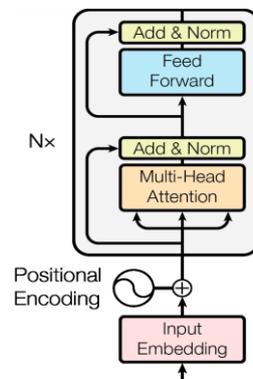


**Scientific Computing**

**Up to 99.9% sparse matrices**

**Graph Neural Networks**

Layer1    Layer2
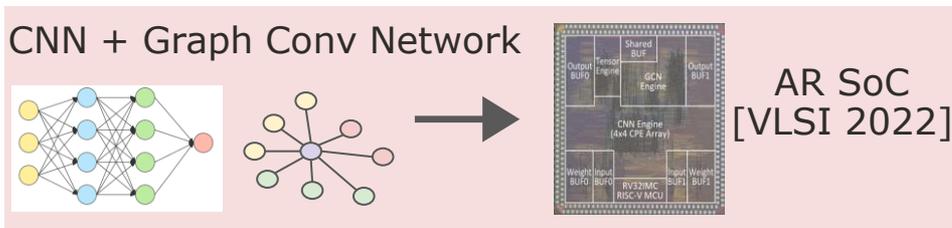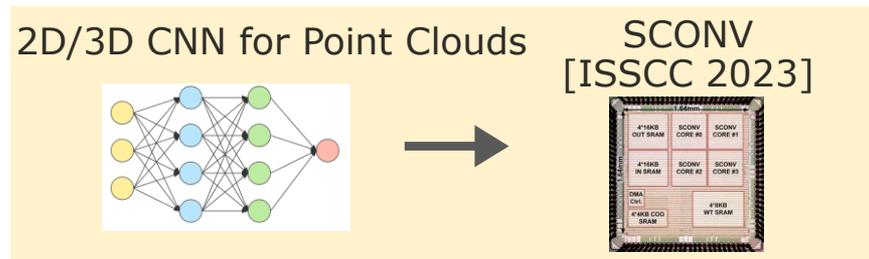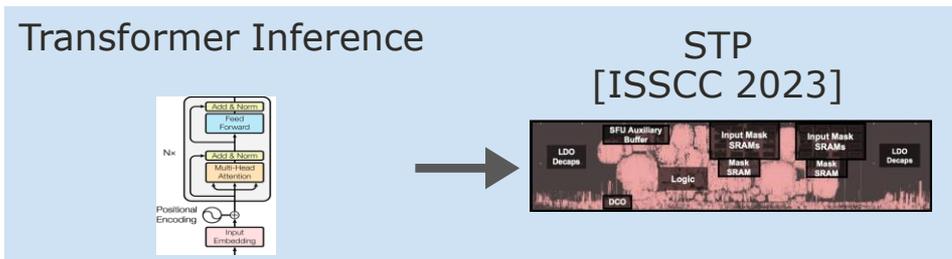
**Up to 99.9% sparse matrices**

**Sparse Transformers**

**Up to 93% sparse matrices**

# End-to-End Sparse Accelerators

- Hardware accelerators for end-to-end applications exploiting sparsity are fixed function and become obsolete as applications evolve



Transformer Inference — STP [ISSCC 2023]

2D/3D CNN for Point Clouds — SCONV [ISSCC 2023]

CNN + Graph Conv Network — AR SoC [VLSI 2022]

3D Navigation — GPPU [VLSI 2023]

# End-to-End Sparse Accelerators
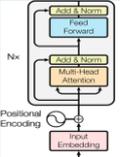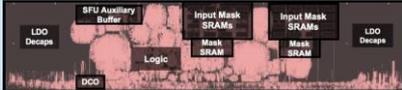
- Hardware accelerators for end-to-end applications exploiting sparsity are fixed function and become obsolete as applications evolve
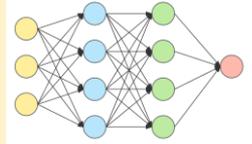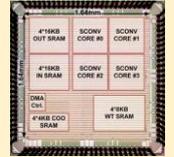


Transformer Inference → STP [ISSCC 2023]
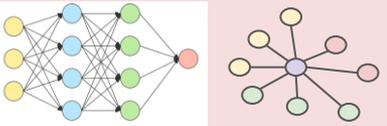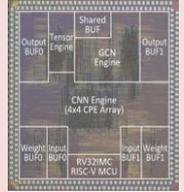
2D/3D CNN for Point Clouds → SCONV [ISSCC 2023]
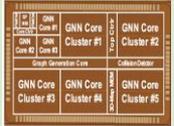
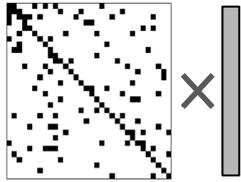CNN + Graph Conv Network → AR SoC [VLSI 2022]

3D Navigation → GPPU [VLSI 2023]

✅ Energy and area efficient
❌ Accelerate a fixed type of application
❌ Do not adapt to evolving application domains

# Kernel-Level Sparse Accelerators

- Sparse kernel accelerators focus only on sparse matrix-matrix or matrix-vector multiplication

Supported



matrix times vector

matrix times matrix

# Kernel-Level Sparse Accelerators

- Sparse kernel accelerators focus only on sparse matrix-matrix or matrix-vector multiplication
- These accelerators leave out support for high-dimensional tensors, multi-input complex expressions and fusion

Supported                                    Unsupported



matrix times vector


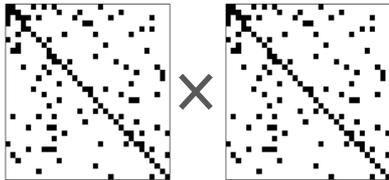
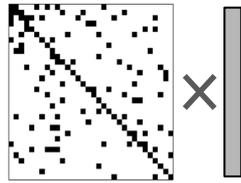tensor times matrix

matrix times matrix
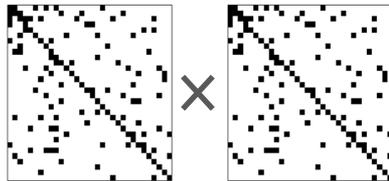
sampled matrix times matrix

# Kernel-Level Sparse Accelerators

- Sparse kernel accelerators focus only on sparse matrix-matrix or matrix-vector multiplication
- These accelerators leave out support for high-dimensional tensors, multi-input complex expressions and fusion
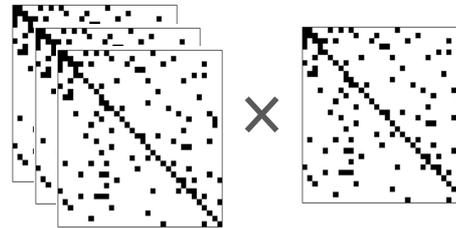
Supported

Unsupported

matrix times vector

matrix times matrix

tensor times matrix

sampled matrix times matrix

Require a highly programmable accelerator!

# Programmable Accelerators

- But most programmable accelerators **only** focus on dense applications



Amber CGRA
[VLSI 2022]



MIMO CGRA
[JSSC 2020]

✅ Energy and area efficiency approaching ASICs

✅ Accelerate a domain of applications

✅ Adaptable to evolving dense application domains

# Onyx

- Programmable accelerator for <u>any</u> tensor algebra expression



Graph Neural Networks

Scientific Computing

Sparse Transformers

**Sparse/ Dense Compiler**

CNNs/ MLPs

Image Processing

# Onyx

- Programmable accelerator for <u>any</u> tensor algebra expression



Graph Neural Networks

Scientific Computing

Sparse Transformers

**Sparse/ Dense Compiler**

CNNs/ MLPs

Image Processing

$$A = Bc + a \qquad a = Bc$$
$$A = B \odot C \qquad A = B + C \qquad a = \alpha Bc + \beta a$$
$$A = \alpha B \quad A = 0 \quad A = BC$$
$$A = BCd$$
$$a = b \odot c \qquad A = B \odot (CD)$$
$$A_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj} \qquad A = B^T \qquad a = B^T Bc$$
$$A_{ik} = \sum_j B_{ijk} c_j \quad A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij}$$
$$A_{ijk} = \sum_l B_{ikl} C_{lj} \qquad A_{ij} = (\sum_k B_{ijk} C_{ijk}) + D_{ij}$$
$$C = \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \, \overline{P_{il}} \qquad \tau = \sum_i z_i (\sum_j z_j \theta_{ij})(\sum_k z_k \theta_{ik})$$
$$a = \sum_{ijklmnop} M_{ij} P_{jk} M_{kl} P_{lm} \overline{M_{nm}} P_{no} \overline{M_{po}} \, \overline{P_{ip}}$$

**Arbitrary sparse/ dense tensor algebra expressions**

Sparse abstract machine (SAM)/CoreIR dataflow graph

# Onyx

- Programmable accelerator for <u>any</u> tensor algebra expression



Graph Neural Networks

Scientific Computing

Sparse Transformers

**Sparse/ Dense Compiler**

CNNs/ MLPs

Image Processing

$$A = Bc + a \qquad a = Bc$$
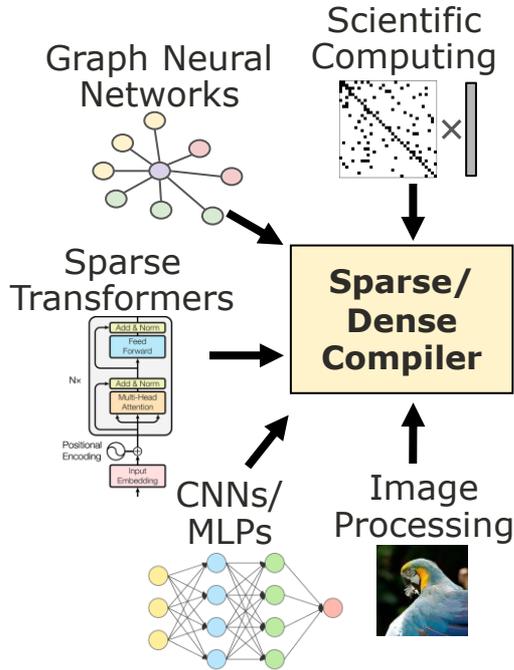$$A = B \odot C \qquad A = B + C \qquad a = \alpha Bc + \beta a$$
$$A = BCd \qquad A = \alpha B \quad A = 0 \quad A = BC$$
$$a = b \odot c \qquad A = B \odot (CD)$$
$$A_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj} \qquad A = B^T \quad a = B^T Bc$$
$$A_{ik} = \sum_{j} B_{ijk} c_j \quad A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij}$$
$$A_{ijk} = \sum_{l} B_{ikl} C_{lj} \qquad A_{ij} = (\sum_{k} B_{ijk} C_{ijk}) + D_{ij}$$
$$C = \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \, \overline{P_{il}} \qquad \tau = \sum_{i} z_i (\sum_{j} z_j \theta_{ij})(\sum_{k} z_k \theta_{ik})$$
$$a = \sum_{ijklmnop} M_{ij} P_{jk} M_{kl} P_{lm} \overline{M_{nm}} P_{no} \overline{M_{po}} \, \overline{P_{ip}}$$

**Arbitrary sparse/ dense tensor algebra expressions**

Sparse abstract machine (SAM)/CoreIR dataflow graph

**Onyx SoC**
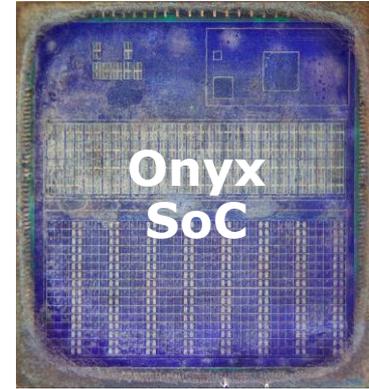
- ☑ Energy and area efficient
- ☑ Accelerates an application domain
- ☑ Configurable for sparse + dense
- ☑ Adaptable to fast evolving domains

# Onyx Architecture

# Onyx Architecture

# Onyx Architecture

# Contributions

# Contributions



Sparse Acceleration Hardware

Composable primitives for accelerating arbitrary sparse tensor algebra kernels

# Contributions



## Sparse Acceleration Hardware

Composable primitives for accelerating arbitrary sparse tensor algebra kernels

PE Tile

U

n

MEM Tile

## Dense Acceleration Improvements

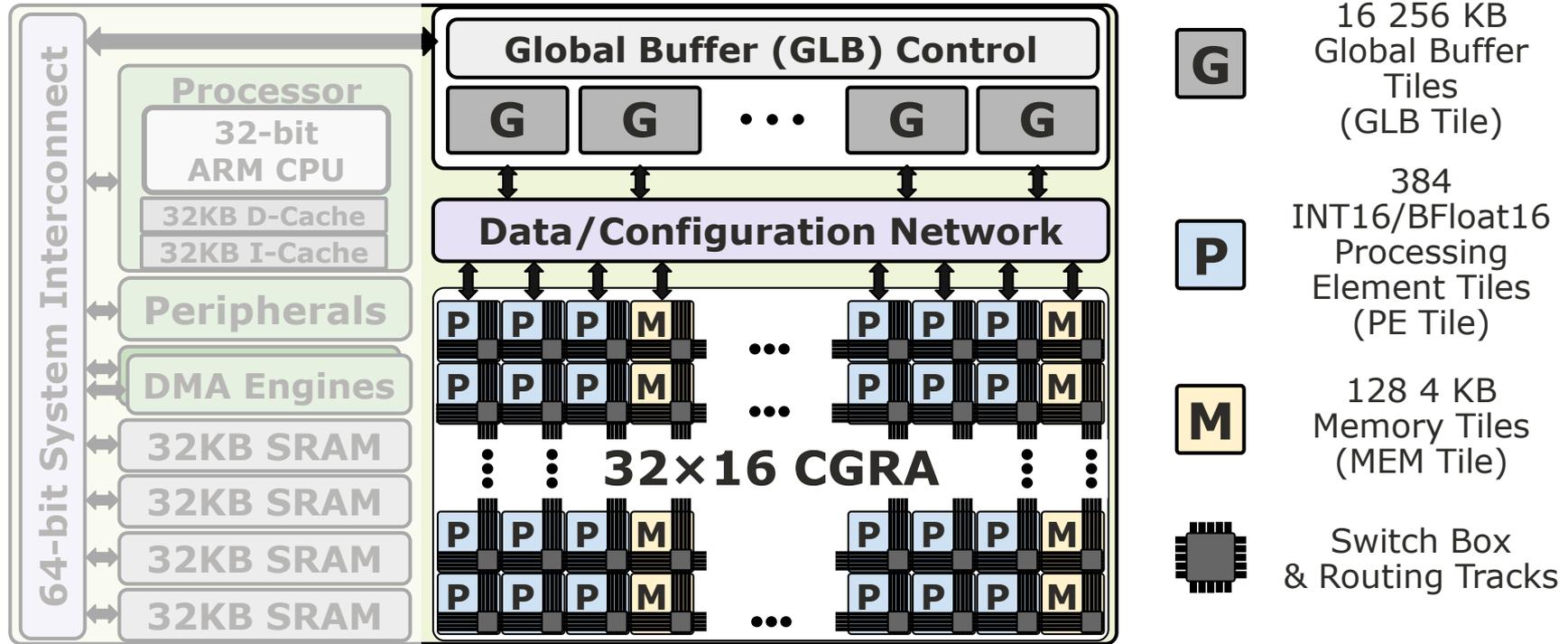Compute and memory controller optimizations for dense applications

MAC

⋮

CROP

New PE OPs

Onyx

Amber

PE Usage per App

Active

Gated

Gated

clk

read

gclk

Dynamic-/Static- GLB Clock Gating

# Contributions



## Sparse Acceleration Hardware

Composable primitives for accelerating arbitrary sparse tensor algebra kernels

## Dense Acceleration Improvements

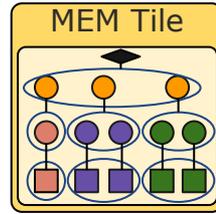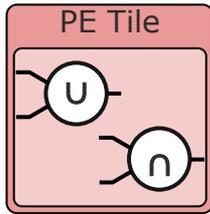Compute and memory controller optimizations for dense applications

# Sparse Tensor Abstraction

- Sparse tensors can be represented as fibertree [1] structures

**Tensor**



[1] V. Sze et al., M&C'2020.

# Sparse Tensor Abstraction

- Sparse tensors can be represented as fibertree [1] structures

**Tensor**

*dimension $j$*

| | | | |
|---|---|---|---|
| 0 | 1 | 0 | 0 |
| 2 | 0 | 3 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 4 | 0 | 5 |

*dimension $i$* — 0, 1, 2, 3

**Fibertree**

level

*coordinate*

fiber

$i$ ( 0 ) ( 1 ) ( 3 )

[1] V. Sze et al., M&C'2020.

# Sparse Tensor Abstraction

- Sparse tensors can be represented as fibertree [1] structures

**Tensor**

*dimension $j$*

|   | 0 | **1** | 2 | **3** |
|---|---|---|---|---|
| 0 | 0 | 1 | 0 | 0 |
| 1 | 2 | 0 | 3 | 0 |
| 2 | 0 | 0 | 0 | 0 |
| **3** | 0 | 4 | 0 | 5 |

*dimension $i$*

**Fibertree**

level

*coordinate*

fiber

reference

$i$ : 0 1 3

$j$ : 1 | 0 2 | 1 3

[1] V. Sze et al., M&C'2020.

# Sparse Tensor Abstraction

- Sparse tensors can be represented as fibertree [1] structures



**Tensor**

**Fibertree**

[1] V. Sze et al., M&C'2020.

# Stream and Storage Format

- Fibertree in stream format

**Fibertree**

level

$i$  0  1  3

$j$  1  0  2  1  3

val  1  2  3  4  5

*coord*
*fiber*
*ref*

*coordinates: coord*
*references: ref*
*segments: seg*        *values: val*

**Stream Format**

MEM Tile

$i$ coord    D S$_0$ **3 1 0**

S:stop
D:done

PE Tile

# Stream and Storage Format

- Fibertree in stream format

**Fibertree**



coordinates: coord
references: ref
segments: seg          values: val

**Stream Format**

# Stream and Storage Format

- Fibertree in storage format



**Fibertree**

level

*coord*

*fiber*

*ref*

*i*  0  1  3

*j*  1 | 0  2 | 1  3

val  1 | 2  3 | 4  5

*coordinates: coord*
*references: ref*
*segments: seg*  *values: val*

**Stream Format**

*i* {
coord  $D\ S_0\ \mathbf{3\ 1\ 0}$   S:stop
ref  $D\ S_0\ 2\ 1\ 0$   D:done
}

*j* {
coord  $D\ S_1\ 3\ 1$ | $S_0\ 2\ 0$ | $S_0\ 1$
ref  $D\ S_1\ 4\ 3$ | $S_0\ 2\ 1$ | $S_0\ 0$
val  $D\ S_1\ 5\ 4$ | $S_0\ 3\ 2$ | $S_0\ 1$
}

← time

**Storage Format**

*denote fiber length

*i* {
| 0 | 3 |   seg array
| **0** | **1** | **3** |   coord array
}

*j* {
| 0 | 1 | 3 | 5 |
| 1 | 0 | 2 | 1 | 3 |
}

| 1 | 2 | 3 | 4 | 5 |

# Stream and Storage Format

- Fibertree in storage format



**Fibertree**

level

$i$ — 0 1 3 *coord* ···· fiber ···· ref

$j$ — 1 | 0 2 | 1 3

val — 1 | 2 3 | 4 5

*coordinates: coord*
*references: ref*
*segments: seg*

**Stream Format**

$i$ {
 coord   D $S_0$ 3 1 0    S:stop
 ref     D $S_0$ 2 1 **0**    D:done
}

$j$ {
 coord   D $S_1$ 3 1 | $S_0$ 2 0 | $S_0$ 1
 ref     D $S_1$ 4 3 | $S_0$ 2 **1** | $S_0$ 0
}
 val     D $S_1$ 5 4 | $S_0$ 3 2 | $S_0$ 1

time

*values: val*

**Storage Format**

*denote fiber length

$i$ {
 | 0 | 3 |    seg array
 | 0 | 1 | 3 |    coord array
}

$j$ {
 | **0** | **1** | 3 | 5 |
 | 1 | 0 | 2 | 1 | 3 |
}
 | 1 | **2** | **3** | 4 | 5 |

# Sparse Acceleration Hardware

- MEMs: sparse memory controller logic to convert between the stream and storage representation of fibertrees

- PEs: contains logic to eliminate all ineffectual compute and support tensor algebra



MEM Tile

Level Writer
Level Scanner
Level Buffer

PE Tile

Intersecter
Unioner
Coordinate Dropper
Reducer
Repeater

# Primitives in MEM Tile - Level Writer
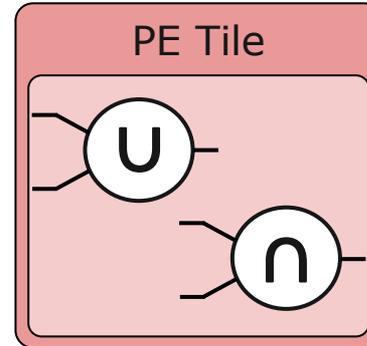
- The level writer converts a stream to the storage format



**Level Writer**

coord_i

*Stream Format*

$DS_0 310$

coord

Segment Gen (Coordinate Counter) — seg

Writer FSM

Coordinate Addr Gen — coord / addr

Segment Addr Gen — seg / addr

data

addr

**Stream Format** → **Storage Format**

coord  $D \ S_0 3 \ 1 \ 0$

*denote fiber length

*i*

| 0 | 3 | seg array |
| 0 | 1 | 3 | coord array |

# Primitives in MEM Tile - Level Writer

- The level writer converts a stream to the storage format



**Level Writer**

coord_i

*Stream Format*

$DS_0 310$

coord

Segment Gen
(Coordinate Counter)

seg

Writer FSM

Coordinate Addr Gen — coord / addr

Segment Addr Gen — seg / addr

data
addr

**Stream Format** ⟶ **Storage Format**

coord  $D\ S_0 3\ 1\ 0$

*denote fiber length
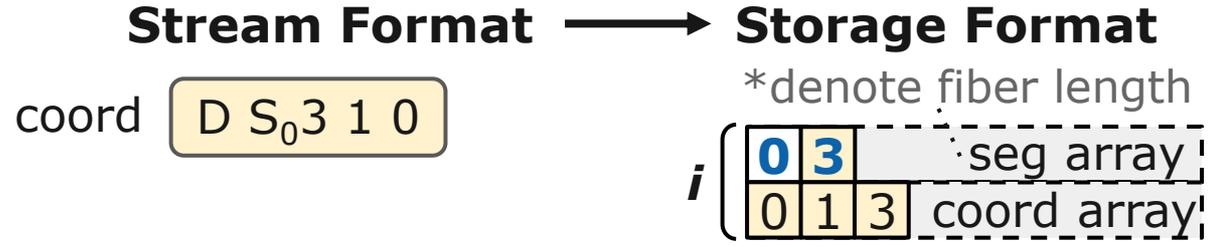
$i$ | 0 3 — seg array
| 0 1 3 — coord array

# Primitives in MEM Tile - Level Writer

- The level writer converts a stream to the storage format



**Level Writer**

coord_i

*Stream Format*

$DS_03 1 0$

coord

Segment Gen
(Coordinate Counter) — seg

Writer FSM

**Coordinate Addr Gen** — coord / addr

**Segment Addr Gen** — seg / addr

data

addr

**Stream Format** ⟶ **Storage Format**

coord $\quad$ $D\ S_0 3\ 1\ 0$

*denote fiber length

*i*

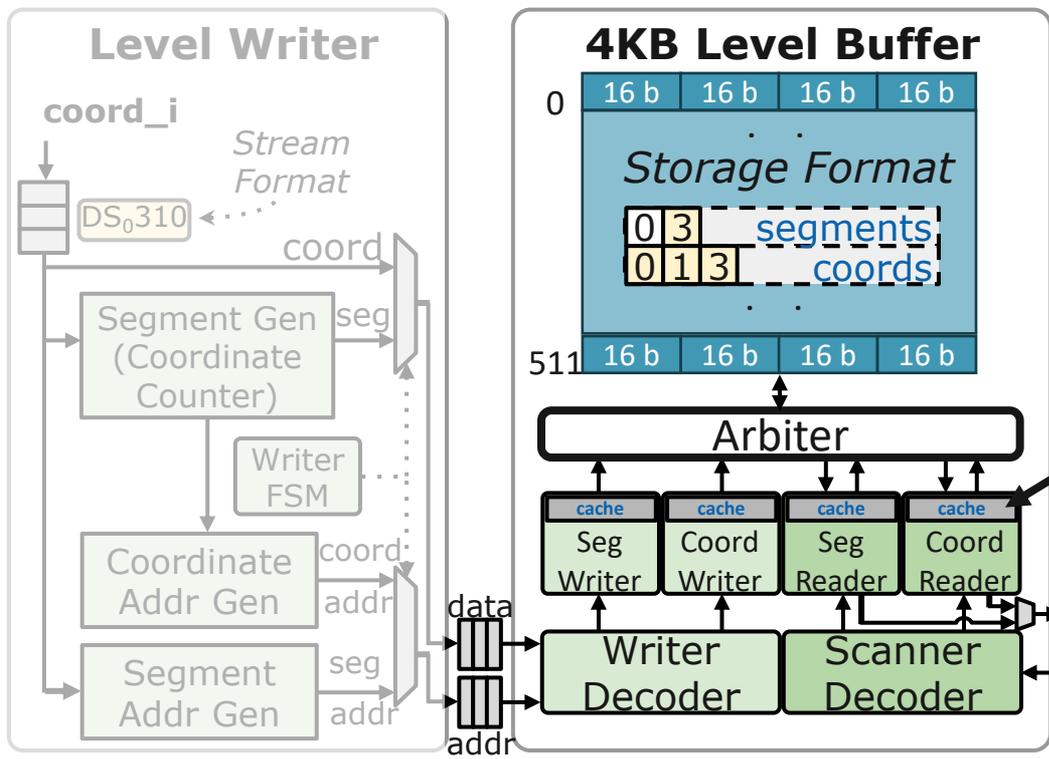| 0 | 3 | seg array |
| 0 | 1 | 3 | coord array |

# Primitives in MEM Tile - Level Buffer

- The level buffer arbitrates writes/reads with support for caching
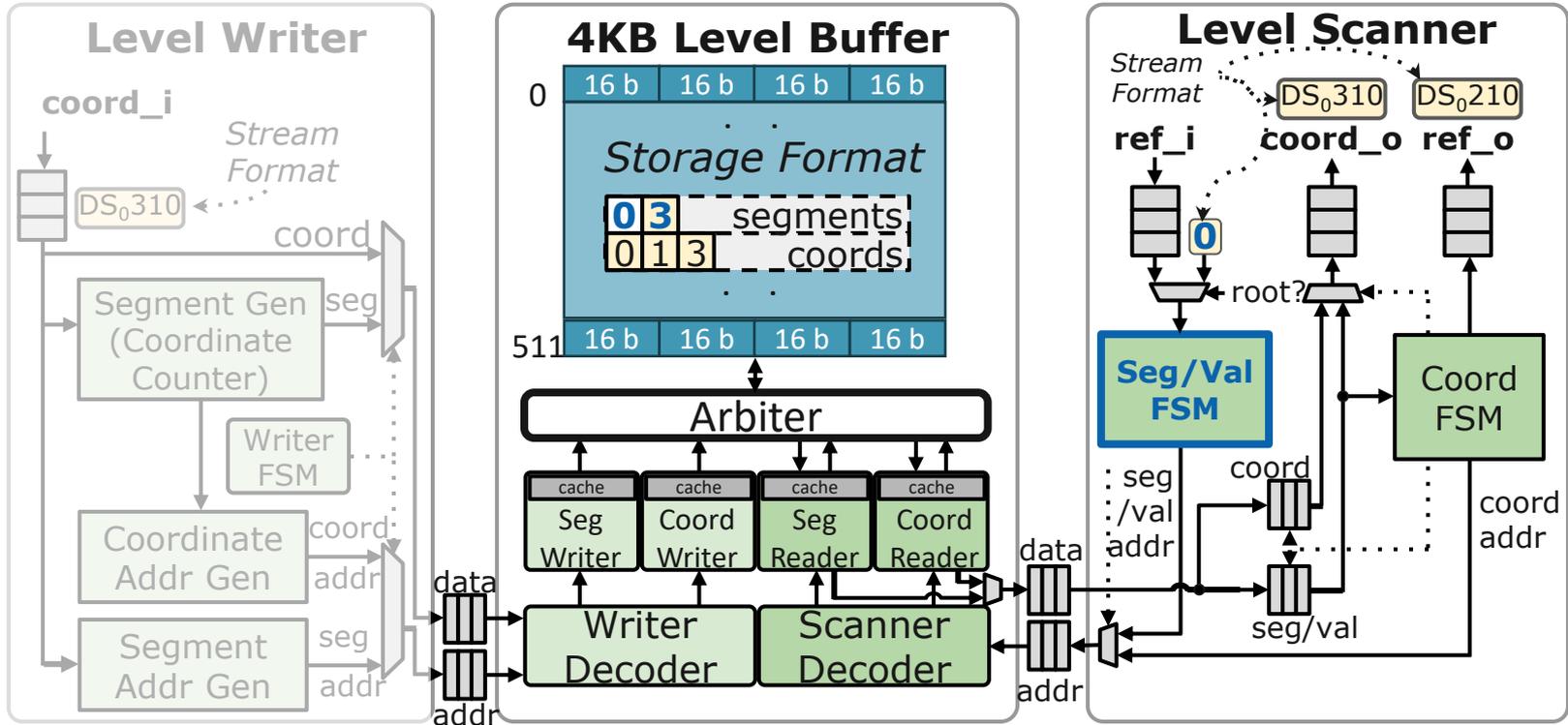
# Primitives in MEM Tile – Level Buffer

- The level buffer arbitrates writes/reads with support for caching



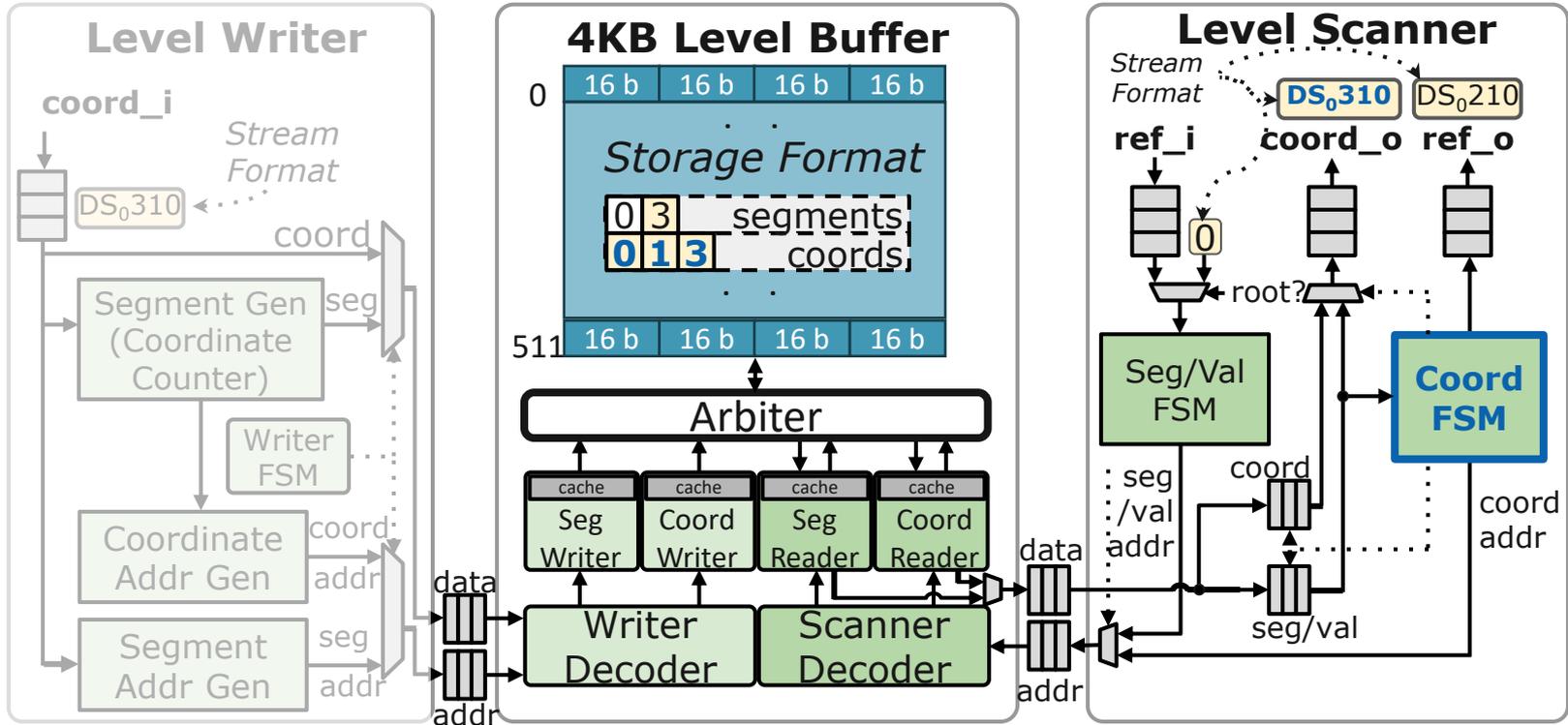Avoid redundant reads of the same 4-element word
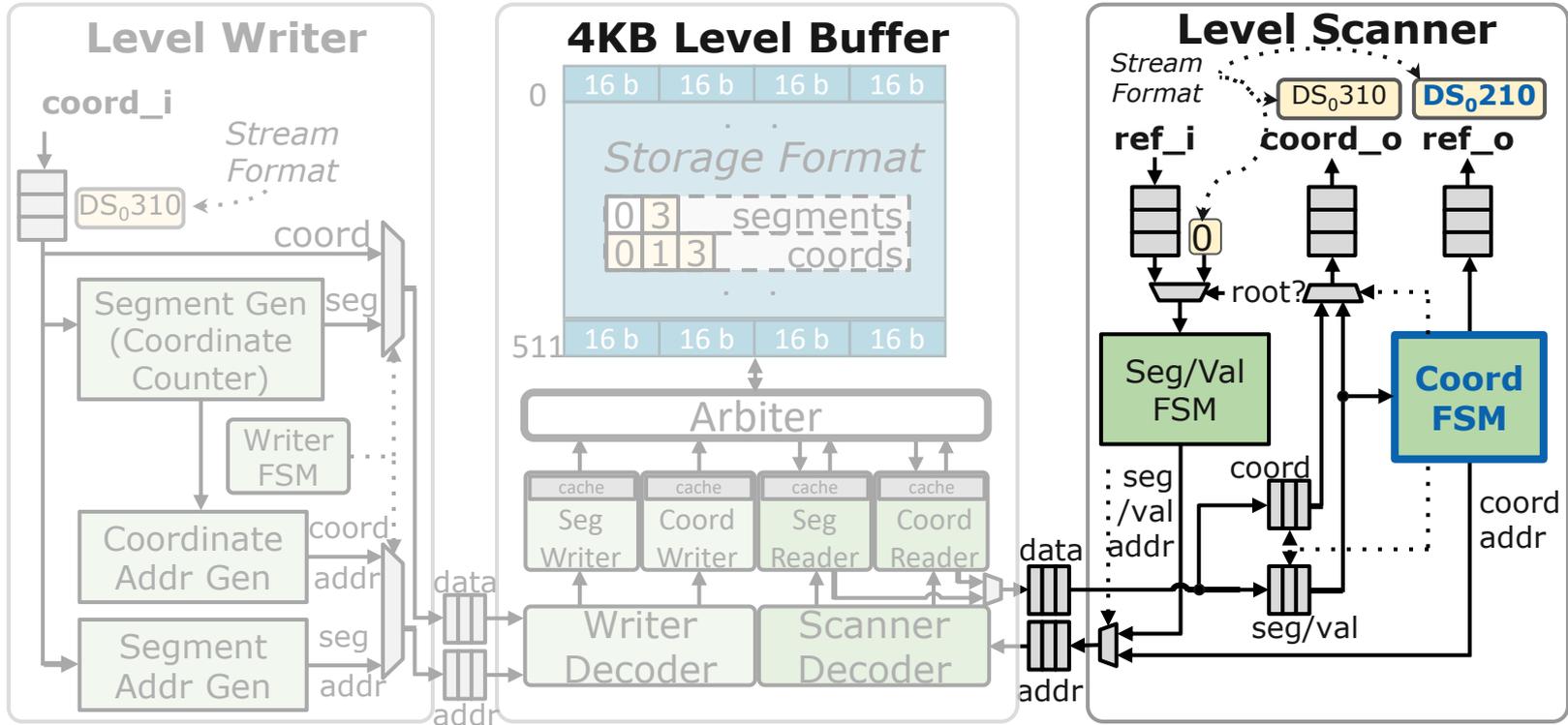
# Primitives in MEM Tile - Level Scanner

- The level scanner produces the next level fiber for a reference

# Primitives in MEM Tile - Level Scanner

- The level scanner produces the next level fiber for a reference

# Primitives in MEM Tile - Level Scanner

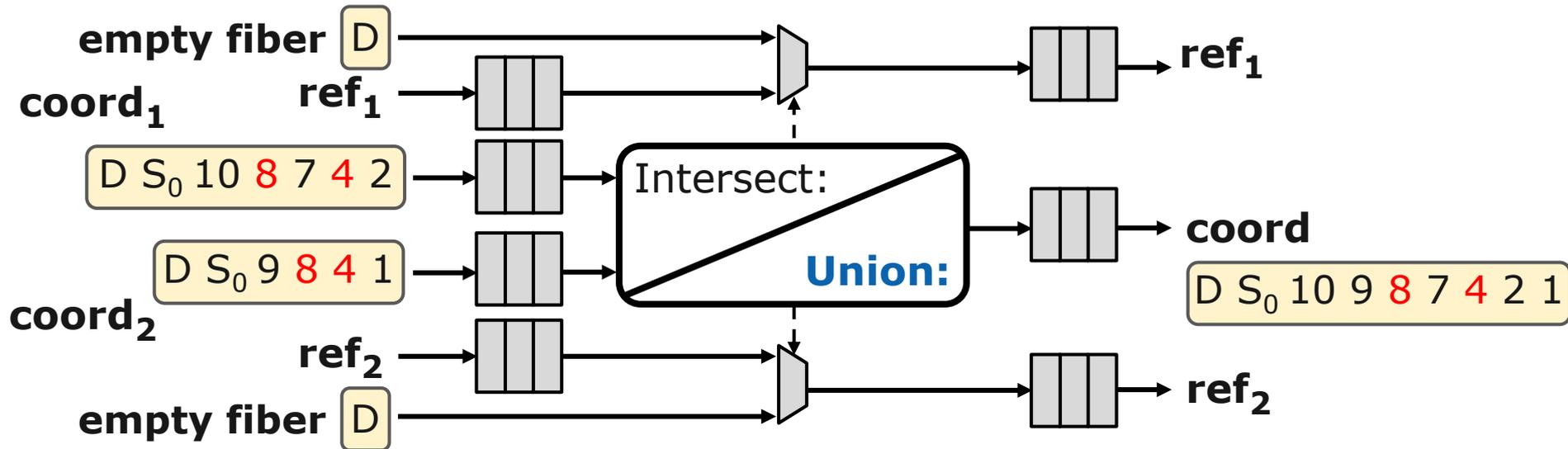- The level scanner produces the next level fiber for a reference

# Primitives in MEM Tile

- Store and process any-dimensional tensors

# Primitives in PE Tile - <u>Intersecter</u> / Unioner

- The intersecter compares and keeps equivalent coordinates



empty fiber D → ref$_1$

coord$_1$ ref$_1$

D S$_0$ 10 8 7 4 2

Intersect:

Union:

D S$_0$ 9 8 4 1

coord$_2$

ref$_2$

empty fiber D → ref$_2$

D S$_0$ 8 4

coord

# Primitives in PE Tile - Intersecter / <u>Unioner</u>

- The unioner collects all incoming coordinates



**empty fiber** D → ref$_1$

**coord$_1$**

ref$_1$

D S$_0$ 10 8 7 4 2

Intersect:

D S$_0$ 9 8 4 1

**Union:**

**coord$_2$**

ref$_2$

**empty fiber** D

→ coord

D S$_0$ 10 9 8 7 4 2 1

→ ref$_2$

# Primitives in PE Tile - Coordinate Dropper

- The coordinate dropper removes empty fibers after an intersect



**inner coord**

D $S_0$ 3 2 1 0

D $S_1$ 3 1 $S_0$ $S_0$ 2 0 $S_0$ 1

**outer coord**

== S ?

**inner coord**

D $S_0$ 3 1 0
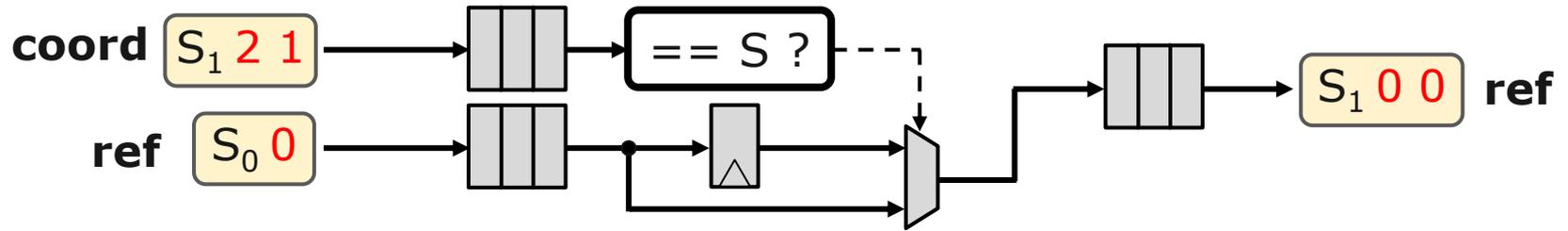
D $S_1$ 3 1 $S_0$ 2 0 $S_0$ 1

**outer coord**

# Primitives in PE Tile - Reducer

- The reducer accumulates over a fiber

# Primitives in PE Tile - Repeater
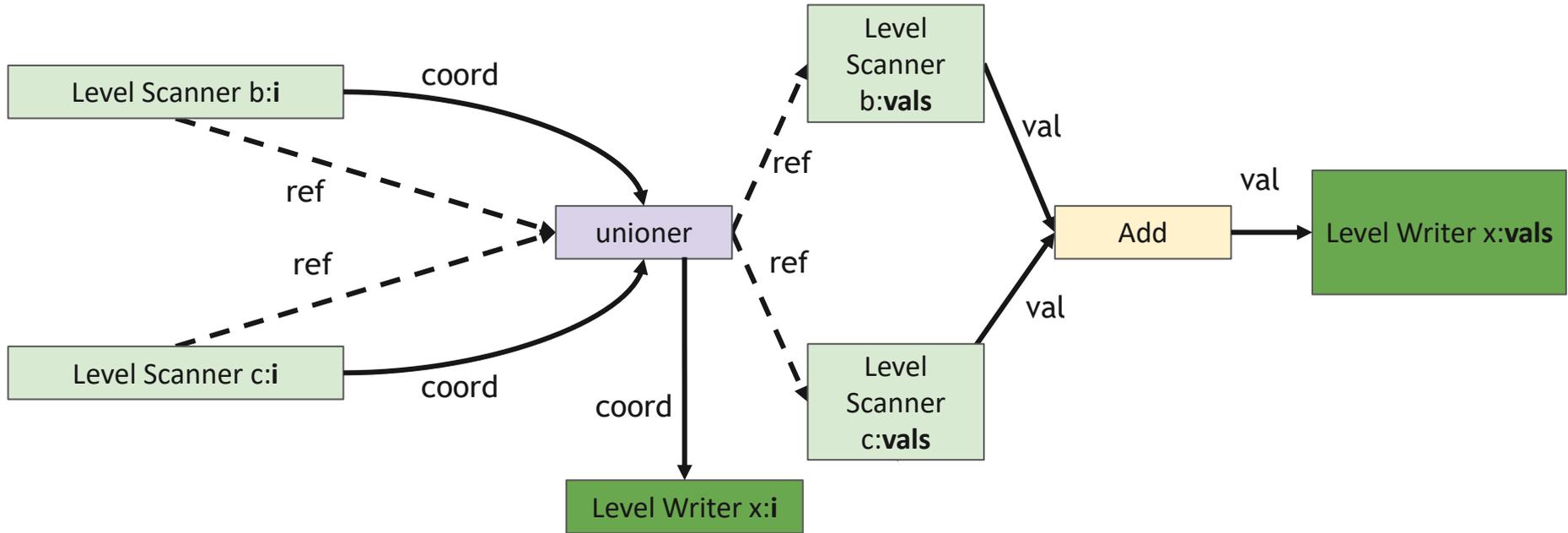
- The repeater broadcasts a stream over another



**coord** $S_1$ 2 1

**ref** $S_0$ 0

== S ?

$S_1$ 0 0 **ref**

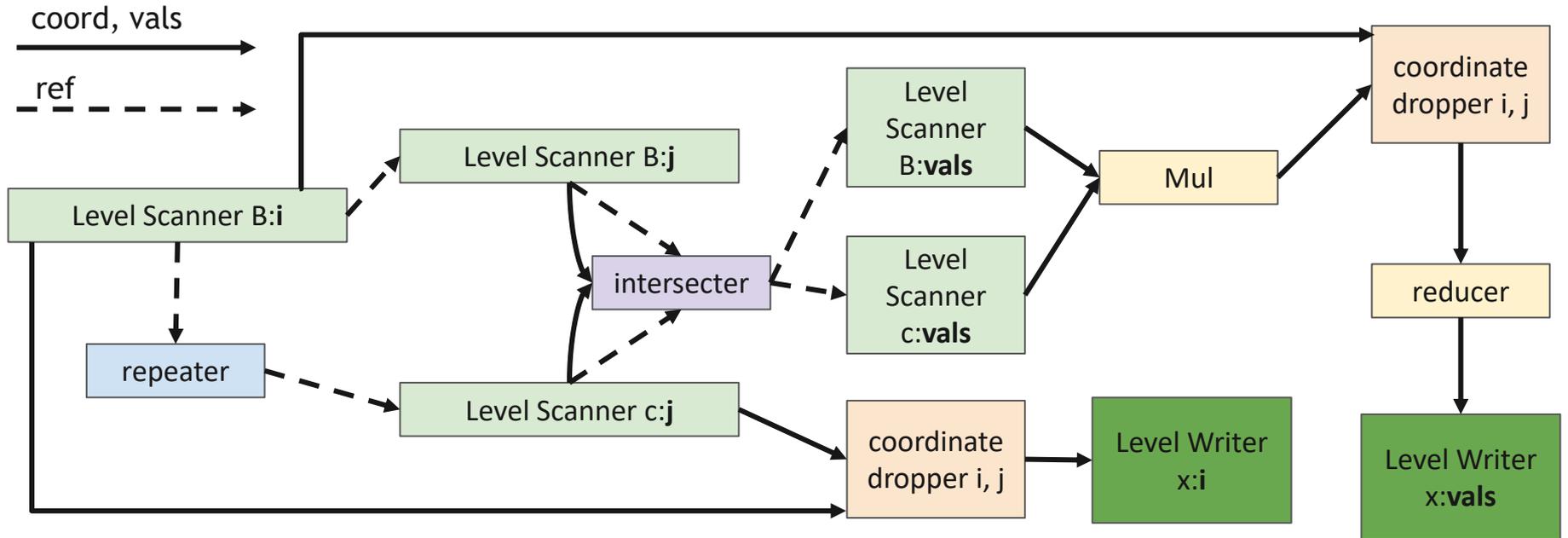Ex: Scalar Broadcasted over a Vector

5 * 0 3 4 0 → 3 * 5 and 4 * 5

# Composable Primitives

- With these primitives we support all of sparse tensor algebra
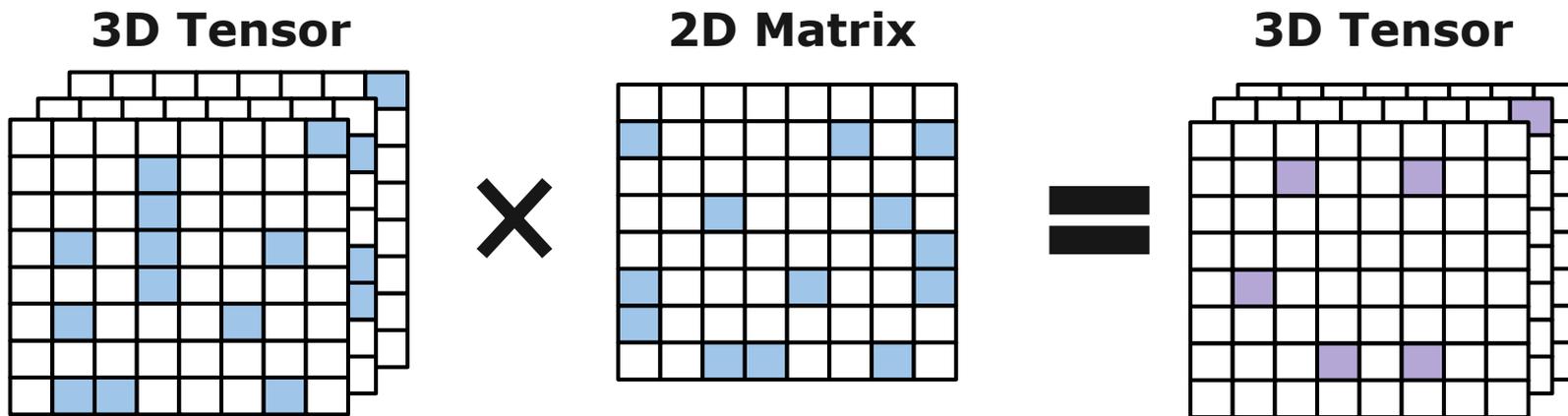
Vector-Vector Element Add $(X_i = b_i + c_i)$

# Composable Primitives

- With these primitives we support all of sparse tensor algebra
  Ex: Matrix-Vector Multiplication ($X_i = \sum_j B_{ij} c_j$)

coord, vals

ref

Level Scanner B:**i** → Level Scanner B:**j**

repeater

Level Scanner c:**j**

intersecter

Level Scanner B:**vals**

Level Scanner c:**vals**

Mul

coordinate dropper i, j

reducer

coordinate dropper i, j → Level Writer x:**i**
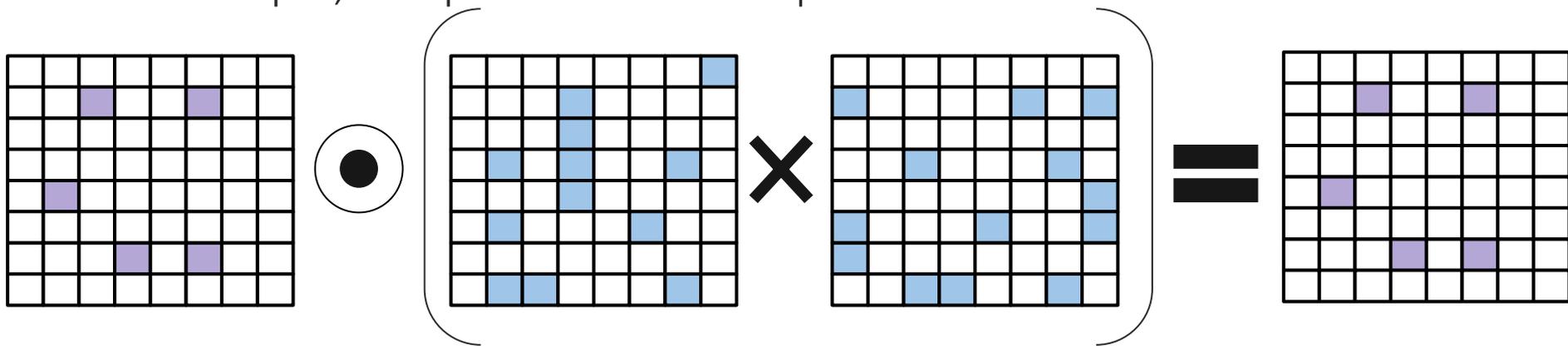
Level Writer x:**vals**

# Composable Primitives

- Utilizing these composable primitives, we support **high-dimensional tensors** and multiple inputs with fusion
- For example, tensor times matrix:
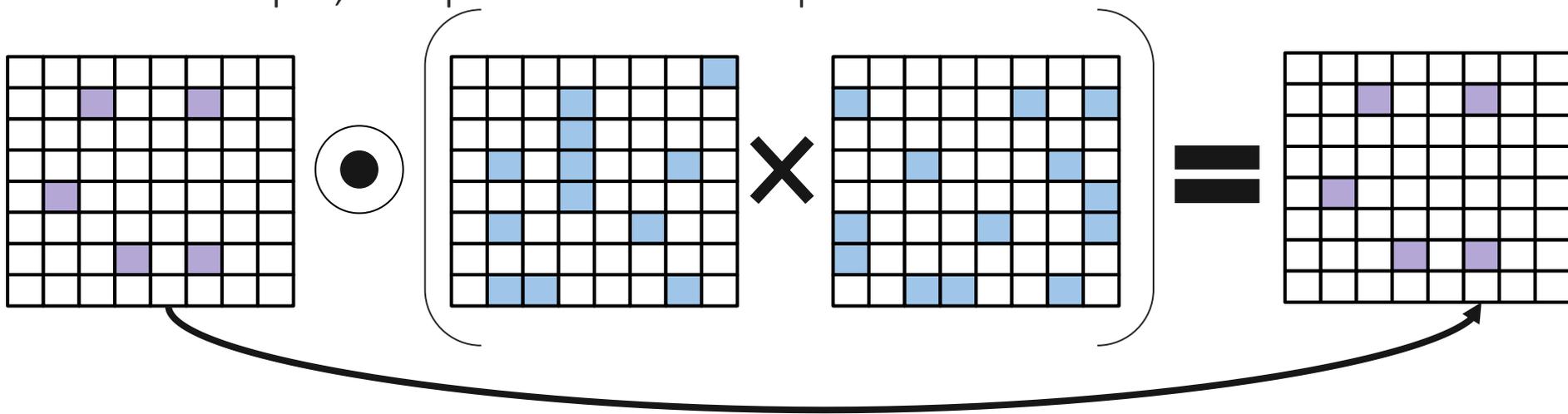
**3D Tensor**     **2D Matrix**     **3D Tensor**

# Composable Primitives

- Utilizing these composable primitives, we support high-dimensional tensors and **multiple inputs with fusion**
- For example, sampled matrix multiplication:

# Composable Primitives

- Utilizing these composable primitives, we support high-dimensional tensors and **multiple inputs with fusion**
- For example, sampled matrix multiplication:
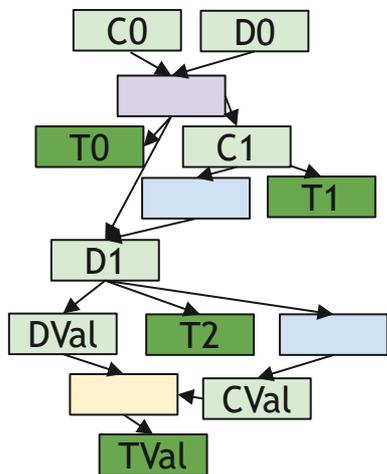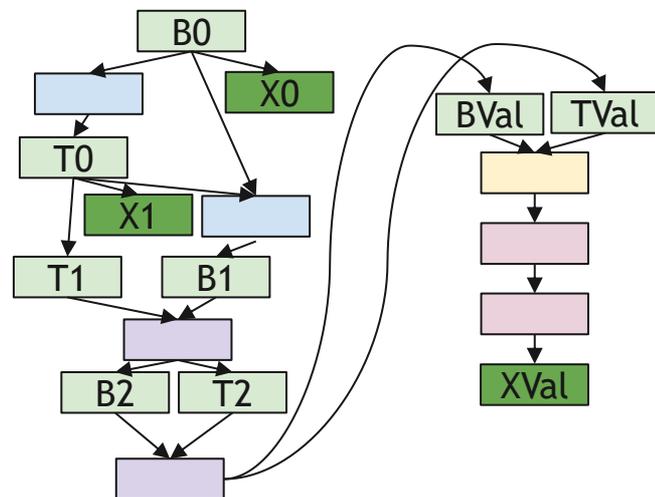


Only calculate necessary indices!

# Without Fusion

- For example, for a higher-dimensional, multi-input expression like **Matricized Tensor Times Khatri-Rao Product (**MTTKRP: $X_{ij}=\sum_{kl}B_{ikl}C_{jk}D_{jl}$)
- An accelerator with 2-input support splits MTTKRP into Kernel 1: $T_{jkl}=C_{jk}D_{jl}$ and Kernel 2: $X_{ij}=\sum_{kl}B_{ikl}T_{jkl}$



Kernel 1: $T_{jkl}=C_{jk}D_{jl}$

Kernel 2: $X_{ij}=\sum_{kl}B_{ikl}T_{jkl}$

Level Scanner
Level Writer
Intersect
Repeat
Mul
Reduce

# Without Fusion

- For example, for a higher-dimensional, multi-input expression like **Matricized Tensor Times Khatri-Rao Product (**MTTKRP: $X_{ij} = \sum_{kl} B_{ikl} C_{jk} D_{jl}$)
- An accelerator with 2-input support splits MTTKRP into Kernel 1: $T_{jkl} = C_{jk} D_{jl}$ and Kernel 2: $X_{ij} = \sum_{kl} B_{ikl} T_{jkl}$
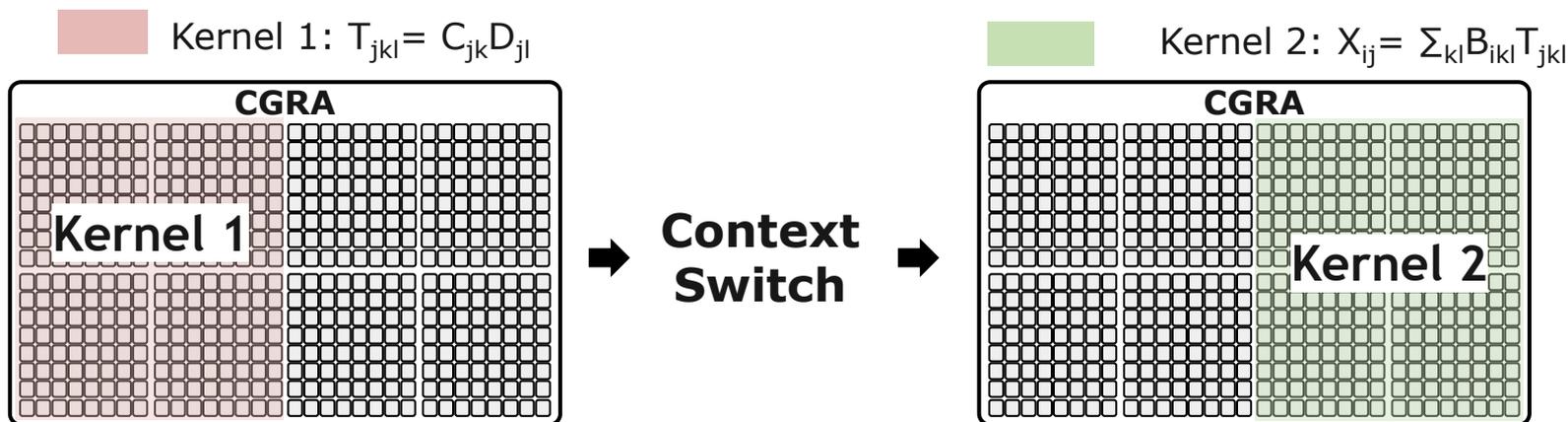


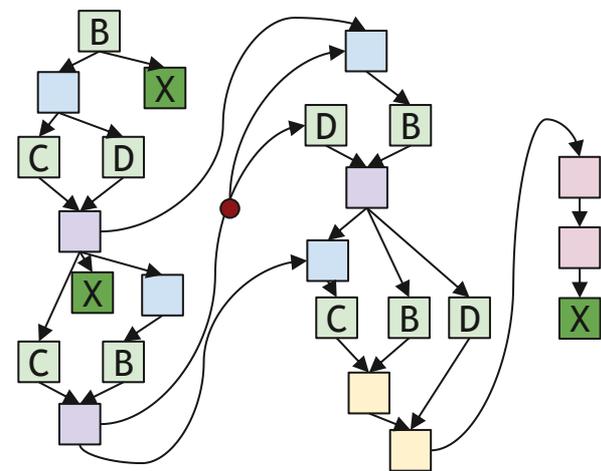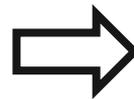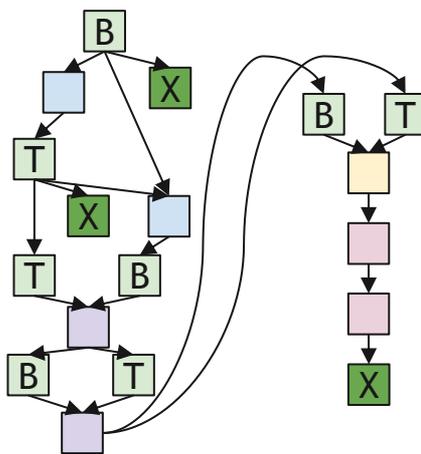Kernel 1: $T_{jkl} = C_{jk} D_{jl}$     Kernel 2: $X_{ij} = \sum_{kl} B_{ikl} T_{jkl}$
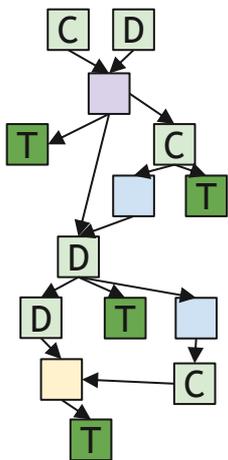
*Time*

# Fusion Results

- With Onyx we use a single **fused** kernel on the array



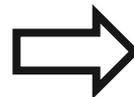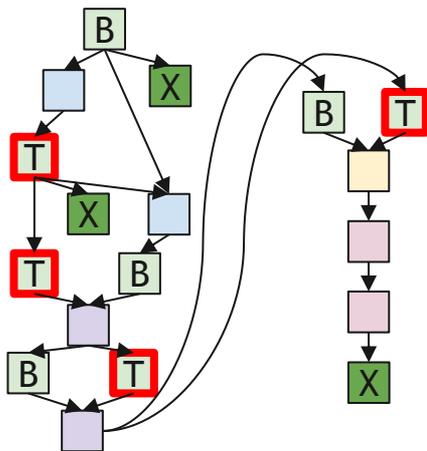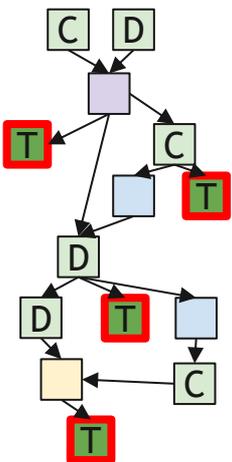Kernel 1: $T_{jkl} = C_{jk}D_{jl}$  Kernel 2: $X_{ij} = \Sigma_{kl}B_{ikl}T_{jkl}$  Fused Kernel : $X_{ij} = \Sigma_{kl}B_{ikl}C_{jk}D_{jl}$

Level Scanner
Level Writer
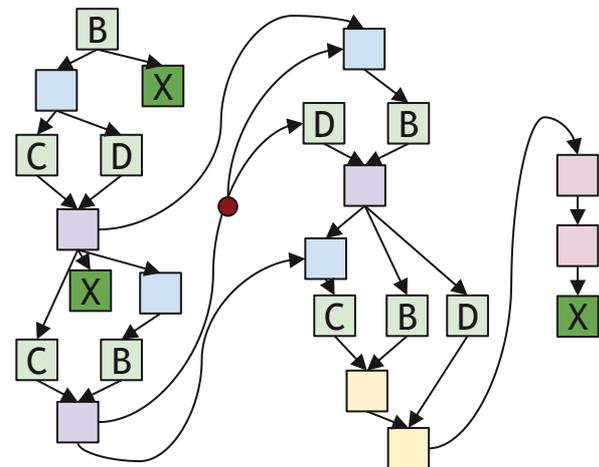Intersecter
Repeater
Multiply
Reducer

# Fusion Results

- With Onyx we use a single <u>fused</u> kernel on the array
- Eliminate intermediate storage of T



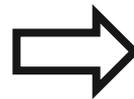Kernel 1: $T_{jkl} = C_{jk}D_{jl}$   Kernel 2: $X_{ij} = \Sigma_{kl}B_{ikl}T_{jkl}$   Fused Kernel : $X_{ij} = \Sigma_{kl}B_{ikl}C_{jk}D_{jl}$

Level Scanner
Level Writer
Intersecter
Repeater
Multiply
Reducer

# Fusion Results

- With Onyx we use a single <u>fused</u> kernel on the array
- Perform all intersections before compute



Kernel 1: $T_{jkl}= C_{jk}D_{jl}$   Kernel 2: $X_{ij}= \Sigma_{kl}B_{ikl}T_{jkl}$   Fused Kernel : $X_{ij}=\Sigma_{kl}B_{ikl}C_{jk}D_{jl}$

Legend:
- Level Scanner
- Level Writer
- Intersecter
- Repeater
- Multiply
- Reducer

# MTTKRP Fusion Results

- Onyx supports a single **<u>fused</u>** kernel on the array
- Avoids intermediate storage and eliminates unnecessary compute



Runtime (Cycles)

**6.6x Improvement**

Kernel 1 ■ Kernel 2 ■ Fused Kernel

# Contributions

## Sparse Acceleration Hardware

Composable primitives accelerating arbitrary sparse tensor algebra kernels



## Dense Acceleration Improvements

Compute and memory controller optimizations for dense applications



New PE OPs

PE Usage per App

Dynamic-/Static-GLB Clock Gating

# Challenge 1: Compute Density

- Compute only 37.8% of PE area
- Since each app requires several simple PEs, it is difficult to further unroll applications - low CGRA PE utilization

## Amber [2] PE Tile Area Breakdown
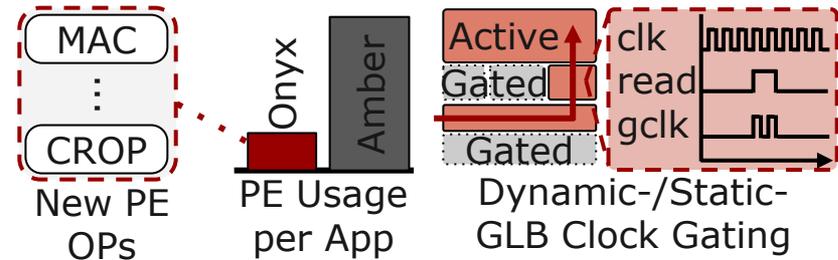


Config 2.1%
Interconnect 37.9%
ALU/FPU 37.8%
Register File 22.3%

**Further unrolling exceeds PE limit!**

## PE usage on Amber [2]



■ PE  — Total number of PEs

[2] K. Feng et al., JSSC'2023.

# Optimized Compute

- Utilized the Automated PE Exploration tool (APEX [3])
- Step 1: Mine frequent subgraphs from each application



**Camera Pipeline**     **Unsharp**     **Gaussian**     **ML**     **Harris**

[3] J. Melchert et al., ASPLOS'2023.

# Optimized Compute

- Step 1: Mine frequent subgraphs from each application
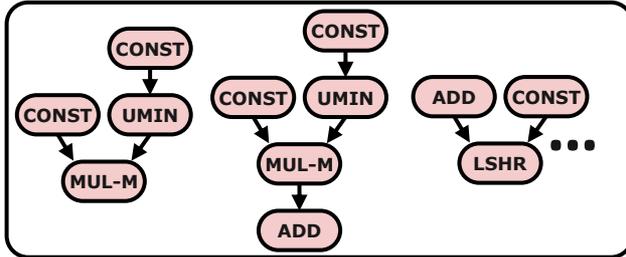- Step 2: Order subgraphs by # non-overlapping occurrences using maximal independent set analysis



**Camera Pipeline**

**Unsharp**

**Gaussian**

**ML**

**Harris**

High Occurrences    Low Occurrences

Operations in Baseline PE    Frequent subgraphs merged in Onyx PE    Frequent subgraphs not merged due to large area overhead

# Optimized Compute

- Step 3: Merge subgraphs into PE



**Baseline PE**

Regfile Output

PE Instruction Set:

INT16/BIT: ADC SBC ABS MUX MULT0 MULT1 MULT2 SHR SHL OR AND NOT XOR MIN MAX EQ CMP

BFloat16: ADD SUB CMP MUL GETMAN* ADDIEXP* SUBEXP* EXP2F* F2INT* GETFR* INT2F* (*Used in complex ops)

# Optimized Compute

- Step 3: Merge subgraphs into PE



PE Instruction Set:

INT16/BIT: ADC SBC ABS MUX MULT0
MULT1 MULT2 SHR SHL OR AND NOT XOR
MIN MAX EQ CMP

MAC variants (MULADD MULSUB) TADD
variants (ADDADD ADDSUB SUBADD
SUBSUB) MINMAX MULSHR

BFloat16: ADD SUB CMP MUL GETMAN*
ADDIEXP* SUBEXP* EXP2F* F2INT*
GETFR* INT2F* (*Used in complex ops)

# Optimized Compute

- Step 3: Merge subgraphs into PE



Onyx PE

PE Instruction Set:

INT16/BIT: ADC SBC ABS MUX MULT0 MULT1 MULT2 SHR SHL OR AND NOT XOR MIN MAX EQ CMP

MAC variants (MULADD MULSUB) TADD variants (ADDADD ADDSUB SUBADD SUBSUB) MINMAX MULSHR

BFloat16: ADD SUB CMP MUL GETMAN* ADDIEXP* SUBEXP* EXP2F* F2INT* GETFR* INT2F* (*Used in complex ops)

# Optimized Compute

- Step 3: Merge subgraphs into PE



**Onyx PE**

PE Instruction Set:

INT16/BIT: ADC SBC ABS MUX MULT0
MULT1 MULT2 SHR SHL OR AND NOT XOR
MIN MAX EQ CMP

MAC variants (MULADD MULSUB) TADD
variants (ADDADD ADDSUB SUBADD
SUBSUB) MINMAX MULSHR

BFloat16: ADD SUB CMP MUL GETMAN*
ADDIEXP* SUBEXP* EXP2F* F2INT*
GETFR* INT2F* (*Used in complex ops)

# Optimized Compute

- Step 3: Merge subgraphs into PE



PE Instruction Set:

INT16/BIT: ADC SBC ABS MUX MULT0 MULT1 MULT2 SHR SHL OR AND NOT XOR MIN MAX EQ CMP

MAC variants (MULADD MULSUB) TADD variants (ADDADD ADDSUB SUBADD SUBSUB) MINMAX MULSHR

BFloat16: ADD SUB CMP MUL GETMAN* ADDIEXP* SUBEXP* EXP2F* F2INT* GETFR* INT2F* (*Used in complex ops)

# Optimized Compute

- Step 3: Merge subgraphs into PE
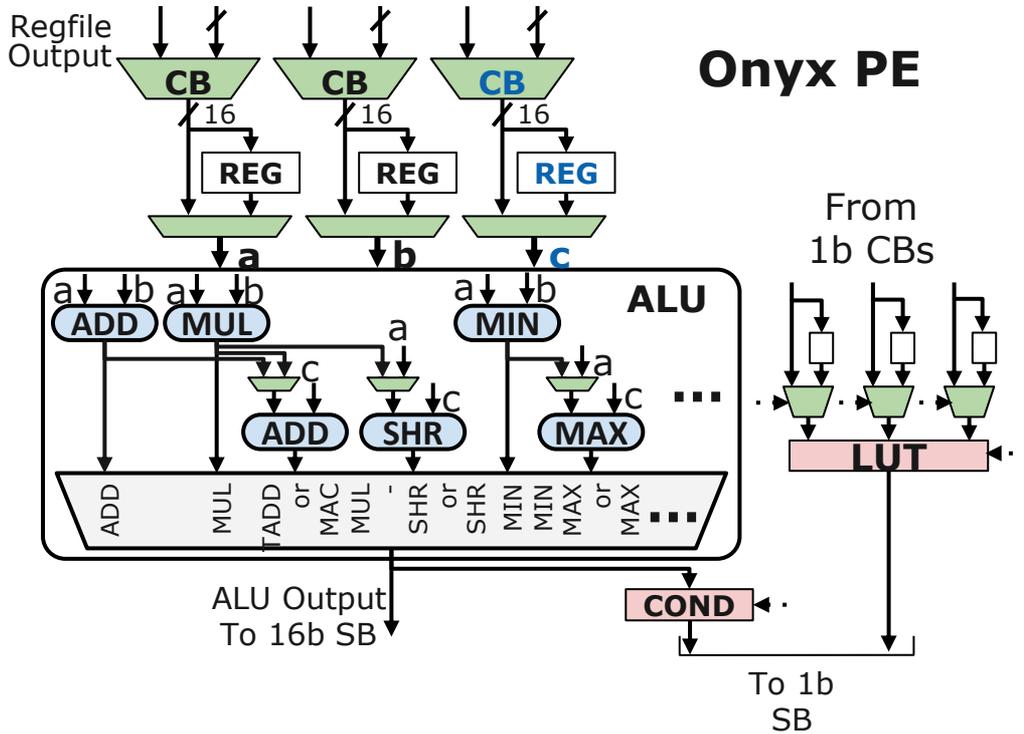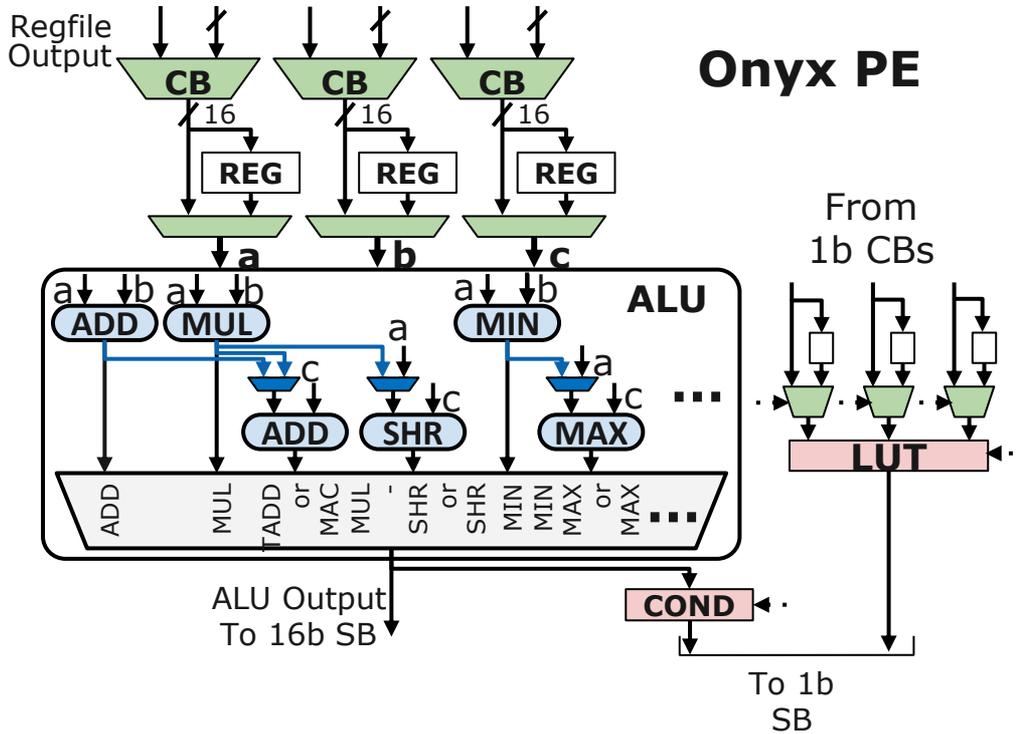


PE Instruction Set:

INT16/BIT: ADC SBC ABS MUX MULT0 MULT1 MULT2 SHR SHL OR AND NOT XOR MIN MAX EQ CMP

**MAC variants (MULADD MULSUB) TADD variants (ADDADD ADDSUB SUBADD SUBSUB) MINMAX MULSHR**

BFloat16: ADD SUB CMP MUL GETMAN* ADDIEXP* SUBEXP* EXP2F* F2INT* GETFR* INT2F* (*Used in complex ops)

# Optimized Compute Results

- Increased compute allows for further parallelization, reducing runtime



Iso-Frequency Runtime (ms/frame)

■ Amber ■ Onyx

50% — Harris
25% — Unsharp
12% — Blur
50% — ResNet Layer

# Challenge 2: Area Breakdown of MEM

- SRAMs only 32.1% of MEM Tile

# Dense Memory Controller

# Dense Memory Controller



AG Address Generator
SG Schedule Generator
ID Iteration Domain
SIPO Single-In Parallel-Out
PISO Parallel-In Single-Out

# Optimized Dense Memory Controller



❶ Simplified write controllers by replacing address generators (AGs) with delay blocks for read-modify-write operations

# Optimized Dense Memory Controller



❶ Simplified write controllers by replacing address generators (AGs) with delay blocks for read-modify-write operations

❷ Halved SIPO FIFO depth

# Optimized Dense Memory Controller



❶ Simplified write controllers by replacing address generators (AGs) with delay blocks for read-modify-write operations

❷ Halved SIPO FIFO depth

❸ Reduced counter and configuration register bit widths in iteration domains (IDs)

# Dense Memory Controller Results

- Memory controller and configuration reduced by 24%



[2] K. Feng et al., JSSC'2023.

# Chip Specifications and Die Photo



| | Onyx |
|---|---|
| Architecture | SoC with CGRA |
| Node | GF 12 nm |
| Area (mm$^2$) | 23 |
| Precision | BF16, INT16 |
| SRAM (MiB) | 4.5 |
| Voltage (V) | 0.78 |
| Freq (MHz) | 970 |
| Peak GOPS | 571 |
| GOPS/W | 756 |

# Application Compiler

- Designed end-end compiler to map Sparse Tensor Algebra Expressions onto CGRA

*Sparse Tensor Algebra Expression*

Compute (stmt):

```
X(i,j) = B(i,k) * C(k,j);
```

Scheduling and Format:

```
Format sp = Sparse;
Tensor X({sp,sp});
Tensor B({sp,sp});
Tensor C({sp,sp},{1,0});
stmt.reorder({i,j,k})
```

# Application Compiler

- Designed end-end compiler to map Sparse Tensor Algebra Expressions onto CGRA



Sparse Tensor Algebra Expression → Compiler: *Custard* [4] → Sparse Dataflow Graph (*SAM*)

Compute (stmt):

```
X(i,j) = B(i,k) * C(k,j);
```

Scheduling and Format:

```
Format sp = Sparse;
Tensor X({sp,sp});
Tensor B({sp,sp});
Tensor C({sp,sp},{1,0});
stmt.reorder({i,j,k})
```

LS: Level Scanner
LW: Level Writer
R: Repeater

→ Reference
→ Coordinate
→ Value

[4] O. Hsu et al., ASPLOS'2023.

Stanford University

# Application Compiler

- Designed end-end compiler to map Sparse Tensor Algebra Expressions onto CGRA

[4] O. Hsu et al., ASPLOS'2023.

# Application Compiler

- Designed end-end compiler to map Sparse Tensor Algebra Expressions onto CGRA



Compute (stmt):

```
X(i,j) = B(i,k) * C(k,j);
```

Scheduling and Format:

```
Format sp = Sparse;
Tensor X({sp,sp});
Tensor B({sp,sp});
Tensor C({sp,sp},{1,0});
stmt.reorder({i,j,k})
```

LS: Level Scanner
LW: Level Writer
R: Repeater

→ Reference
→ Coordinate
→ Value

[4] O. Hsu et al., ASPLOS'2023.

# Benchmark Application Suite

Dense applications written in Halide compared against CPU, GPU, FPGA, and CGRA:

- Image processing
  - Blur: image blur
  - Unsharp: enhances local contrast by smoothing an image
  - Camera pipeline: processes raw data from an image sensor into a color image

- Computer vision
  - Harris: detects corners

- Machine learning
  - ResNet-18: image classification

# Benchmark Application Suite

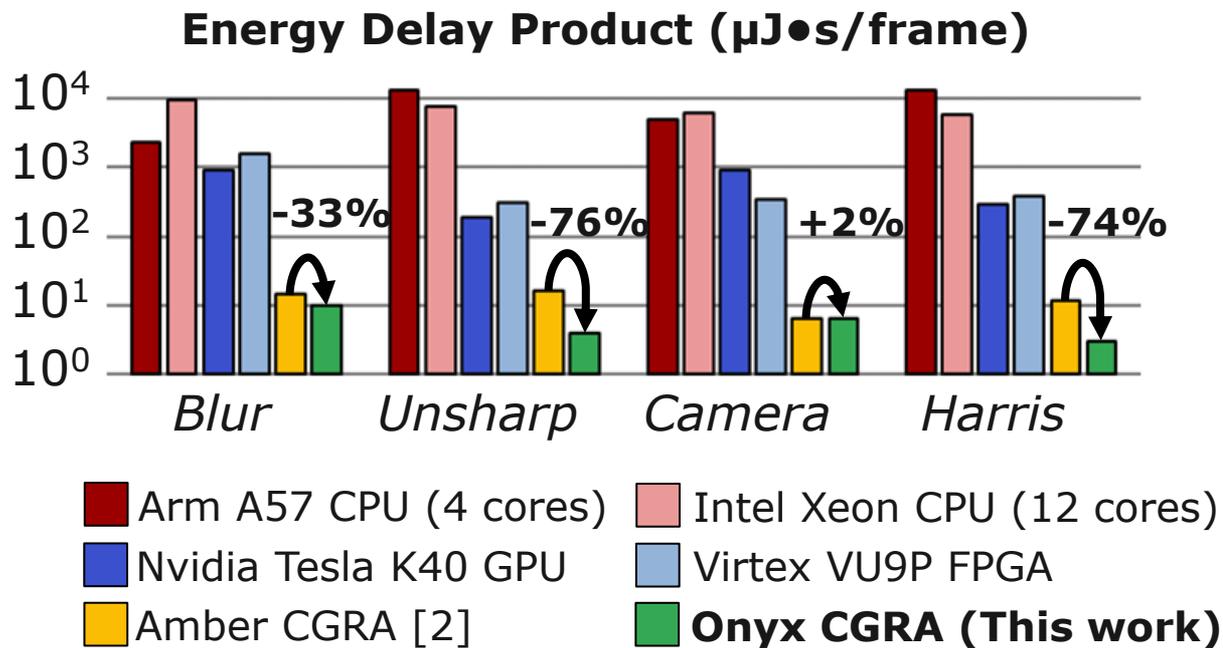Sparse tensor algebra expressions compiled using Custard compared against CPU:

- 2D Dataset: Randomly selected SuiteSparse Matrices from linear programming, computational fluid dynamics, chemical process simulation, undirected weighted random graphs, etc.
  - Matrix Dimensions: 821x1876 - 2021x2021

- 3D Dataset: Uniform random generated sparse (67%) tensors
  - Tensors Dimensions: 8x37x10 - 28x35x54

# Dense Image Processing Results

- Up to 76% better EDP versus the state-of-the-art CGRA
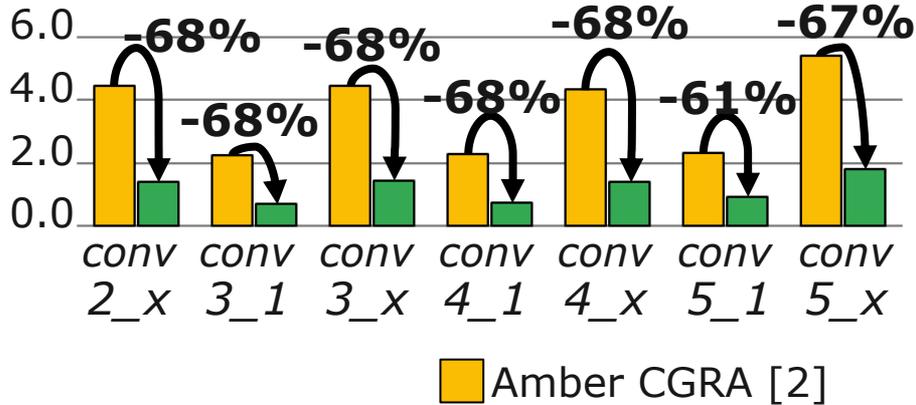
**Energy Delay Product (μJ•s/frame)**



Legend:
- Arm A57 CPU (4 cores)
- Intel Xeon CPU (12 cores)
- Nvidia Tesla K40 GPU
- Virtex VU9P FPGA
- Amber CGRA [2]
- **Onyx CGRA (This work)**
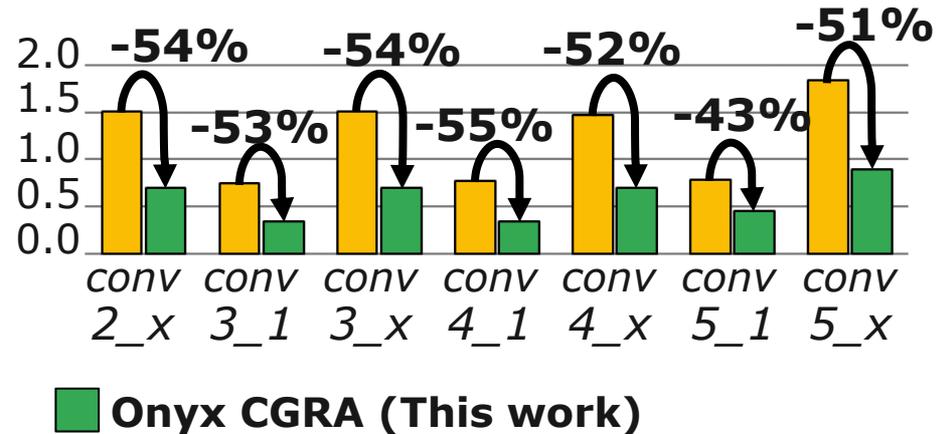
[2] K. Feng et al., JSSC'2023.

# Dense Image ResNet Results

- Up to 55% better energy versus the state-of-the-art CGRA

**ResNet-18 Layers Runtime (ms)**
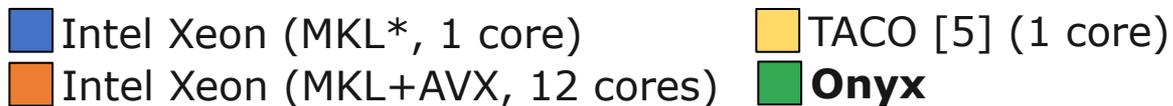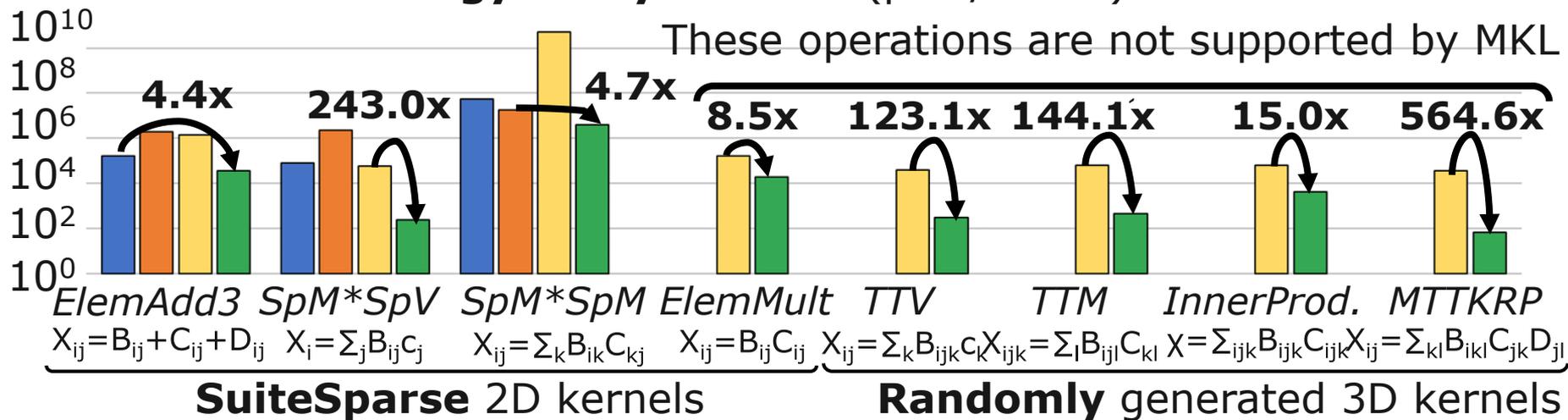
**ResNet-18 Layers Energy (mJ)**



■ Amber CGRA [2]  ■ **Onyx CGRA (This work)**

[2] K. Feng et al., JSSC'2023.

# Sparse Tensor Algebra Results

- Up to 564.6x better EDP versus a CPU

**Energy Delay Product** (pJ•s/frame)

These operations are not supported by MKL



**4.4x** **243.0x** **4.7x** **8.5x** **123.1x** **144.1x** **15.0x** **564.6x**

| ElemAdd3 | SpM*SpV | SpM*SpM | ElemMult | TTV | TTM | InnerProd. | MTTKRP |

$X_{ij}=B_{ij}+C_{ij}+D_{ij}$  $X_i=\sum_j B_{ij}c_j$  $X_{ij}=\sum_k B_{ik}C_{kj}$  $X_{ij}=B_{ij}C_{ij}$  $X_{ij}=\sum_k B_{ijk}c_k$  $X_{ijk}=\sum_l B_{ijl}C_{kl}$  $X=\sum_{ijk}B_{ijk}C_{ijk}$  $X_{ij}=\sum_{kl}B_{ikl}C_{jk}D_{jl}$

**SuiteSparse** 2D kernels          **Randomly** generated 3D kernels

- ▮ Intel Xeon (MKL*, 1 core)          ▮ TACO [5] (1 core)
- ▮ Intel Xeon (MKL+AVX, 12 cores)          ▮ **Onyx**

[5] F. Kjolstad et al., OOPSLA'2017.
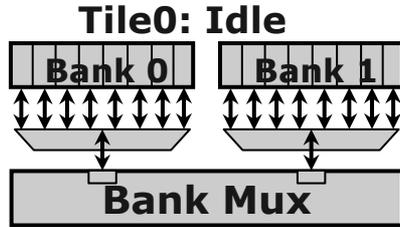
# Summary of Key Contributions

- Onyx is an SoC designed for flexible and efficient acceleration of sparse and dense data applications

  - **Sparse**: composable primitives for arbitrary sparse tensor algebra supporting multi-dimensional, multi-input expressions with fusion
  - **Dense**: compute density increased, and memory controllers optimized

- Automatic end-to-end compiler maps applications onto Onyx

- Onyx achieves up to 565x EDP improvement over CPUs with sparse libraries and 85% lower EDP versus the state-of-the-art CGRA on dense applications

- Demonstrates an approach for accelerating fast evolving application domains
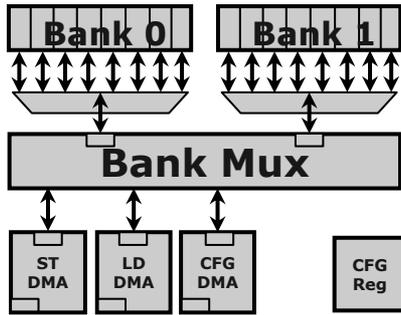
# Backup Slides

# GLB Architecture

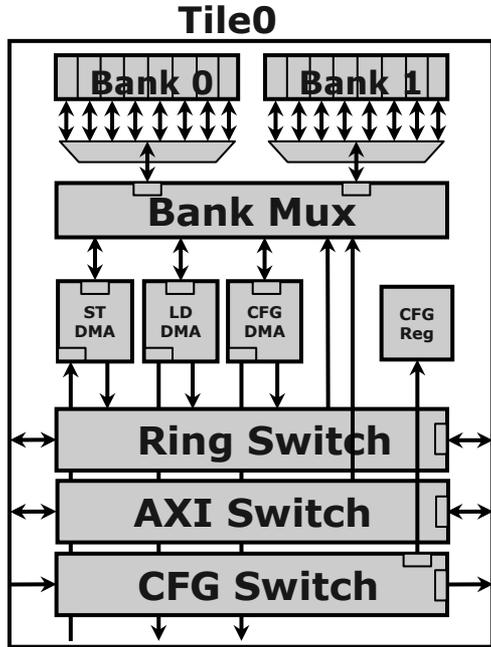- Each GLB tile is composed of two banks

# GLB Architecture

- Has Store, Load, Configuration DMAs and Configuration registers
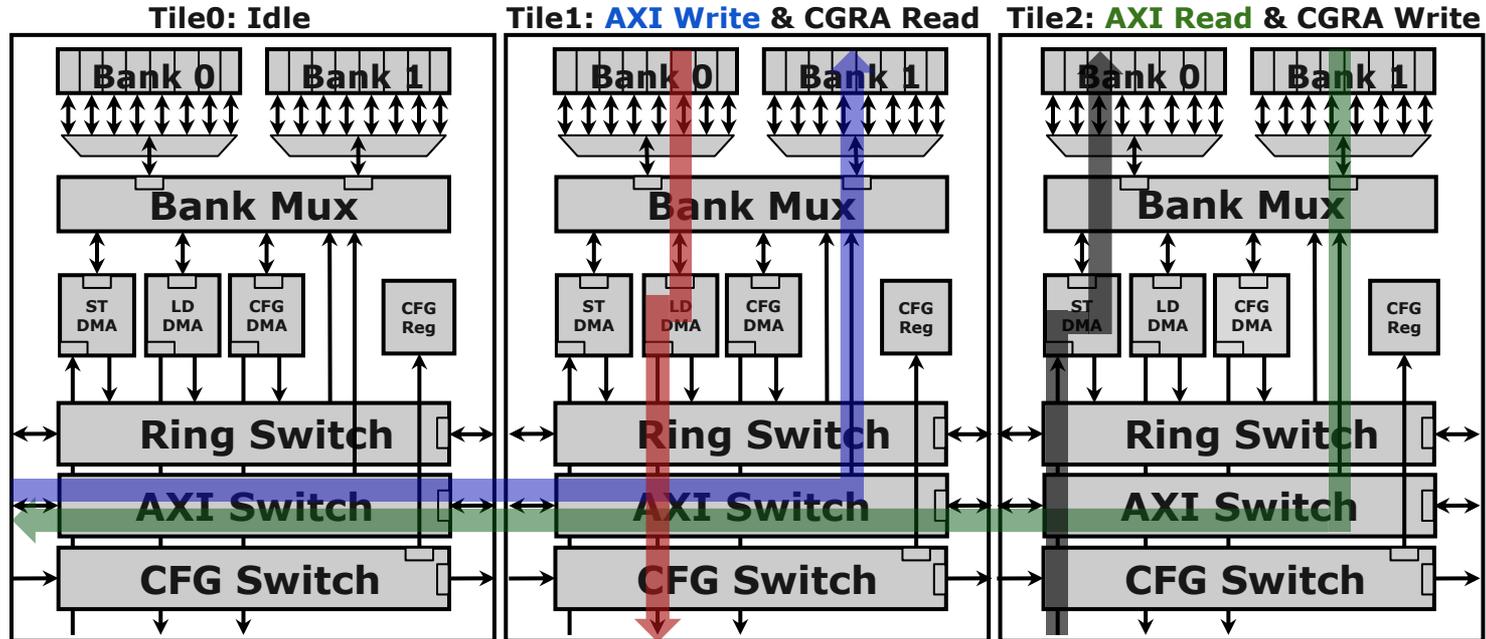
Stanford University

# GLB Architecture

- Switches for CGRA and AXI communication and configuration
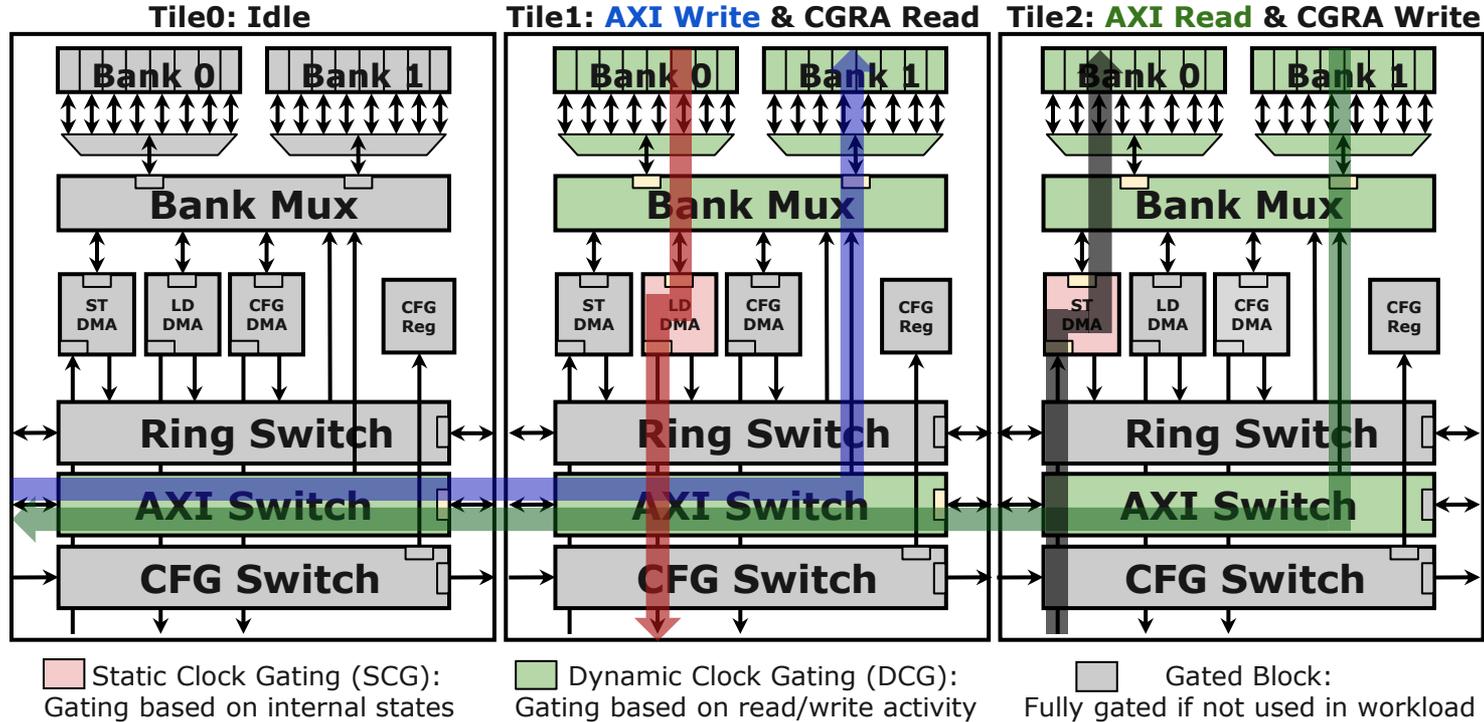
# GLB Dynamic-/Static- Clock Gating

- Double buffering workload: what blocks are active?

# GLB Dynamic-/Static- Clock Gating

- Only parts of the blocks are active during read/write transactions

# GLB Results

- 24% power saved on example double buffering workload