# Custard and The Sparse Abstract Machine: Compiling Sparse Applications to Coarse-Grained Reconfigurable Arrays

Olivia Hsu

Stanford University

# Dataflow hardware can accelerate sparse tensor algebra

# Dataflow hardware can accelerate sparse tensor algebra

$$a = Bc + a \qquad a = Bc \qquad A = B + C$$

$$a = B^T c + d$$
$$\qquad\qquad a = B^T c \qquad A = \alpha B \qquad a = Bc + b$$
$$\qquad\qquad\qquad a = b \odot c \qquad a = B(c + d)$$

$$A = B + C + D \qquad A = BC$$
$$A = B \odot (CD)$$

$$A = B \odot C \qquad A = 0 \qquad A = BCd \qquad A = B^T$$
$$\qquad\qquad\qquad\qquad\qquad\qquad\qquad a = B^T Bc$$

$$a = b + c \qquad A = B \qquad K = A^T C A \qquad a = \alpha Bc + \beta a$$

$$A_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj} \qquad A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij}$$

$$A_{lj} = \sum_{ik} B_{ikl} C_{ij} D_{kj} \qquad A_{ij} = \sum_{k} B_{ijk} c_k$$

$$A_{ijk} = \sum_{l} B_{ikl} C_{lj} \qquad A_{ik} = \sum_{j} B_{ijk} c_j$$

$$A_{jk} = \sum_{i} B_{ijk} c_i \qquad A_{ijl} = \sum_{k} B_{ikl} C_{kj}$$

$$C = \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \, \overline{P_{il}} \qquad \tau = \sum_{i} z_i \left( \sum_{j} z_j \theta_{ij} \right) \left( \sum_{k} z_k \theta_{ik} \right)$$

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

**Need Generality to Handle This…**

# Dataflow hardware can accelerate sparse tensor algebra



**Need Generality to Handle This...**

# Dataflow hardware can accelerate sparse tensor algebra



Need Generality to Handle This…

# Dataflow hardware can accelerate sparse tensor algebra



Need Generality to Handle This…

Onyx CGRA

[Koul et al. VLSI, HotChips 2024]

but really any sparse accelerator…

# Performance requires generality in schedules

# Performance requires generality in schedules

Fusion

Dataflow Ordering

# Performance requires generality in schedules

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

SDDMM

Fusion

Dataflow Ordering

# Performance requires generality in schedules

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

SDDMM



Fusion                                          Dataflow Ordering

# Performance requires generality in schedules



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

SDDMM

Cycles

$10^5$

$10^4$

Fusion

Dataflow Ordering

# Performance requires generality in schedules

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

SDDMM



Unfused

Cycles

$10^5$

$10^4$

Fusion

Dataflow Ordering

# Performance requires generality in schedules

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

SDDMM



Fusion

Dataflow Ordering

# Performance requires generality in schedules

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

SDDMM



Cycles

$10^5$

$10^4$

Unfused

Fused

Fusion

Dataflow Ordering

# Performance requires generality in schedules

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

SDDMM



$$X_{ij} = B_{ik} \cdot C_{kj}$$

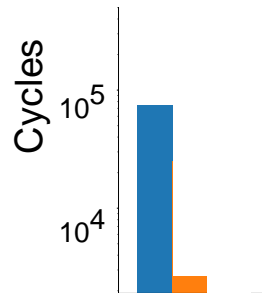SpMSpM



Fusion



jki    kij    kji

Dataflow Ordering

# Performance requires generality in schedules



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

SDDMM

$$X_{ij} = B_{ik} \cdot C_{kj}$$

SpMSpM

Fusion

Dataflow Ordering

# Efficient mapping requires a compiler

# Efficient mapping requires a compiler

$a = Bc + a$  $a = Bc$  $A = B + C$  Linear Algebra

$a = B^T c + d$  $A = \alpha B$  $a = Bc + b$

$a = B^T c$  $a = b \odot c$  $a = B(c + d)$

$A = B + C + D$  $A = BC$  $A = B \odot (CD)$

$A = B \odot C$  $A = 0$  $A = BCd$  $A = B^T$  $a = B^T Bc$

$a = b + c$  $A = B$  $K = A^T C A$  $a = \alpha Bc + \beta a$

$$A_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj} \quad A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij}$$

$$A_{lj} = \sum_{ik} B_{ikl} C_{ij} D_{kj} \quad A_{ij} = \sum_{k} B_{ijk} c_k$$

$$A_{ijk} = \sum_{l} B_{ikl} C_{lj} \quad A_{ik} = \sum_{j} B_{ijk} c_j$$

Data analytics (tensor factorization)

$$A_{jk} = \sum_{i} B_{ijk} c_i \quad A_{ijl} = \sum_{k} B_{ikl} C_{kj}$$

$$C = \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \, \overline{P_{il}} \quad \tau = \sum_{i} z_i (\sum_{j} z_j \theta_{ij})(\sum_{k} z_k \theta_{ik})$$

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

Quantum Chromodynamics

Algorithm (expression)

# Efficient mapping requires a compiler



Algorithm (expression)

Schedule

# Efficient mapping requires a compiler



Linear Algebra

$a = Bc + a$   $a = Bc$   $A = B + C$
$a = B^T c + d$   $a = Bc + b$
$a = B^T c$   $A = \alpha B$   $a = Bc + b$
$A = B + C + D$   $A = BC$   $a = B(c + d)$
$A = B \odot C$   $A = 0$   $A = BCd$   $A = B^T$
$a = b + c$   $A = B$   $K = A^T C A$   $a = \alpha Bc + \beta a$

$A_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj}$   $A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij}$

$A_{lj} = \sum_{ik} B_{ikl} C_{ij} D_{kj}$   $A_{ij} = \sum_{k} B_{ijk} c_k$

$A_{ijk} = \sum_{l} B_{ikl} C_{lj}$   $A_{ik} = \sum_{j} B_{ijk} c_j$

Data analytics (tensor factorization)

$A_{jk} = \sum_{i} B_{ijk} c_i$   $A_{ijl} = \sum_{k} B_{ikl} C_{kj}$

$C = \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \, \overline{P_{il}}$   $\tau = \sum_{i} z_i (\sum_{j} z_j \theta_{ij})(\sum_{k} z_k \theta_{ik})$

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

Quantum Chromodynamics

**Algorithm (expression)**

$\times$

reorder

precompute

parallelize split

map   divide

vectorize unroll

position

**Schedule**

# Efficient mapping requires a compiler

$a = Bc + a$

$a = Bc$    $A = B + C$    Linear Algebra

$a = B^T c + d$    $a = Bc + b$

$a = B^T c$    $A = \alpha B$

$A = B + C + D$    $A = BC$    $a = b \odot c$    $a = B(c + d)$

$A = B \odot C$    $A = 0$    $A = BCd$    $A = B^T$    $A = B \odot (CD)$

$a = b + c$    $A = B$    $K = A^T CA$    $a = \alpha Bc + \beta a$    $a = B^T Bc$

$$A_{ij} = \sum_{kl} B_{ikl}C_{lj}D_{kj} \quad A_{kj} = \sum_{il} B_{ikl}C_{lj}D_{ij}$$

$$A_{lj} = \sum_{ik} B_{ikl}C_{ij}D_{kj} \quad A_{ij} = \sum_{k} B_{ijk}c_k$$

$$A_{ijk} = \sum_{l} B_{ikl}C_{lj} \quad A_{ik} = \sum_{j} B_{ijk}c_j$$

Data analytics (tensor factorization)

$$A_{jk} = \sum_{i} B_{ijk}c_i \quad A_{ijl} = \sum_{k} B_{ikl}C_{kj}$$

$$C = \sum_{ijkl} M_{ij}P_{jk}\overline{M_{lk}}\,\overline{P_{il}} \quad \tau = \sum_{i} z_i(\sum_{j} z_j \theta_{ij})(\sum_{k} z_k \theta_{ik})$$

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

Quantum Chromodynamics

**Algorithm (expression)**

$\times$

reorder

precompute

parallelize  split

map   divide

vectorize  unroll

position

$\Rightarrow$  Onyx

**Schedule**

# Efficient mapping requires a compiler



**Algorithm (expression)**

Linear Algebra

$a = Bc + a$  $a = Bc$  $A = B + C$

$a = B^T c + d$  $a = \alpha B$  $a = Bc + b$

$a = B^T c$  $A = \alpha B$  $a = B(c + d)$

$A = B + C + D$  $A = BC$  $a = b \odot c$  $A = B \odot (CD)$

$A = B \odot C$  $A = 0$  $A = BCd$  $A = B^T$  $a = B^T Bc$

$a = b + c$  $A = B$  $K = A^T CA$  $a = \alpha Bc + \beta a$

$A_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj}$  $A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij}$

$A_{lj} = \sum_{ik} B_{ikl} C_{ij} D_{kj}$  $A_{ij} = \sum_{k} B_{ijk} c_k$

$A_{ijk} = \sum_{l} B_{ikl} C_{lj}$  $A_{ik} = \sum_{j} B_{ijk} c_j$

Data analytics (tensor factorization)

$A_{jk} = \sum_{i} B_{ijk} c_i$  $A_{ijl} = \sum_{k} B_{ikl} C_{kj}$

$C = \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \overline{P_{il}}$  $\tau = \sum_{i} z_i (\sum_{j} z_j \theta_{ij})(\sum_{k} z_k \theta_{ik})$

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

Quantum Chromodynamics

$\times$

**Schedule**

reorder

precompute

parallelize split

map   divide

vectorize unroll

position

**Backend**

EIE

SpArch   OuterSPACE

Gamma   MatRaptor

Plasticine

Sigma

Tensaurus

Onyx   Sparse TPU

SCNN

Eyeriss V2

SpaceA

UCNN

Spada

Fifer   Sparse CGRAs

SPU

ExTensor

Capstan

… and **future**

# Efficient mapping requires a compiler

All of Sparse Tensor Algebra

$a = Bc + a$    $a = Bc$    $A = $ ... Algebra

$a = B^T c + d$

$a = B^T c$    $A = $ ... $+ b$

$a = B(c + d)$

$A = B + C$    ... $\odot c$    $A = B \odot (CD)$

$A = B \odot$    $A = BCd$    $A = B^T$    $a = B^T Bc$

$= B$    $K = A^T CA$    $a = \alpha Bc + \beta a$

$_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj}$    $A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij}$

$A_{lj} = \sum_{ik} B_{ikl} C_{ij} D_{kj}$    $A_{ij} = \sum_{k} B_{ijk} c_k$

$A_{ijk} = \sum_{l} B_{ikl} C_{lj}$    $A_{ik} = \sum_{j} B_{ijk} c_j$

Data analytics (tensor factorization)

$A_{jk} = \sum_{i} B_{ijk} c_i$    $A_{ijl} = \sum_{k} B_{ikl} C_{kj}$

$C = \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \, \overline{P_{il}}$    $\tau = \sum_{i} z_i (\sum_{j} z_j \theta_{ij})(\sum_{k} z_k \theta_{ik})$

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

Quantum Chromodynamics

**Algorithm (expression)**

$\times$

reorder

precompute

parallelize  split

map  divide

vectorize  unroll

position

**Schedule**

EIE

SpArch    OuterSPACE

Gamma    MatRaptor

Plasticine

Sigma

Tensaurus

Onyx    Sparse TPU

SCNN

Eyeriss V2

SpaceA

UCNN

Spada

Fifer    Sparse CGRAs

SPU

ExTensor

Capstan

… and **future**

**Backend**

# Efficient mapping requires a compiler

a = Bc + a   a = Bc   A = ... Algebra

$a = B^T c + d$   ... + b

$a = B^T c$ ...   $a = B(c + d)$

$A = B + C$ ...   ... $\odot c$   $A = B \odot (CD)$

$A = B \odot$ ...   $A = BCd$   $A = B^T$

... $= B$   $K = A^T C A$   $a = B^T Bc$   $a = \alpha Bc + \beta a$

$..._{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj}$   $A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij}$

$A_{lj} = \sum_{ik} B_{ikl} C_{ij} D_{kj}$   $A_{ij} = \sum_{k} B_{ijk} c_k$

$A_{ijk} = \sum_{l} B_{ikl} C_{lj}$   $A_{ik} = \sum_{j} B_{ijk} c_j$

$A_{jk} = \sum_{i} B_{ijk} c_i$   $A_{ijl} = \sum_{k} B_{ikl} C_{kj}$

Data analytics (tensor factorization)

$C = \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \, \overline{P_{il}}$   $\tau = \sum_{i} z_i (\sum_{j} z_j \theta_{ij})(\sum_{k} z_k \theta_{ik})$

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

Quantum Chromodynamics

**All of Sparse Tensor Algebra**

**Algorithm (expression)**

$\times$

reorder

precompute

parallelize  split

map   divide

vectorize  unroll

position

**Schedule**

**Lends itself to multiple backends**

...CE

MatRaptor

...cine   Sigma

Tensaurus

Onyx   Sparse TPU

SCNN

Eyeriss V2

SpaceA

UCNN   Spada

Fifer   Sparse CGRAs

SPU   ExTensor

Capstan

… and **future**

**Backend**

# This work in the DSL-based CGRA flow



Custard Application → Compiler: Custard → Sparse Abstract Machine → Sparse Mapper → Dataflow Graph of SAM Primitives → Placer & Router → CGRA Bit-Stream

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

# This work in the DSL-based CGRA flow

[Hsu et al. ASPLOS 2023]

*Custard* Application → Compiler: *Custard* → Sparse Abstract Machine → Sparse Mapper → Dataflow Graph of *SAM* Primitives → Placer & Router → CGRA Bit-Stream

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

# SAM leverages domain-specific languages and comes with the Custard compiler

Custard Application

Custard Application → Compiler: Custard → SAM Graph → Sparse Mapper → Dataflow Graph of SAM Primitives → Placer & Router → CGRA Bit-Stream
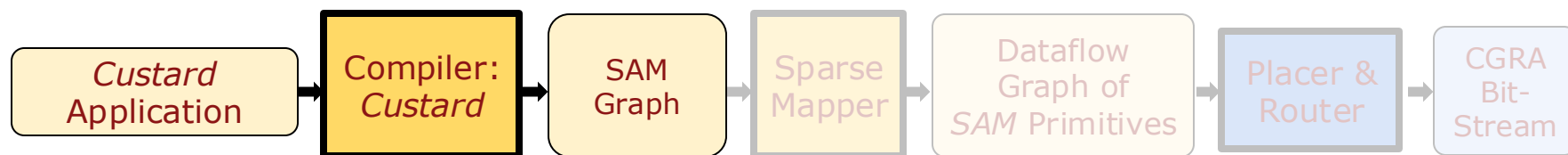
# SAM leverages domain-specific languages and comes with the Custard compiler

*Custard* Application

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

Tensor Index Notation

*Custard* Application → Compiler: *Custard* → SAM Graph → Sparse Mapper → Dataflow Graph of *SAM* Primitives → Placer & Router → CGRA Bit-Stream

# SAM leverages domain-specific languages and comes with the Custard compiler

Custard
Application

```
Format sp = Sparse;
Tensor X({sp,sp});
Tensor B({sp,sp});
Tensor C({sp,sp},{1,0});
```

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

Formats

Tensor Index
Notation

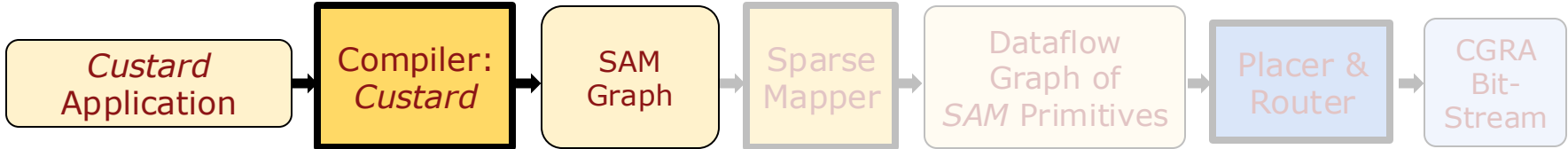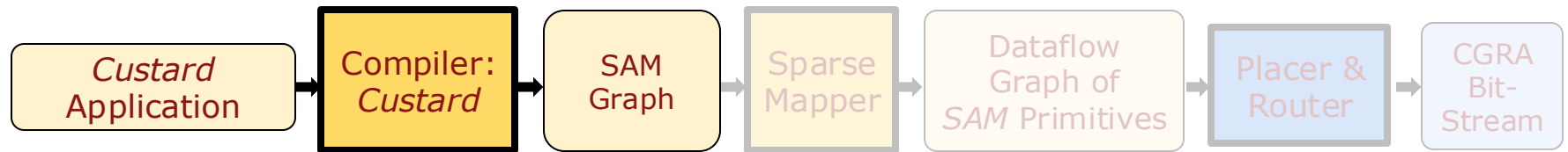| Custard Application | → | Compiler: Custard | → | SAM Graph | → | Sparse Mapper | → | Dataflow Graph of SAM Primitives | → | Placer & Router | → | CGRA Bit-Stream |

# SAM leverages domain-specific languages and comes with the Custard compiler



Custard Application

```
Format sp = Sparse;
Tensor X({sp,sp});
Tensor B({sp,sp});
Tensor C({sp,sp},{1,0});
```

`X(i,j) = B(i,k) * C(k,j);  stmt.reorder({i,j,k});`

Formats     Tensor Index Notation     Schedule

Custard Application → Compiler: Custard → SAM Graph → Sparse Mapper → Dataflow Graph of SAM Primitives → Placer & Router → CGRA Bit-Stream

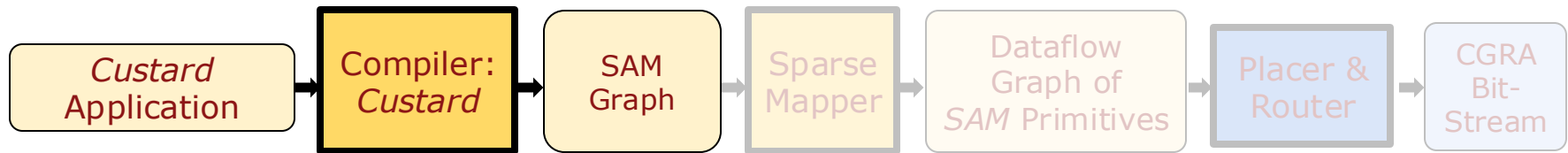# SAM leverages domain-specific languages and comes with the Custard compiler
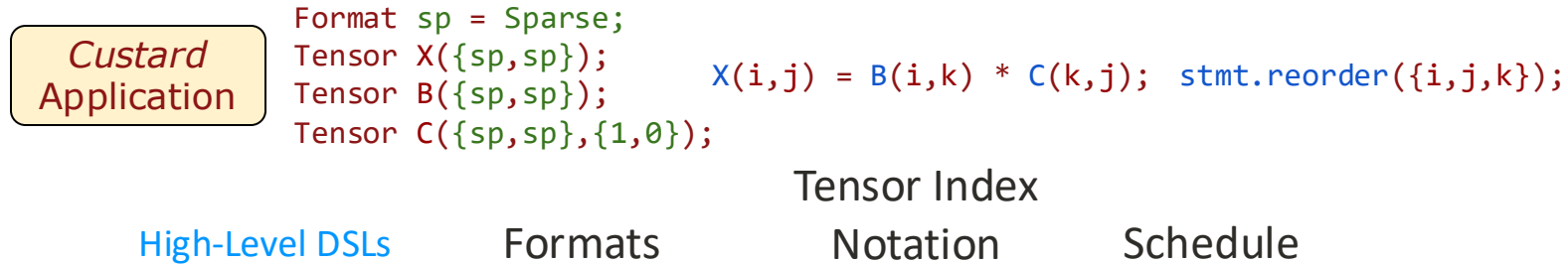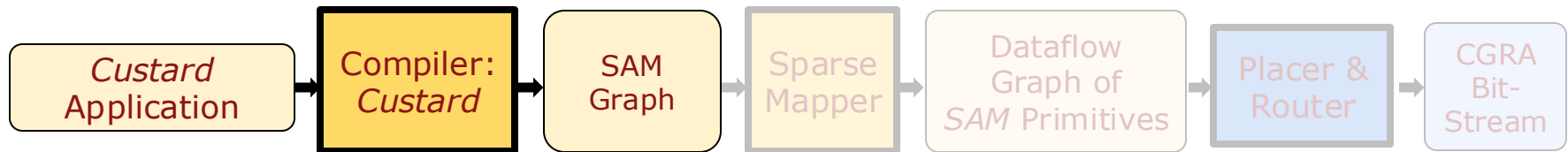


Custard Application

```
Format sp = Sparse;
Tensor X({sp,sp});
Tensor B({sp,sp});
Tensor C({sp,sp},{1,0});
```

`X(i,j) = B(i,k) * C(k,j);  stmt.reorder({i,j,k});`

High-Level DSLs     Formats     Tensor Index Notation     Schedule

Custard Application → Compiler: Custard → SAM Graph → Sparse Mapper → Dataflow Graph of SAM Primitives → Placer & Router → CGRA Bit-Stream

# SAM leverages domain-specific languages and comes with the Custard compiler



```
Format sp = Sparse;
Tensor X({sp,sp});
Tensor B({sp,sp});
Tensor C({sp,sp},{1,0});
```

X(i,j) = B(i,k) * C(k,j);  stmt.reorder({i,j,k});

Custard Application

High-Level DSLs

Formats

Tensor Index Notation

Schedule

Compiler: Custard

**The Sparse Abstract Machine**

Custard Application → Compiler: Custard → SAM Graph → Sparse Mapper → Dataflow Graph of SAM Primitives → Placer & Router → CGRA Bit-Stream

# Representing dataflow in SAM

SAM represents:

1. Wires carrying data through streams
2. Modules that compute on the data through primitives

# Representing dataflow in SAM

SAM represents:

1. Wires carrying data through streams

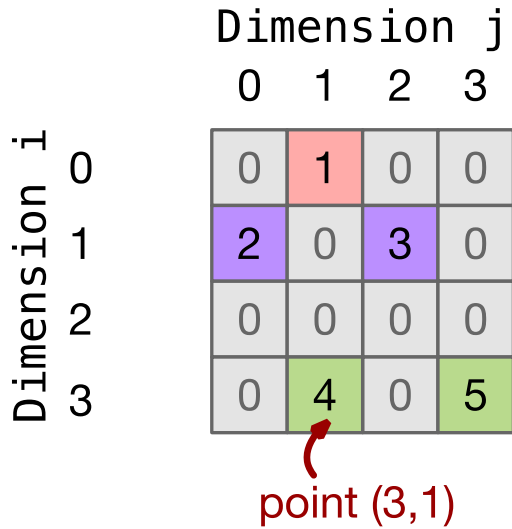2. Modules that compute on the data through primitives

Input Streams  ➡️                              ➡️  Output Streams

# Representing dataflow in SAM

SAM represents:

1. Wires carrying data through streams
2. Modules that compute on the data through primitives

Input Streams → **SAM Primitive** (State Machine) → Output Streams

# Representing tensors in SAM

point (3,1)

# Representing tensors in SAM

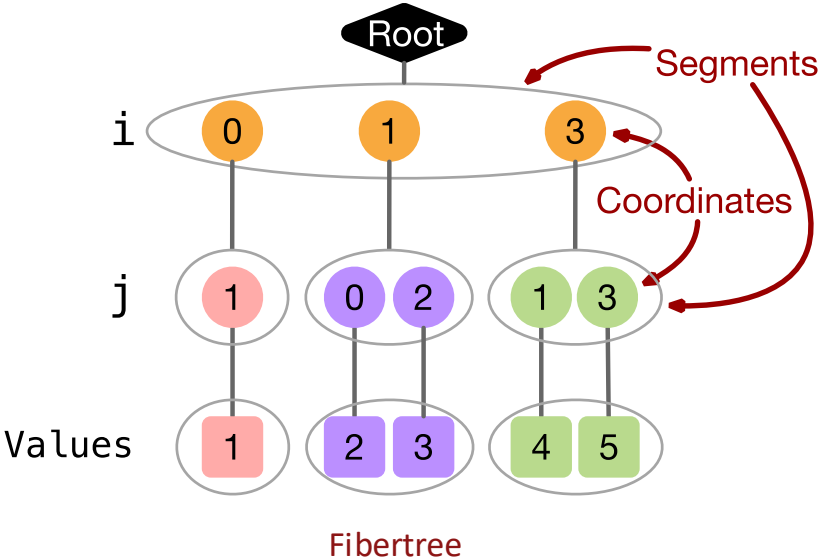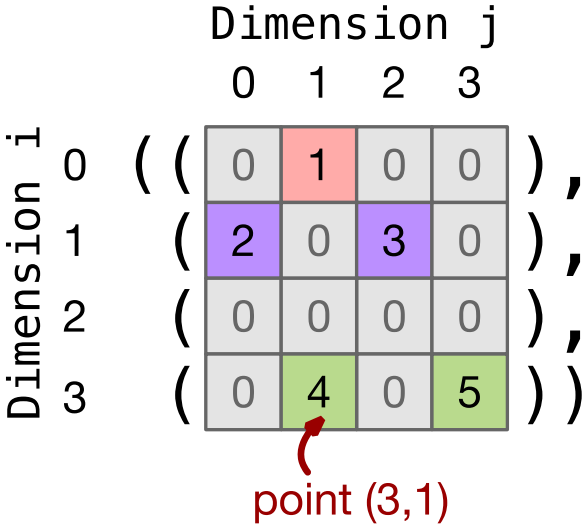# Representing tensors in SAM

Fibertree
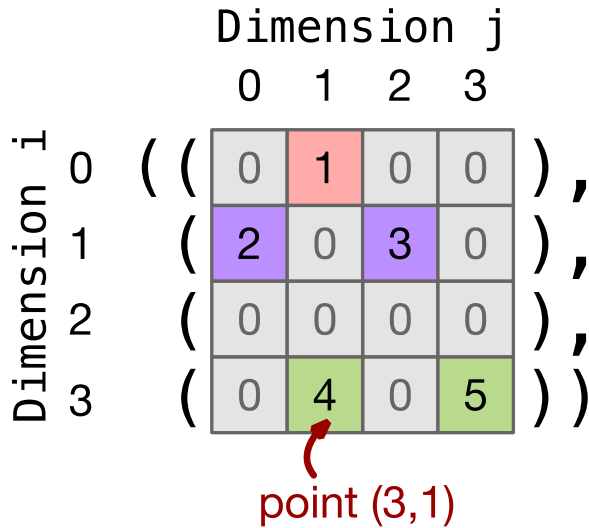
# Representing tensors in SAM

Dimension j

0   1   2   3

Dimension i

0  ( ( 0   1   0   0 ),

1    ( 2   0   3   0 ),

2    ( 0   0   0   0 ),

3    ( 0   4   0   5 ) )

point (3,1)

Root

i   0   1   3

Segments

(0, 1, 3)

Coordinates

j   1   0   2   1   3

((1), (0, 2), (1, 3))

Values   1   2   3   4   5

((1), (2, 3), (4, 5))

Fibertree

# As data structures and flattened streams

(0, 1, 3)

((1), (0, 2), (1, 3))

((1), (2, 3), (4, 5))

# As data structures and flattened streams

(0, 1, 3)

((1), (0, 2), (1, 3))

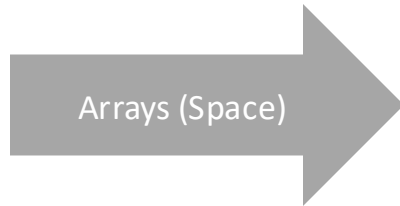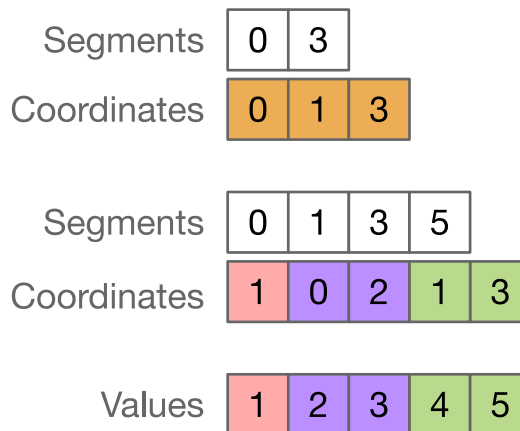((1), (2, 3), (4, 5))

Arrays (Space)

# As data structures and flattened streams
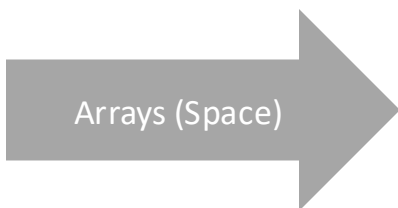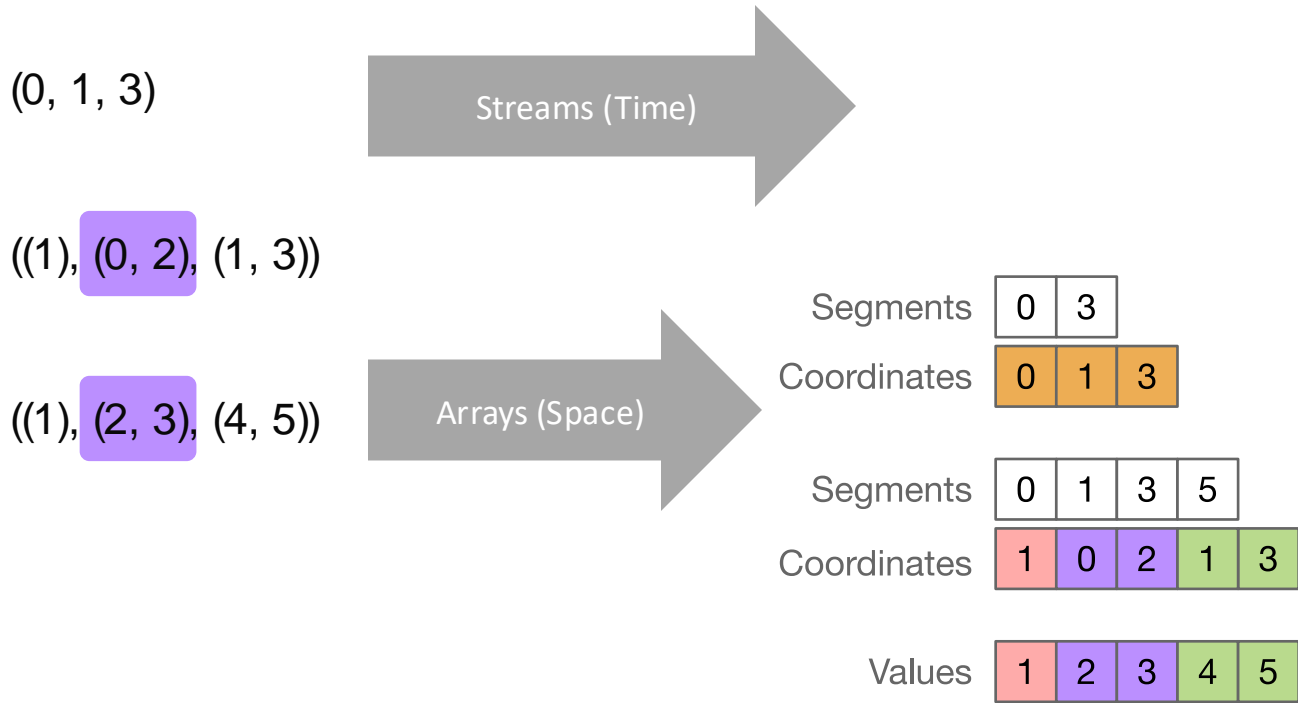
(0, 1, 3)

((1), (0, 2), (1, 3))

((1), (2, 3), (4, 5))        Arrays (Space) →

Segments | 0 | 3 |

Coordinates | 0 | 1 | 3 |

Segments | 0 | 1 | 3 | 5 |

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

Values | 1 | 2 | 3 | 4 | 5 |

Stanford University

9

# As data structures and flattened streams

(0, 1, 3)

Streams (Time)

((1), (0, 2), (1, 3))

((1), (2, 3), (4, 5))

Arrays (Space)

| Segments | 0 | 3 |
|---|---|---|

| Coordinates | 0 | 1 | 3 |
|---|---|---|---|

| Segments | 0 | 1 | 3 | 5 |
|---|---|---|---|---|

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

| Values | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|

# As data structures and flattened streams

$D, S_0, 3, 1, 0$ →

(0, 1, 3)

Streams (Time) →

((1), (0, 2), (1, 3))

((1), (2, 3), (4, 5))

Arrays (Space) →

| Segments | 0 | 3 | |
|---|---|---|---|
| Coordinates | 0 | 1 | 3 |

| Segments | 0 | 1 | 3 | 5 |
|---|---|---|---|---|
| Coordinates | 1 | 0 | 2 | 1 | 3 |

| Values | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|

# As data structures and flattened streams

$D, S_0, 3, 1, 0$

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

$(0, 1, 3)$

**Streams (Time)**

$((1), (0, 2), (1, 3))$

$((1), (2, 3), (4, 5))$

**Arrays (Space)**

| Segments | 0 | 3 |
|---|---|---|

| Coordinates | 0 | 1 | 3 |
|---|---|---|---|

| Segments | 0 | 1 | 3 | 5 |
|---|---|---|---|---|

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

| Values | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|

**Stanford University**

9

# As data structures and flattened streams

$D, S_0, 3, 1, 0$

$(0, 1, 3)$

**Streams (Time)**

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

$((1), (0, 2), (1, 3))$

$D, S_1, 5, 4, S_0, 3, 2, S_0, 1$

| Segments | 0 | 3 | |
|---|---|---|---|
| Coordinates | 0 | 1 | 3 |

$((1), (2, 3), (4, 5))$

**Arrays (Space)**

| Segments | 0 | 1 | 3 | 5 |
|---|---|---|---|---|
| Coordinates | 1 | 0 | 2 | 1 | 3 |

| Values | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|

**Stanford University**

9

# As data structures and flattened streams



SAM Graph

Control tokens

$D, S_0, 3, 1, 0$

$(0, 1, 3)$

Streams (Time)

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

Control tokens

$((1), (0, 2), (1, 3))$

$D, S_1, 5, 4, S_0, 3, 2, S_0, 1$

$((1), (2, 3), (4, 5))$

Arrays (Space)

| Segments | 0 | 3 |
|---|---|---|

| Coordinates | 0 | 1 | 3 |
|---|---|---|---|

| Segments | 0 | 1 | 3 | 5 |
|---|---|---|---|---|

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

| Values | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|

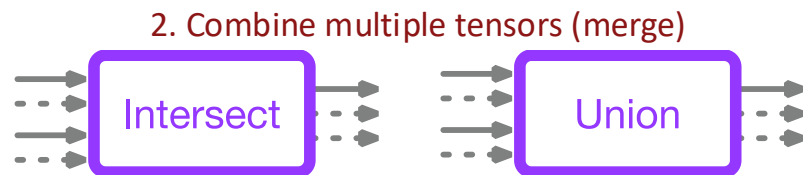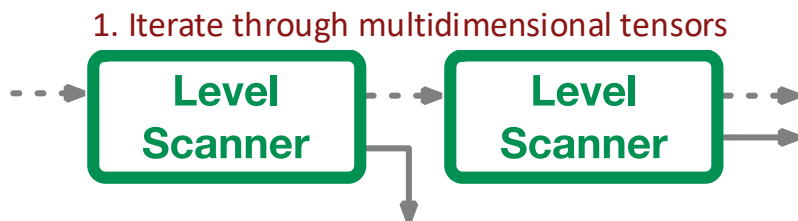# SAM supports all of tensor algebra

The sparse abstract machine has clean interfaces defined for each feature of sparse tensor algebra, called primitives
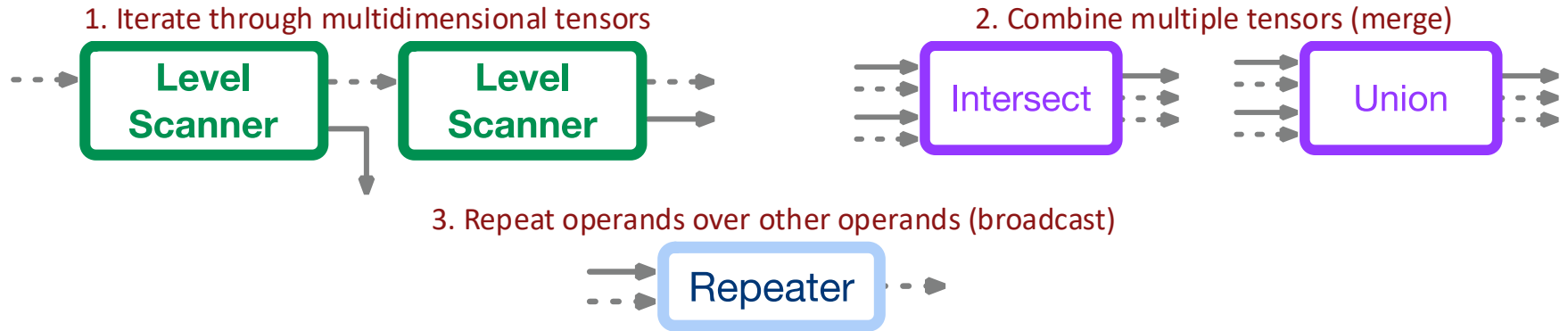
# SAM supports all of tensor algebra

The sparse abstract machine has clean interfaces defined for each feature of sparse tensor algebra, called primitives

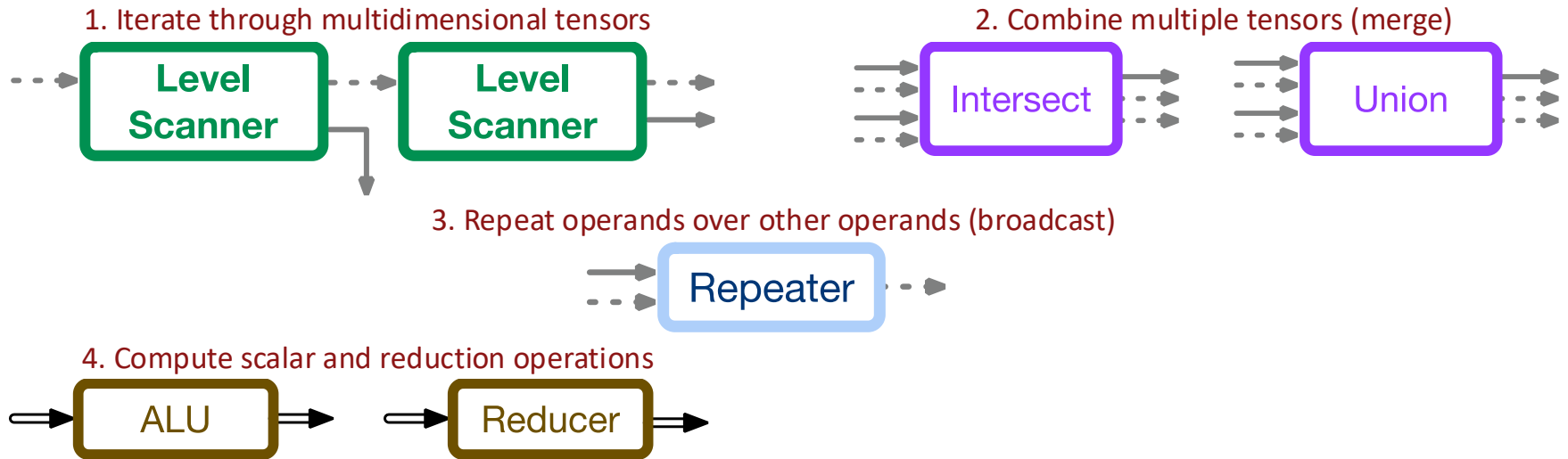1. Iterate through multidimensional tensors

# SAM supports all of tensor algebra

The sparse abstract machine has clean interfaces defined for each feature of sparse tensor algebra, called primitives



1. Iterate through multidimensional tensors

**Level Scanner** → **Level Scanner**

2. Combine multiple tensors (merge)
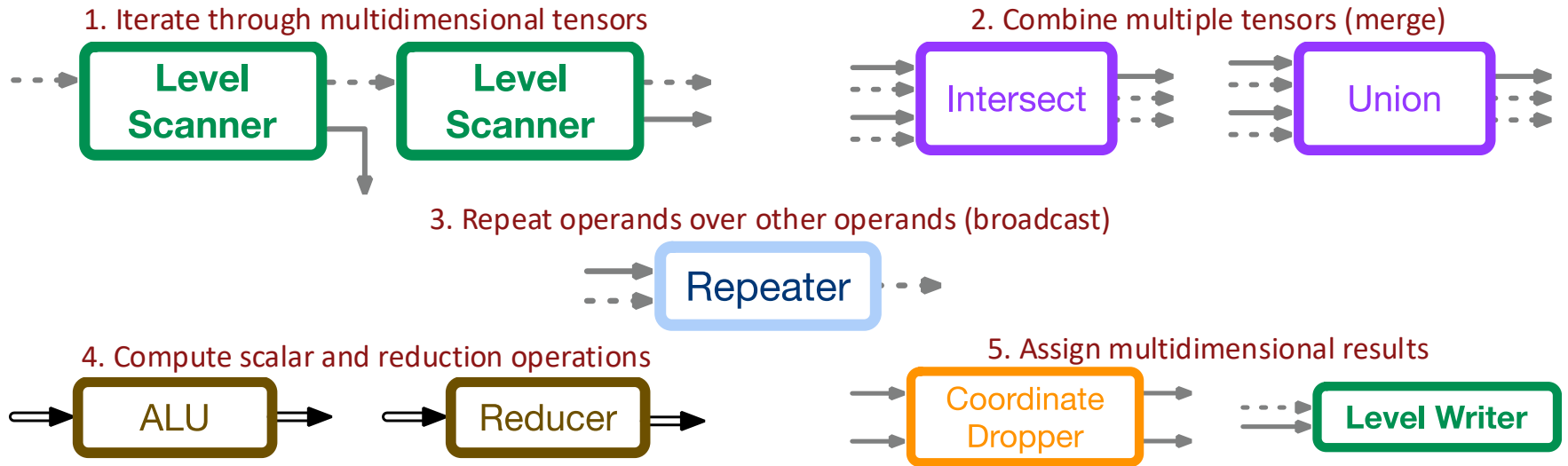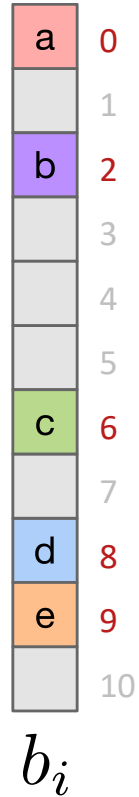
Intersect    Union

# SAM supports all of tensor algebra

The sparse abstract machine has clean interfaces defined for each feature of sparse tensor algebra, called primitives



1. Iterate through multidimensional tensors

**Level Scanner** → **Level Scanner**

2. Combine multiple tensors (merge)

Intersect   Union

3. Repeat operands over other operands (broadcast)

Repeater

# SAM supports all of tensor algebra

The sparse abstract machine has clean interfaces defined for each feature of sparse tensor algebra, called primitives
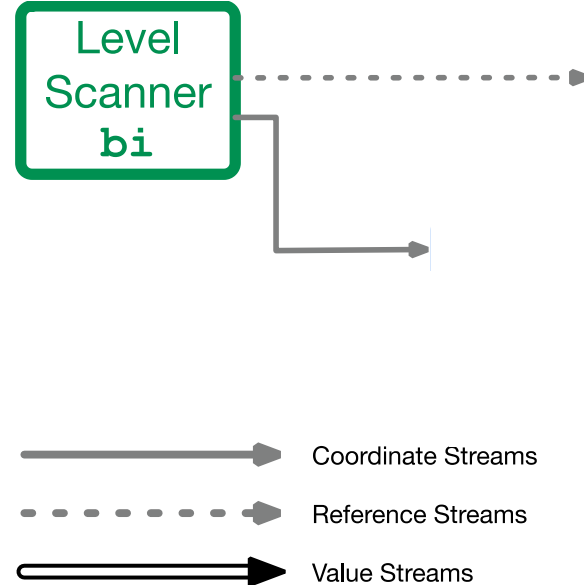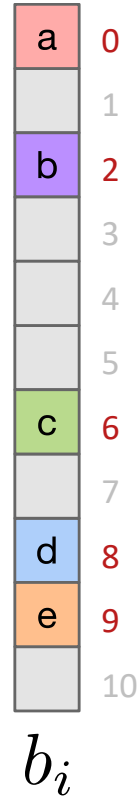


1. Iterate through multidimensional tensors

**Level Scanner** → **Level Scanner**

2. Combine multiple tensors (merge)

Intersect    Union

3. Repeat operands over other operands (broadcast)

Repeater

4. Compute scalar and reduction operations

ALU    Reducer

# SAM supports all of tensor algebra

The sparse abstract machine has clean interfaces defined for each feature of sparse tensor algebra, called primitives



1. Iterate through multidimensional tensors

**Level Scanner** → **Level Scanner**

2. Combine multiple tensors (merge)

Intersect   Union

3. Repeat operands over other operands (broadcast)

Repeater

4. Compute scalar and reduction operations

ALU → Reducer

5. Assign multidimensional results
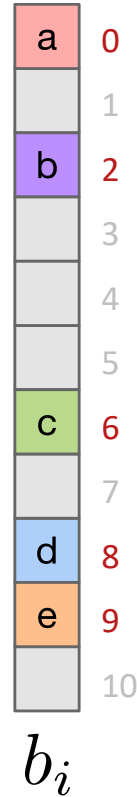
Coordinate Dropper   **Level Writer**

# Vector scaling example with repeaters

$$x_i = \underline{b_i} \cdot c$$
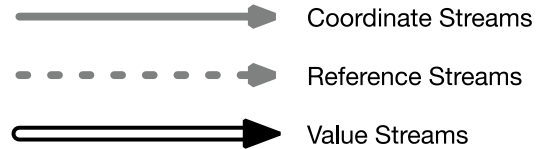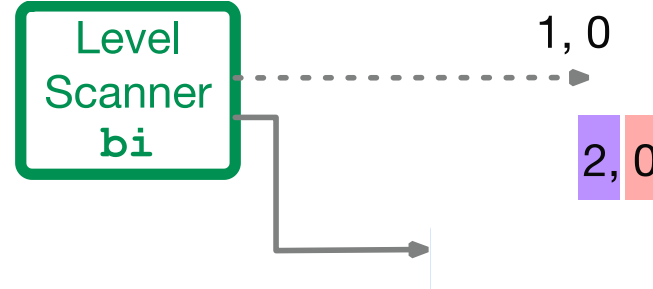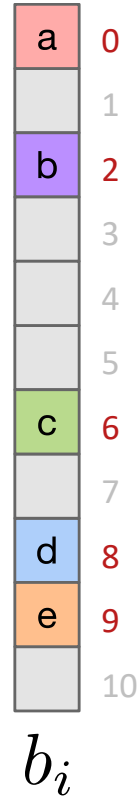


$b_i$

# Vector scaling example with repeaters

$$x_i = \underline{b_i} \cdot c$$

Level Scanner **bi**

Coordinate Streams

Reference Streams

Value Streams

$b_i$

# Vector scaling example with repeaters

$$x_i = \underline{b_i} \cdot c$$



| | |
|---|---|
| a | 0 |
| | 1 |
| b | 2 |
| | 3 |
| | 4 |
| | 5 |
| c | 6 |
| | 7 |
| d | 8 |
| e | 9 |
| | 10 |

$b_i$

Level Scanner **bi**

0

0

Coordinate Streams

Reference Streams

Value Streams

# Vector scaling example with repeaters

$$x_i = \underline{b_i} \cdot c$$

a  0
   1
b  2
   3
   4
   5
c  6
   7
d  8
e  9
   10

$b_i$

Level Scanner **bi**

2, 1, 0

6, 2, 0

→ Coordinate Streams

⇢ Reference Streams

⇒ Value Streams

# Vector scaling example with repeaters

$$x_i = \underline{b_i} \cdot c$$



Level Scanner **bi**

3, 2, 1, 0

8, 6, 2, 0

$b_i$

Coordinate Streams

Reference Streams

Value Streams

# Vector scaling example with repeaters

$$x_i = \underline{b_i} \cdot c$$



Level Scanner **bi**

4, 3, 2, 1, 0
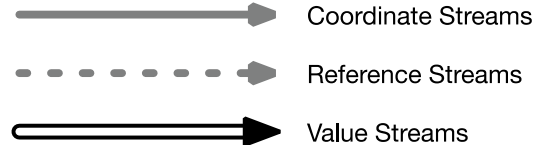
9, 8, 6, 2, 0

$b_i$

→ Coordinate Streams

⇢ Reference Streams

⇒ Value Streams

# Vector scaling example with repeaters

$$x_i = \underline{b_i} \cdot c$$



| | | |
|---|---|---|
| a | 0 | |

Level Scanner **bi**

$S_0$, 4, 3, 2, 1, 0

$S_0$, 9, 8, 6, 2, 0
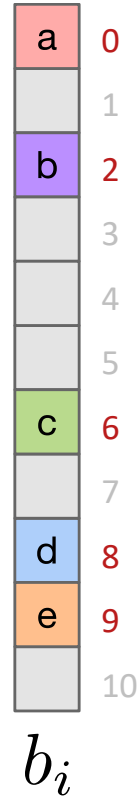
$b_i$

Coordinate Streams

Reference Streams

Value Streams

# Vector scaling example with repeaters

$$x_i = \underline{b_i} \cdot c$$



D, S$_0$, 4, 3, 2, 1, 0

Level
Scanner
**bi**

D, S$_0$, 9, 8, 6, 2, 0

→ Coordinate Streams

⇢ Reference Streams

⇒ Value Streams

$b_i$

# Vector scaling example with repeaters

$$x_i = \underline{b_i \cdot c}$$

$$c \cdot$$



Level Scanner **bi**

D, $S_0$, 4, 3, 2, 1, 0

D, $S_0$, 9, 8, 6, 2, 0

Repeater **c**

D, 0

$b_i$

Coordinate Streams

Reference Streams

Value Streams

# Vector scaling example with repeaters

$$x_i = \underline{b_i} \cdot c$$

$c \cdot$

$c \cdot$



$b_i$

**Level Scanner** `bi`

D, S$_0$, 4, 3, 2, 1, 0

D, S$_0$, 9, 8, 6, 2, 0

**Repeater** `c`

D, 0

0

Coordinate Streams

Reference Streams

Value Streams

# Vector scaling example with repeaters

$$x_i = \underline{b_i} \cdot c$$



$c \cdot$ a   0

$c \cdot$ b   2

$c \cdot$

c · 6

d 8

e 9

$b_i$

Level Scanner **bi**

D, S$_0$, 4, 3, 2, 1, 0

D, S$_0$, 9, 8, 6, 2, 0

Repeater **c**

D, 0

0, 0

Coordinate Streams

Reference Streams

Value Streams

# Vector scaling example with repeaters

$$x_i = \underline{b_i} \cdot c$$



Level Scanner **bi**

D, S$_0$, 4, 3, 2, 1, 0

D, S$_0$, 9, 8, 6, 2, 0

Repeater **c**

D, 0

0, 0, 0, 0

Coordinate Streams

Reference Streams

Value Streams

$b_i$

# Vector scaling example with repeaters

$$x_i = \underline{b_i} \cdot c$$



$c \cdot$ a   0
1
$c \cdot$ b   2
3
4
5
$c \cdot$ $c \cdot$ c   6
7
$c \cdot$ d   8
$c \cdot$ e   9
10

$b_i$

Level Scanner **bi**

D, S$_0$, 4, 3, 2, 1, 0

D, S$_0$, 9, 8, 6, 2, 0

Repeater c

D, 0

D, S$_0$, 0, 0, 0, 0, 0

Coordinate Streams

Reference Streams

Value Streams

# Vector scaling example with repeaters
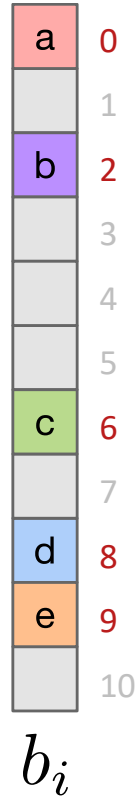
$$x_i = \underline{b_i} \cdot c$$



Level Scanner **bi**

D, $S_0$, 4, 3, 2, 1, 0

D, $S_0$, 9, 8, 6, 2, 0

Repeater c

D, 0

D, $S_0$, 0, 0, 0, 0, 0

Coordinate Streams

Reference Streams

Value Streams

$b_i$

# Primitives compose to compute expressions: SpM*SpM

# Primitives compose to compute expressions: SpM*SpM



$$\forall_i \forall_k \forall_j X_{ij} = B_{ik} \cdot C_{kj}$$

**Stanford University**

# Primitives compose to compute expressions: SpM*SpM

$$\forall_i \forall_k \forall_j X_{ij} = B_{ik} \cdot C_{kj}$$

**Stanford University**

# Primitives compose to compute expressions: SpM*SpM

SAM Graph

$$\forall_i \forall_k \forall_j X_{ij} = B_{ik} \cdot C_{kj}$$

**Stanford University**

# Primitives compose to compute expressions: SpM*SpM

$$\forall_i \forall_k \forall_j X_{ij} = B_{ik} \cdot C_{kj}$$

**Stanford University**

25

# Primitives compose to compute expressions: SpM*SpM

$$\forall_i \forall_k \forall_j X_{ij} = B_{ik} \cdot C_{kj}$$

**Stanford University**

# Primitives compose to compute expressions: SpM*SpM



SAM Graph

Format agnostic
Hierarchical sparse iteration

Hierarchical Streams

Data flows through computation spatially

$$\forall_i \forall_k \forall_j X_{ij} = B_{ik} \cdot C_{kj}$$

# Primitives compose to compute expressions: SpM*SpM

Format agnostic
Hierarchical sparse iteration

Hierarchical Streams

Data flows through computation spatially

Stream Merging

Dropped

**Level Scanner Bi** compressed

**Xi coordinate stream**

**Level Scanner Bk** compressed

**Repeater Ci**

**Level Scanner Ck** compressed

**Intersecter**

**Level Scanner Cj** compressed

**xj coordinate stream**

**Repeater Bj**

Array B vals

Array C vals

**Xi coordinate stream**
**Xj coordinate stream**

**Coordinate Dropper**

Multiplier

Vector Reducer

**Level Writer Xi** compressed

**Level Writer Xj** compressed

**Level Writer X vals** compressed

D, 0

D, 0

$$\forall_i \forall_k \forall_j X_{ij} = B_{ik} \cdot C_{kj}$$

**Stanford University**

25

# Primitives compose to compute expressions: SpM*SpM

Format agnostic
Hierarchical sparse iteration

Tensor broadcasting

Hierarchical Streams

Data flows through computation spatially

Stream Merging

$$\forall_i \forall_k \forall_j X_{ij} = B_{ik} \cdot C_{kj}$$

**Stanford University**

# Primitives compose to compute expressions: SpM*SpM

Format agnostic
Hierarchical sparse iteration

Tensor broadcasting

**Level Scanner Bi** compressed

**Xi coordinate stream**

**Level Scanner Bk** compressed

**Intersecter**

**Repeater Bj**

**Xi coordinate stream**
**Xj coordinate stream**

**Coordinate Dropper**

**Level Writer Xi** compressed

D, 0

D, 0

**Repeater Ci**

**Level Scanner Ck** compressed

**Level Scanner Cj** compressed

**Xj coordinate stream**

**Array B vals**

**Array C vals**

**Multiplier**

**Vector Reducer**

**Level Writer Xj** compressed

**Level Writer X vals** compressed

Hierarchical Streams

Data flows through computation spatially

Dropped

Stream Merging

Computation

$$\forall_i \forall_k \forall_j X_{ij} = B_{ik} \cdot C_{kj}$$

**Stanford University**

25

# Primitives compose to compute expressions: SpM*SpM



Format agnostic
Hierarchical sparse iteration

Tensor broadcasting

Hierarchical Streams

Data flows through computation spatially

Stream Merging

Computation

Tensor Construction

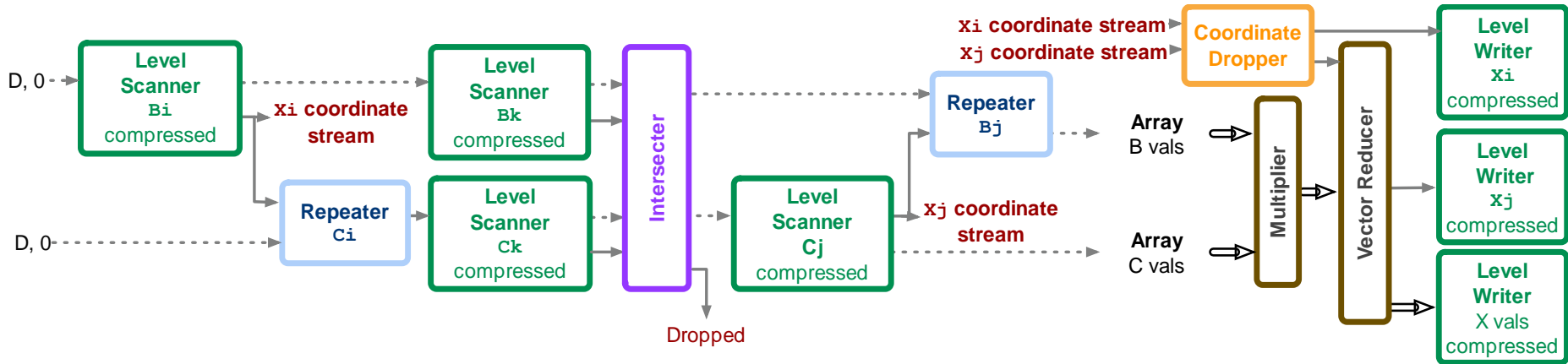$$\forall_i \forall_k \forall_j X_{ij} = B_{ik} \cdot C_{kj}$$

**Stanford University**

25

# Primitives compose to compute expressions: SpM*SpM

$$\forall_i \forall_k \forall_j X_{ij} = B_{ik} \cdot C_{kj}$$

**Stanford University**

# Primitives compose to compute expressions: SpM*SpM



SAM Graph

Format agnostic
Hierarchical sparse iteration

Tensor broadcasting

Sparse outer-level accumulation

D, 0 → **Level Scanner Bi** compressed

**Xi coordinate stream**

**Level Scanner Bk** compressed

**Intersecter**

**Xi coordinate stream**
**Xj coordinate stream**

**Coordinate Dropper**

**Level Writer Xi** compressed

D, 0 → **Repeater Ci**

**Level Scanner Ck** compressed

**Level Scanner Cj** compressed

**xj coordinate stream**

**Repeater Bj**

**Array B vals**

**Multiplier**

**Vector Reducer**

**Level Writer Xj** compressed

**Array C vals**

**Level Writer X vals** compressed

Dropped

Hierarchical Streams

Data flows through computation spatially

Stream Merging

Input Iteration and Stream Merging

Computation

Tensor Construction

$$\forall_i \forall_k \forall_j X_{ij} = B_{ik} \cdot C_{kj}$$

**Stanford University**

25

# Primitives compose to compute expressions: SpM*SpM

Format agnostic
Hierarchical sparse iteration

Tensor broadcasting

Sparse outer-level accumulation

Hierarchical Streams

Data flows through computation spatially

Stream Merging

Input Iteration and Stream Merging

Computation

Tensor Construction

$$\forall_i \forall_k \forall_j X_{ij} = B_{ik} \cdot C_{kj}$$

**Stanford University**

25

# Inner-product algorithm in SAM



$$\forall_i \forall_j \forall_k X_{ij} = B_{ik} * C_{kj}$$

# Inner-product algorithm in SAM



$$\forall_i \forall_j \forall_k X_{ij} = B_{ik} * C_{kj}$$

**Stanford University**

# Inner-product algorithm in SAM



Scalar reducer, less memory

$$\forall_i \forall_j \forall_k X_{ij} = B_{ik} * C_{kj}$$

**Stanford University**

# Inner-product algorithm in SAM



Scalar reducer, less memory

Intersection at last dataflow level

$$\forall_i \forall_j \forall_k X_{ij} = B_{ik} * C_{kj}$$

**Stanford University**

# Custard's compiler algorithm to SAM

Expression

Format Language

Schedule

$$X_{ij} = B_{ik} * C_{kj}$$

```
B=({sparse,sparse},
   {mode0,mode1})
C=({sparse, sparse},
   {mode1, mode0})
```

```
reorder(i, k, j)
```

# Custard's compiler algorithm to SAM

Expression

Format Language

$$X_{ij} = B_{ik} * C_{kj}$$

```
B=({sparse,sparse},
   {mode0,mode1})
C=({sparse, sparse},
   {mode1, mode0})
```

Schedule

```
reorder(i, k, j)
```

$$\forall_i \forall_k \forall_j X_{ij} \quad \texttt{+=} \quad (B_{ik} * C_{kj})$$

# Custard's compiler algorithm to SAM

Expression

$X_{ij} = B_{ik} * C_{kj}$

Format Language
B=({sparse,sparse},
{mode0,mode1})
C=({sparse, sparse},
{mode1, mode0})

Schedule
reorder(i, k, j)

$\forall_i \forall_k \forall_j X_{ij}$ **+=** $(B_{ik} * C_{kj})$

$C_k$

$B_i \longrightarrow$ (i) $\cap$ (k) $C_j$ (j)

$B_k$

$X_i$

$X_j$

# Custard's compiler algorithm to SAM

Expression

$$X_{ij} = B_{ik} * C_{kj}$$

Format Language

```
B=({sparse,sparse},
    {mode0,mode1})
C=({sparse, sparse},
    {mode1, mode0})
```

Schedule

```
reorder(i, k, j)
```

$$\forall_i \forall_k \forall_j X_{ij} \; += \; (B_{ik} * C_{kj})$$

Lower to Example Implementation

# Demo: Generating SAM graphs with Custard

```
> ./sparse_demo.sh compile
```
- This runs the applications in `./sam/compiler/sam-kernels.sh` through the Custard compiler
- All SAM graphs generated in `./sam/compiler/sam-outputs/`

- View the SpMSpM kernel `matmul_ijk` in `./sam/compiler/sam-outputs/png/matmul_ijk.png` in VSCode or using `docker cp`
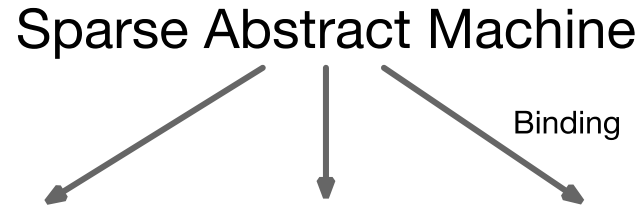
- We will also view a smaller kernel, `mat_elemadd` in `./sam/compiler/sam-outputs/png/mat_elemadd.png,` which we will be using for the rest of the demo

# Although abstract, SAM binds to real hardware

Sparse Abstract Machine

Binding

# Although abstract, SAM binds to real hardware

Sparse Abstract Machine

Binding

Reconfigurable
Accelerators
(RDAs and CGRAs)

# Although abstract, SAM binds to real hardware



Sparse Abstract Machine

Binding

Cycle-approximate
Simulator

Reconfigurable
Accelerators
(RDAs and CGRAs)

# Although abstract, SAM binds to real hardware

Sparse Abstract Machine

Binding

Cycle-approximate
Simulator

Fixed-function
Accelerators

Reconfigurable
Accelerators
(RDAs and CGRAs)

# Although abstract, SAM binds to real hardware



Sparse Abstract Machine

Binding

Cycle-approximate Simulator    Fixed-function Accelerators    Reconfigurable Accelerators (RDAs and CGRAs)

# Programming dataflow requires an abstract machine that a compiler can target

High-Level
Input Languages
(APIs)



Sparse Accelerator

# Programming dataflow requires an abstract machine that a compiler can target

High-Level
Input Languages
(APIs)

Compile



Sparse Accelerator

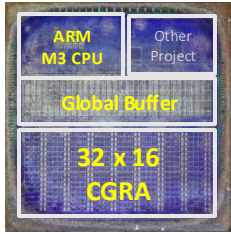# Programming dataflow requires an abstract machine that a compiler can target

High-Level Input Languages (APIs)

Compile



Sparse Accelerator

A clean interface decouples the compiler from specific hardware implementations

# Programming dataflow requires an abstract machine that a compiler can target

**High-Level Input Languages (APIs)** → Compile → Map → Sparse Accelerator
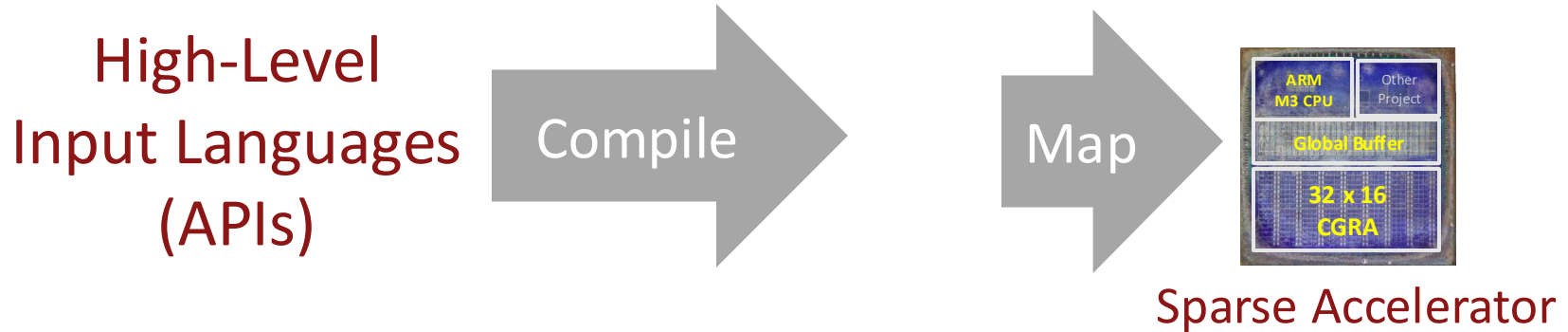
A clean interface decouples the compiler from specific hardware implementations

# Programming dataflow requires an abstract machine that a compiler can target

[Hsu et al. ASPLOS 2023]



High-Level Input Languages (APIs) → Compile → The Sparse Abstract Machine (SAM) → Map → Sparse Accelerator

A clean interface decouples the compiler from specific hardware implementations

# Programming dataflow requires an abstract machine that a compiler can target

[Hsu et al. ASPLOS 2023]



High-Level
Input Languages
(APIs)

Compile

The _ Abstract _ (SAM)
Machine

Clean + unified interface
≈ LLVM IR, Stable ISA

Map

ARM M3 CPU | Other Project
Global Buffer
32 x 16 CGRA

Sparse Accelerator

A clean interface decouples the compiler from specific hardware implementations

# Programming dataflow requires an abstract machine that a compiler can target
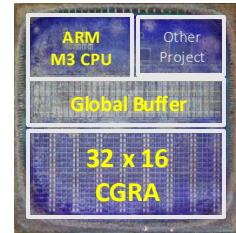
[Hsu et al. ASPLOS 2023]

**High-Level Input Languages (APIs)**

Compile

The Sparse Abstract Machine (SAM)

Map

Clean + unified interface
≈ LLVM IR, Stable ISA

**Sparse Accelerator**
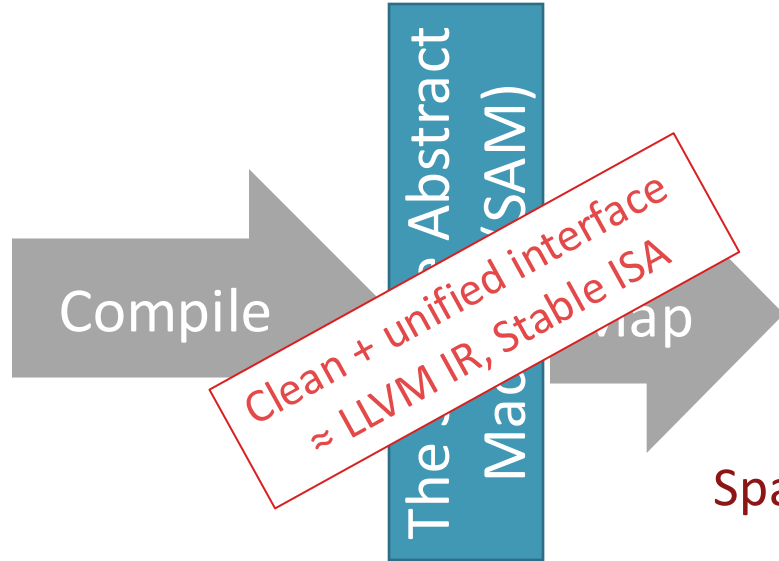
ARM M3 CPU
Other Project
Global Buffer
32 x 16 CGRA

A clean interface decouples the compiler from specific hardware implementations

# Programming dataflow requires an abstract machine that a compiler can target

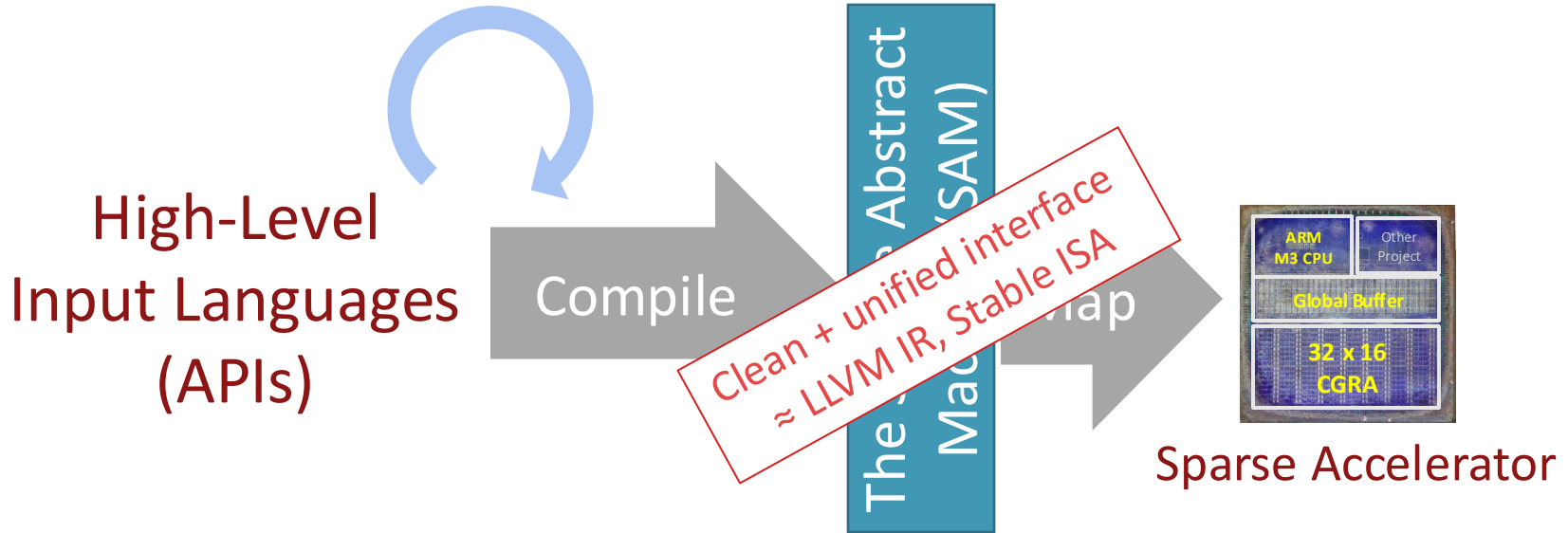[Hsu et al. ASPLOS 2023]
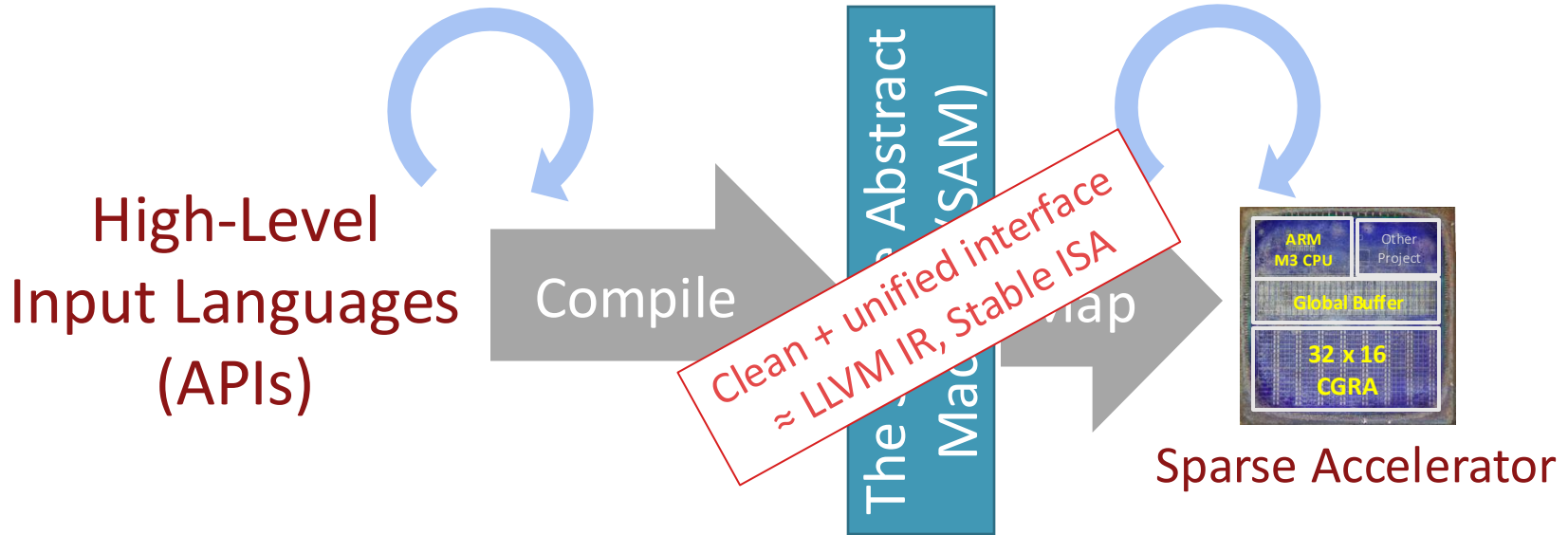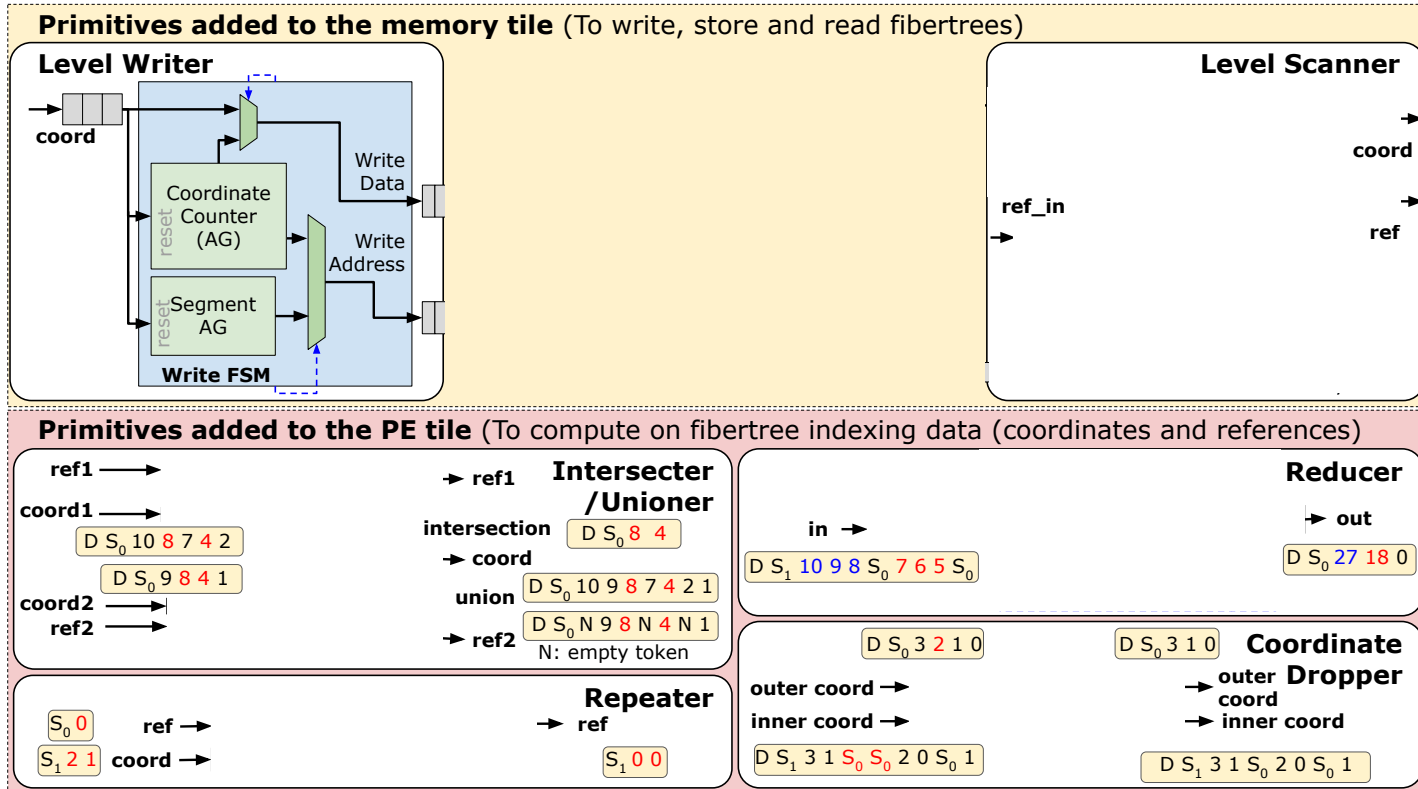


High-Level Input Languages (APIs)

Compile

The Sparse Abstract Machine (SAM)

Map

Clean + unified interface ≈ LLVM IR, Stable ISA

Sparse Accelerator

A clean interface decouples the compiler from specific hardware implementations

# SAM as the architectural specification of our sparse CGRA fabric

**Primitives added to the memory tile** (To write, store and read fibertrees)

**Level Writer**

→
coord

**Level Scanner**

→
coord

ref_in
→

ref

**Primitives added to the PE tile** (To compute on fibertree indexing data (coordinates and references)

ref1 →
coord1 →

D $S_0$ 10 8 7 4 2

D $S_0$ 9 8 4 1

coord2 →
ref2 →

**Intersecter /Unioner**

→ ref1

intersection   D $S_0$ 8 4

→ coord

union   D $S_0$ 10 9 8 7 4 2 1

→ ref2   D $S_0$ N 9 8 N 4 N 1

N: empty token

**Reducer**

in →

→ out

D $S_1$ 10 9 8 $S_0$ 7 6 5 $S_0$

D $S_0$ 27 18 0

D $S_0$ 3 2 1 0

D $S_0$ 3 1 0   **Coordinate Dropper**

outer coord →
inner coord →

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

→ outer coord
→ inner coord

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

**Repeater**

$S_0$ 0

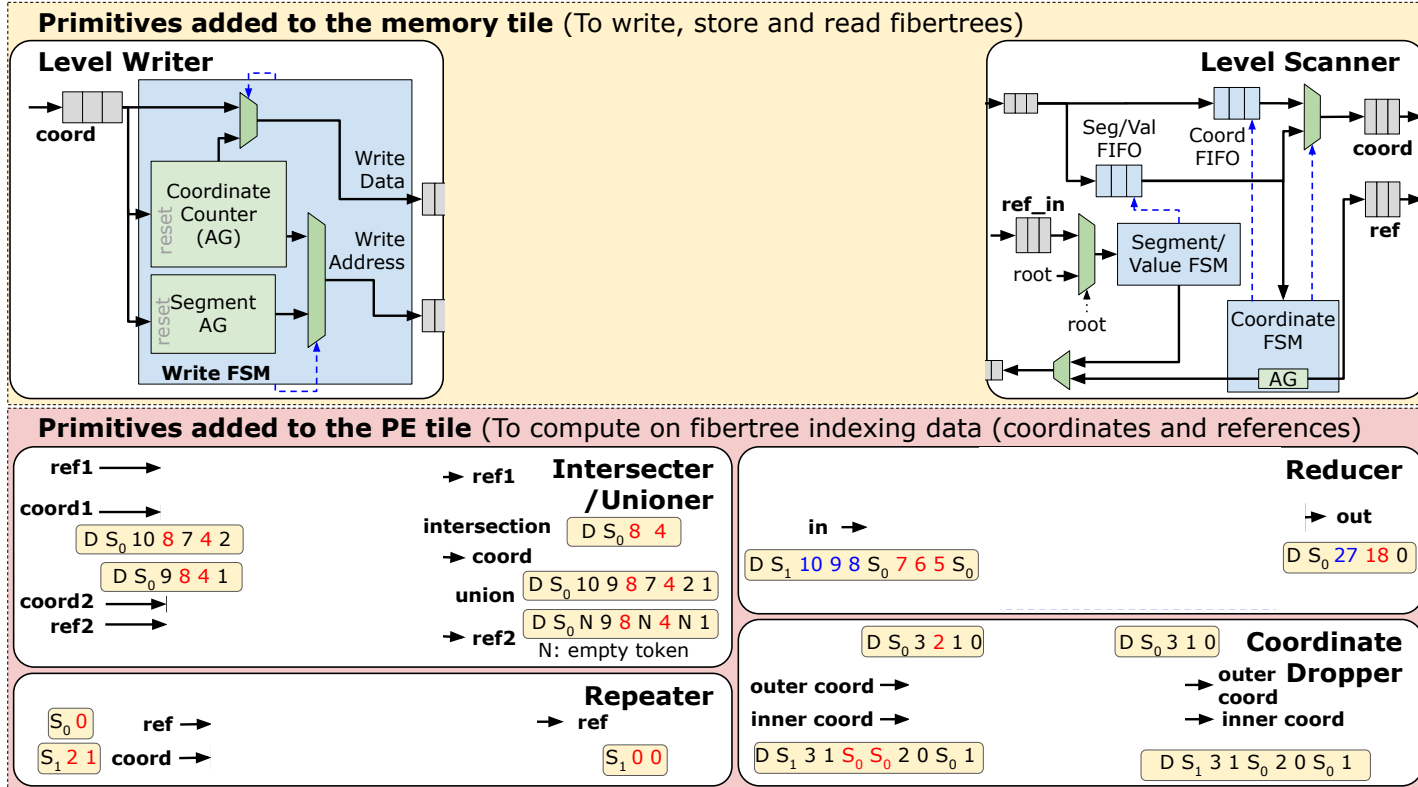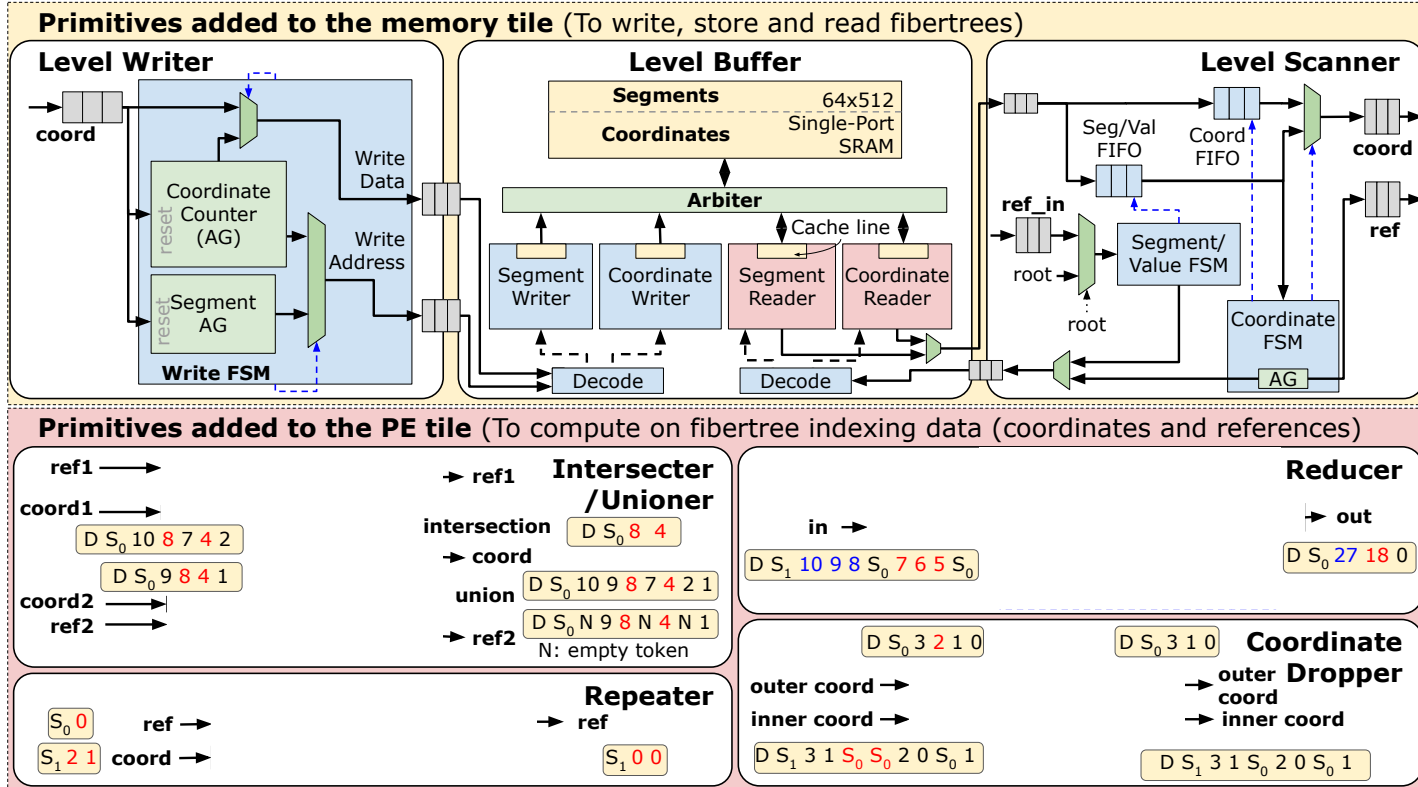$S_1$ 2 1   ref →
coord →

→ ref

$S_1$ 0 0

# SAM as the architectural specification of our sparse CGRA fabric

# SAM as the architectural specification of our sparse CGRA fabric

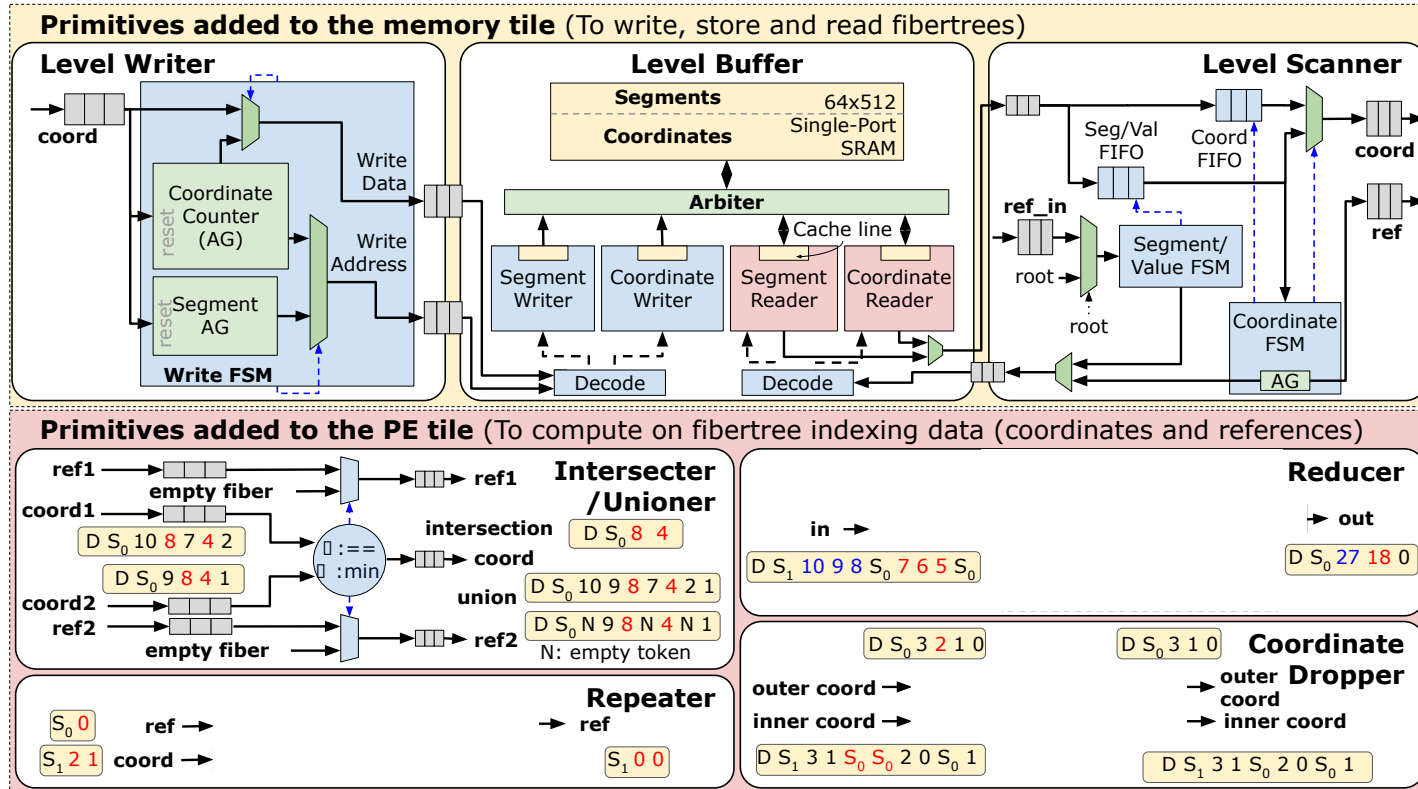# SAM as the architectural specification of our sparse CGRA fabric
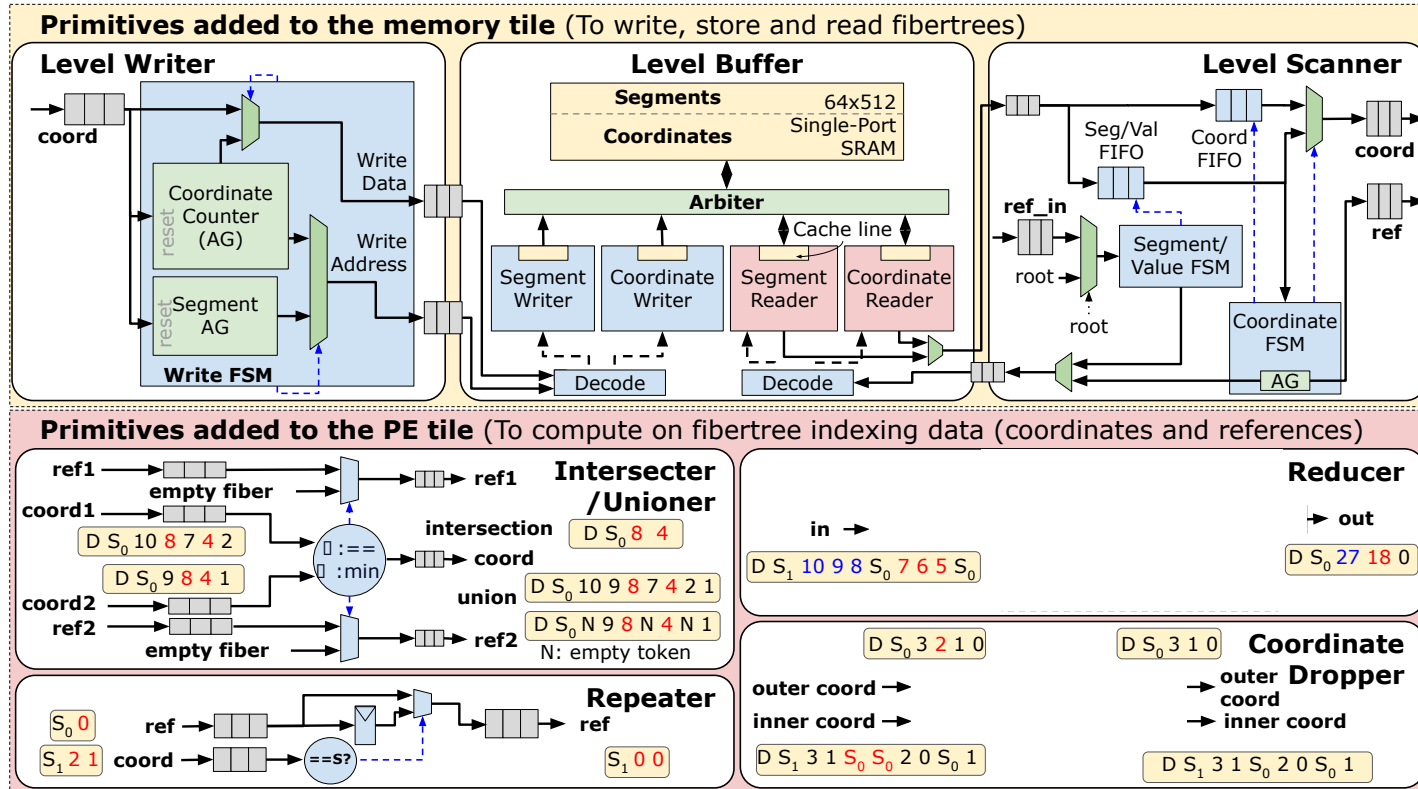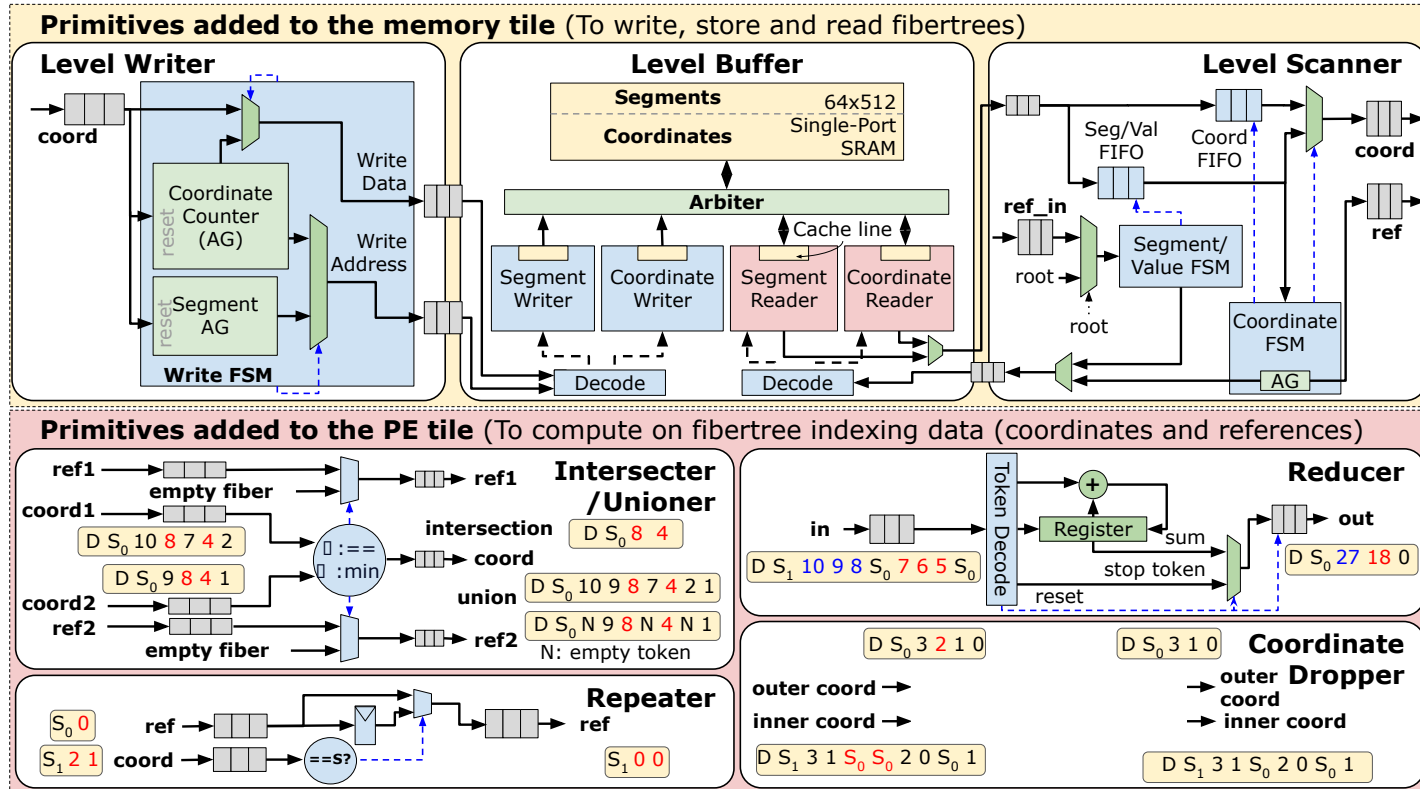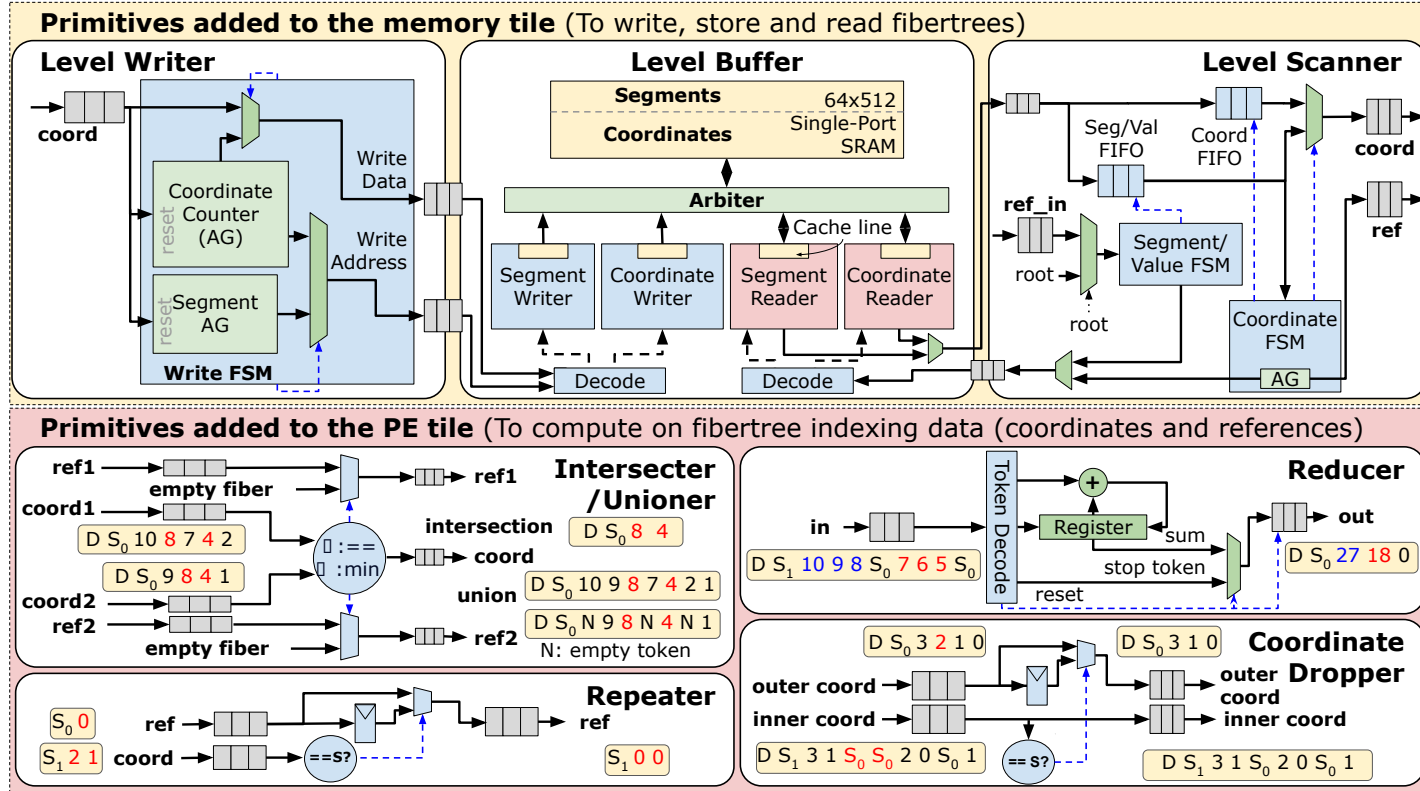
# SAM as the architectural specification of our sparse CGRA fabric

# SAM as the architectural specification of our sparse CGRA fabric

# SAM as the architectural specification of our sparse CGRA fabric

# SAM as the architectural specification of our sparse CGRA fabric

# Hardware-aware sparse dataflow graph

- CGRA architecture and microarchitecture requires more transformations during binding

- Introduce the concept of a hardware-aware SAM graph

- Performs transformations like:
  - Broadcast removal
  - Decomposition of N-joiners to binary joiners
  - Merges Level Scanners and Level Writers
  - Inserts Level Buffers

# Demo: Mapping to CGRA Microarchitecture

`>./sparse_demo.sh lower`

- This runs the SAM graph through the lowering process to produce a hardware-aware sparse dataflow graph

- We can visualize that graph in `/aha/sam/hw_aware_mat_elemadd.png`

# Tool flow that maps SAM to a 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

# Tool flow that maps SAM to a 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

Stanford University

# Demo: Generating CGRA Bitstream for sparse applications

Run the following command:

> `./sparse_demo.sh gen`

- This generates a CGRA bitstream from the hardware-aware graph using tools introduced later
- It also generates a testbench that runs an example matrix through, checking it with gold (written in Numpy)
- Explore output files generated in
  `/aha/garnet/SPARSE_TESTS/mat_elemadd_0/`

- Dataflow hardware, like CGRAs, can speed up sparse computation
- Presented ideas from the Sparse Abstract Machine and Onyx
  - SAM is an abstract IR that represents sparse tensor algebra as dataflow graphs
  - SAM comes with the Custard front-end compiler
- Introduced the AHA flow for sparse applications

# Conclusion