

# Sparse Compiler Design

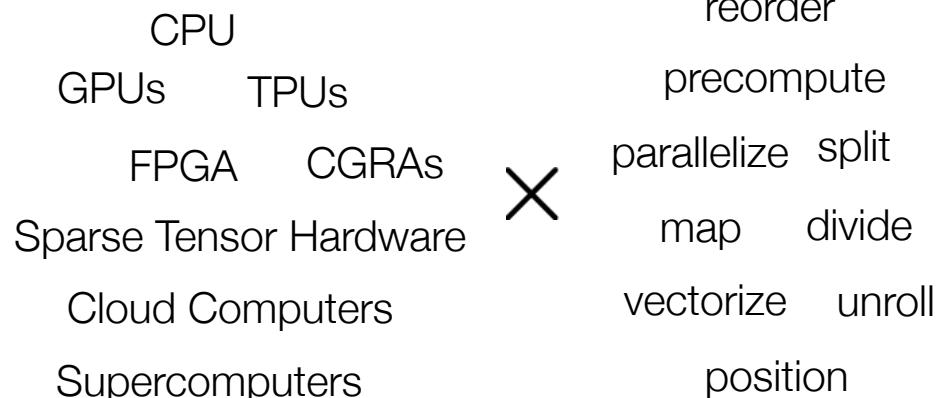
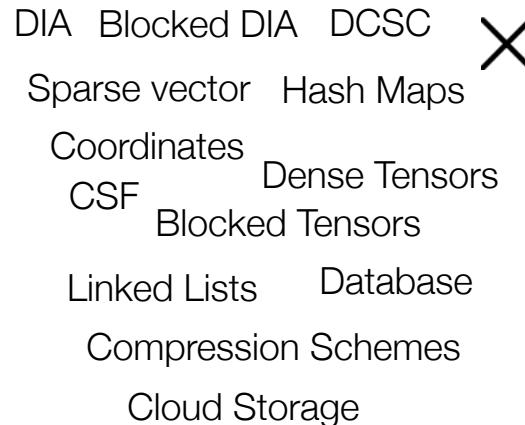
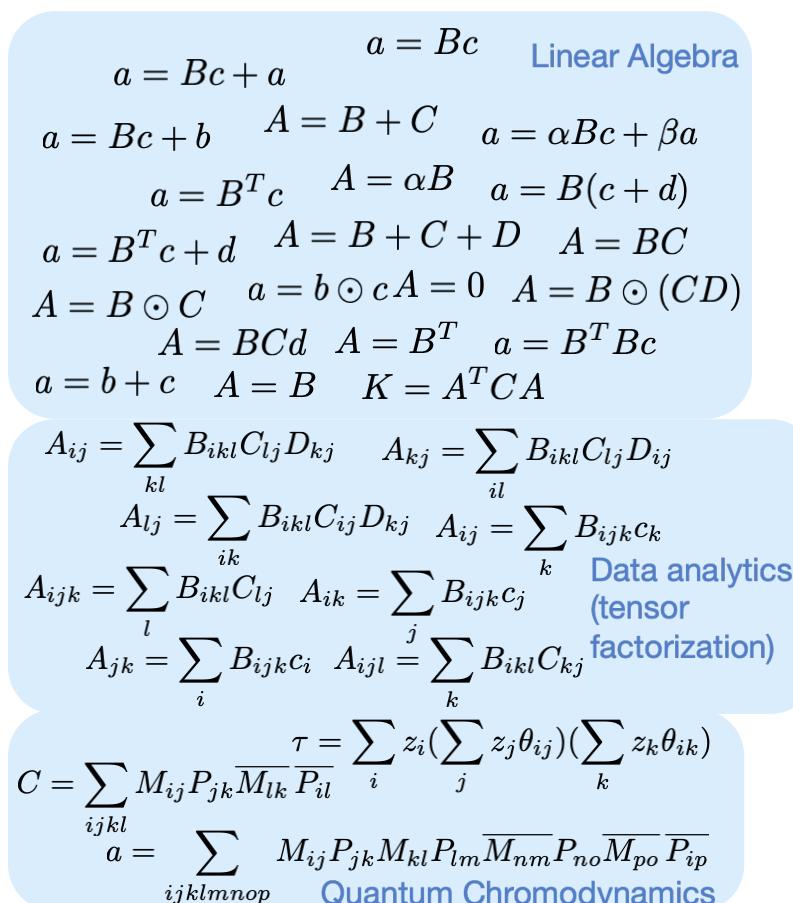
Olivia Hsu

In collaboration with Max Strange, Jaeyeon Won, Kunle Olukotun, Mark Horowitz, Joel Emer, and Fredrik Kjolstad

September 16<sup>th</sup>, 2021

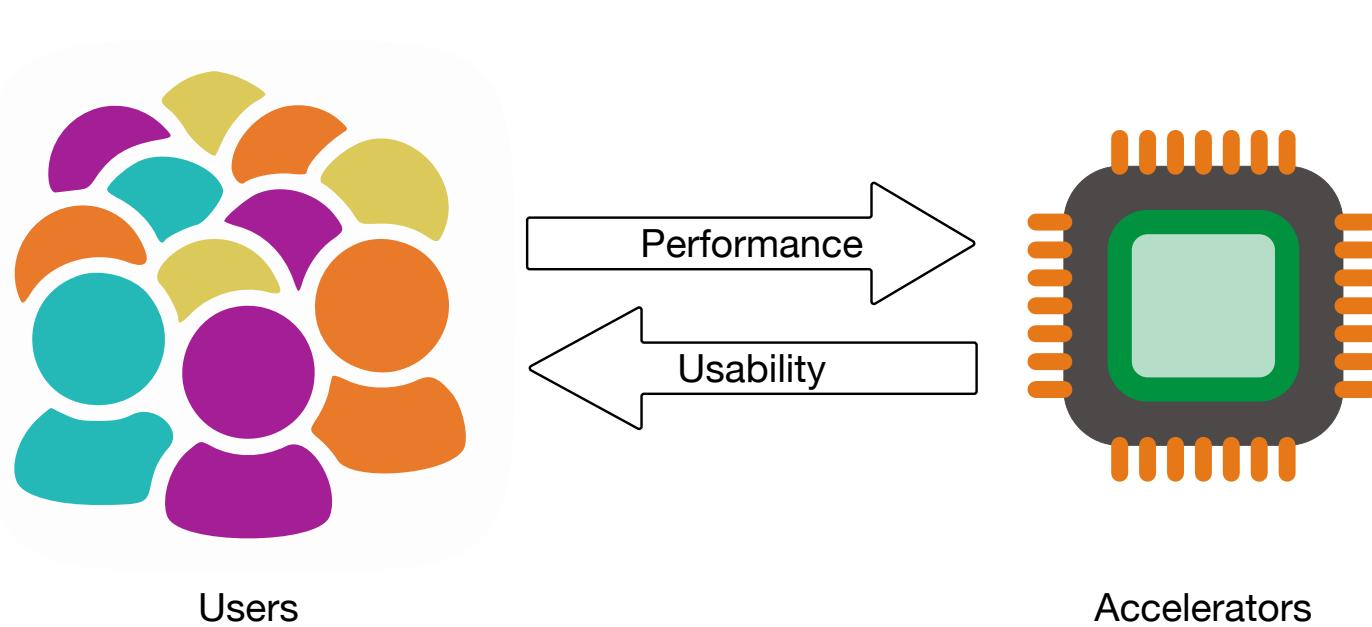
# Motivation

Too many tensor kernels for fixed-function libraries and backends



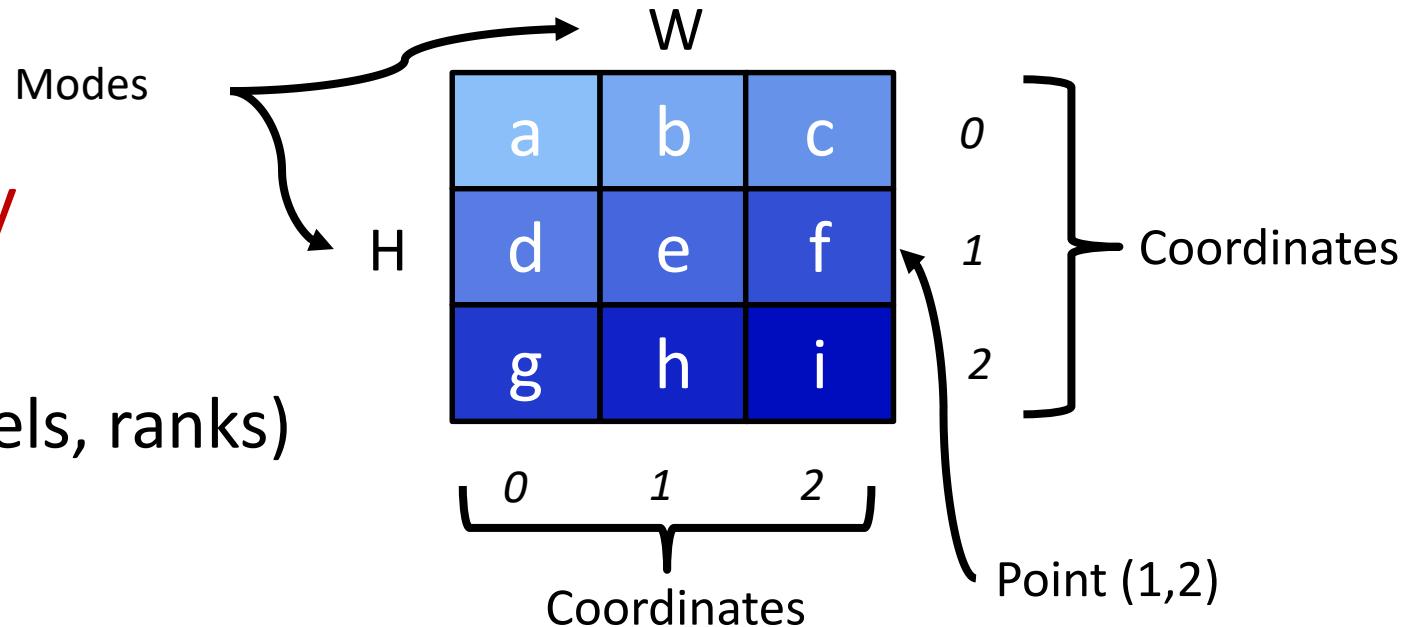
# Benefits

Long Tail of Expressions (algorithm) × Varying Compression Structures (format) × More Performant Backends (platform) × Backend-Specific Transformations (schedules)

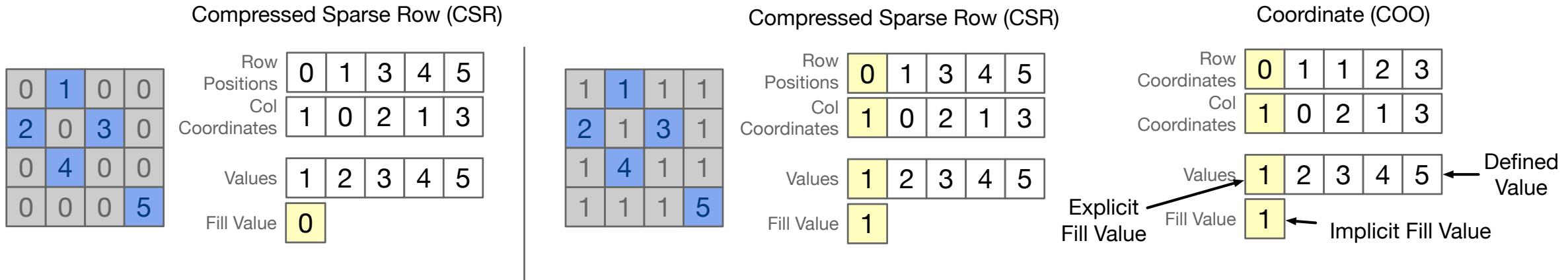


# Background: Tensor Terminology

- **Modes** (aka dimensions, levels, ranks)
  - The order refers to the number of modes
    - Ex. Matrices are 2<sup>nd</sup> order
- Each component of a rank is identified by a **coordinate**
  - Represents the ***data address space***
- Each value element in a tensor is referred to as a “**point**”
  - Points are identified by an ***ordered list*** of coordinates for each rank

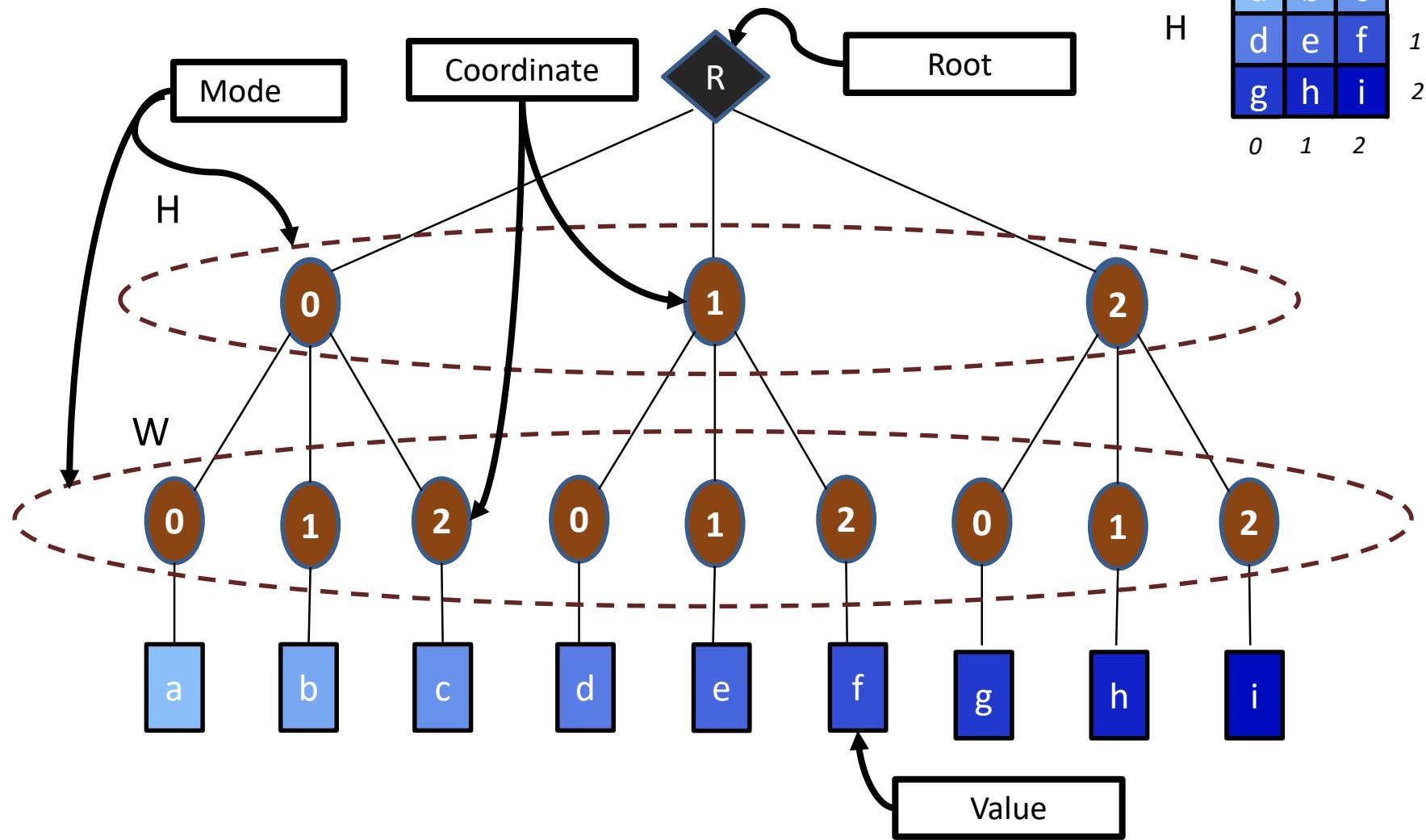


# Background: Compression

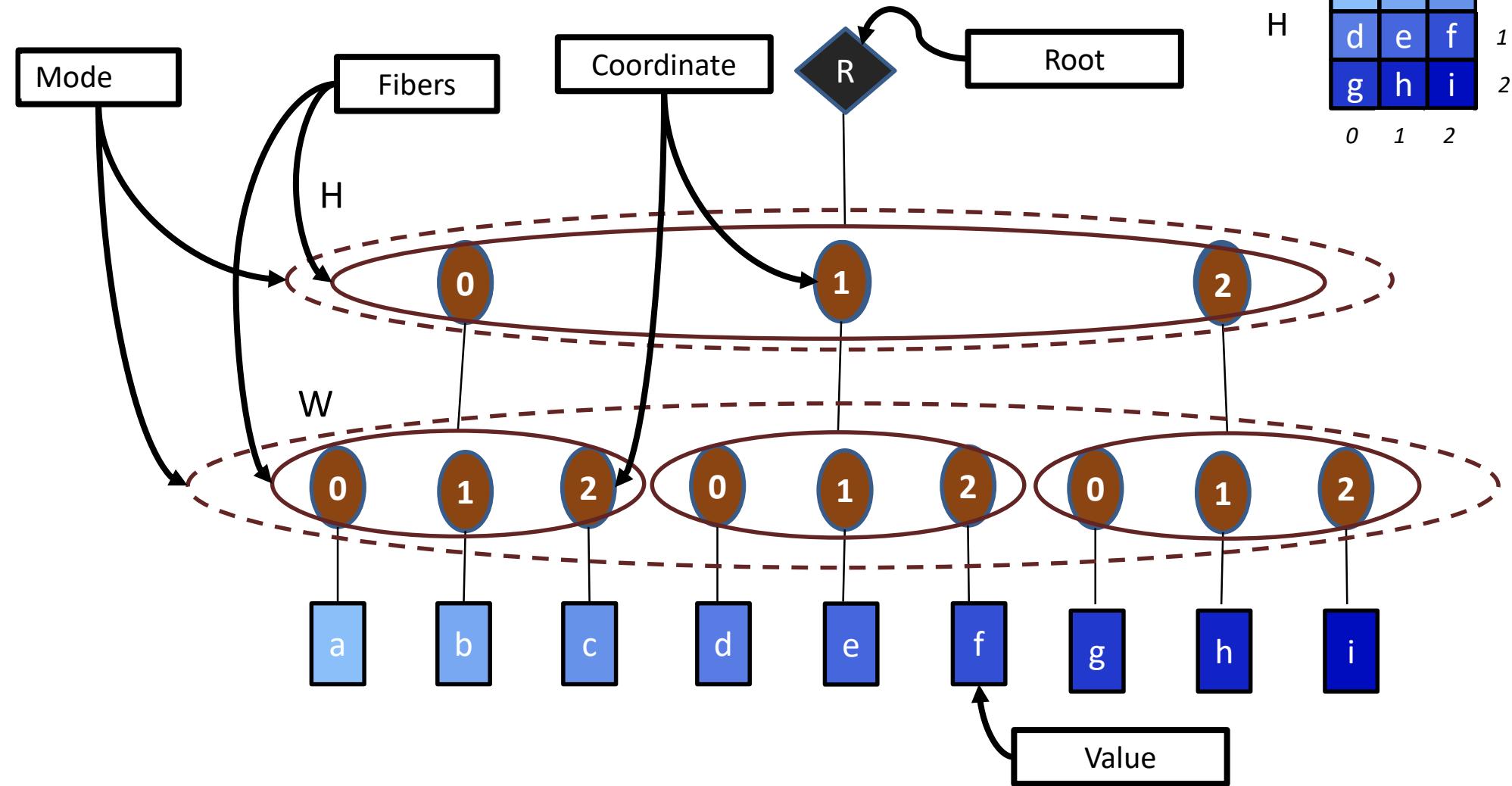


- Each rank (level) has a format attribute that are are ***data structure, representation, and format*** dependent: compressed, uncompressed, etc.
- Uncompressed means the occupancy is ***proportional*** to the max coordinate
- Compressed data has ***metadata overhead*** required to define the scheme
  - Data in a compression array is accessed by ***position***, the ***memory address space***
- ***Sparsity*** refers to the ***data pattern***

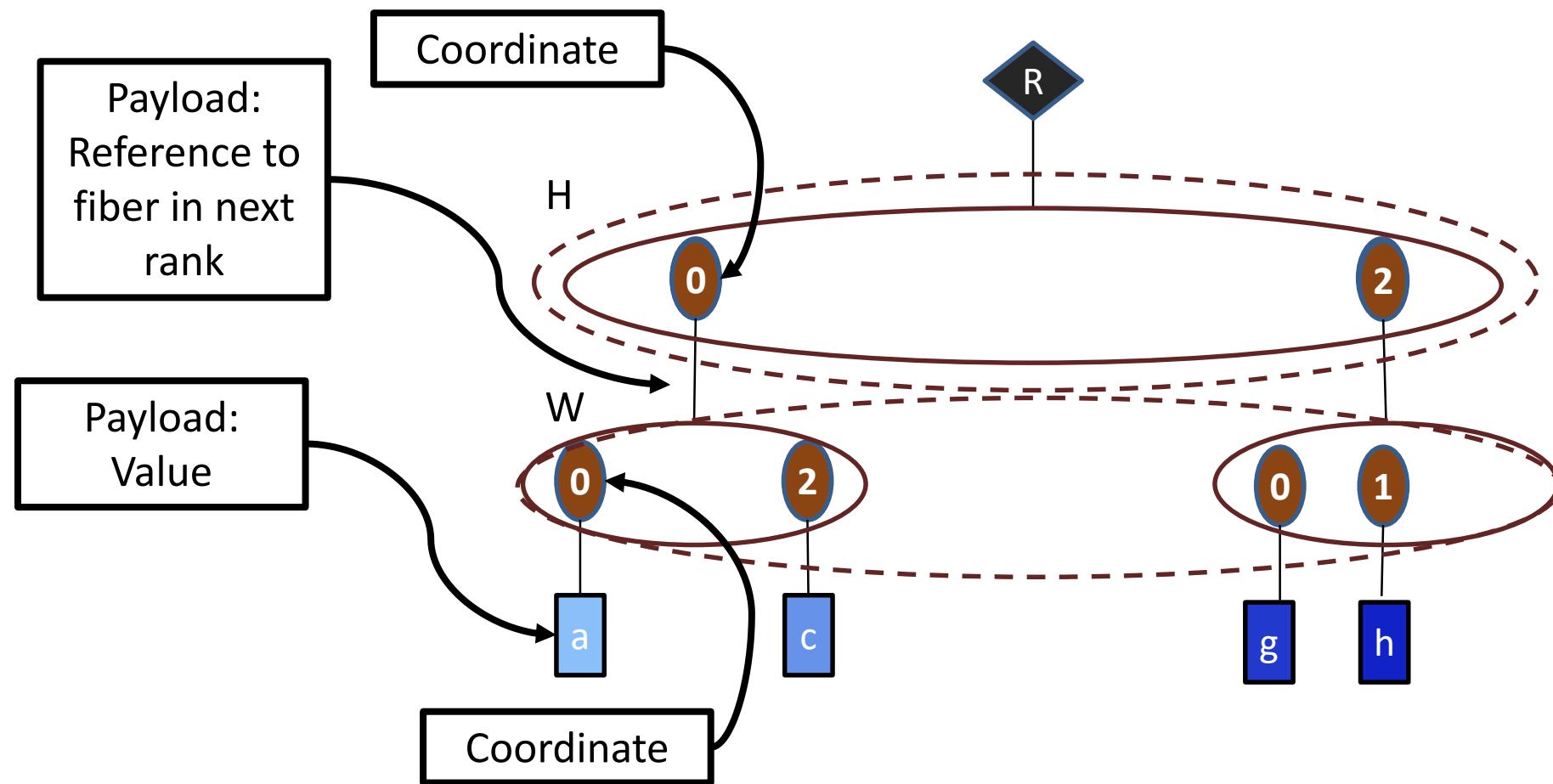
# Background: Tensors as Trees



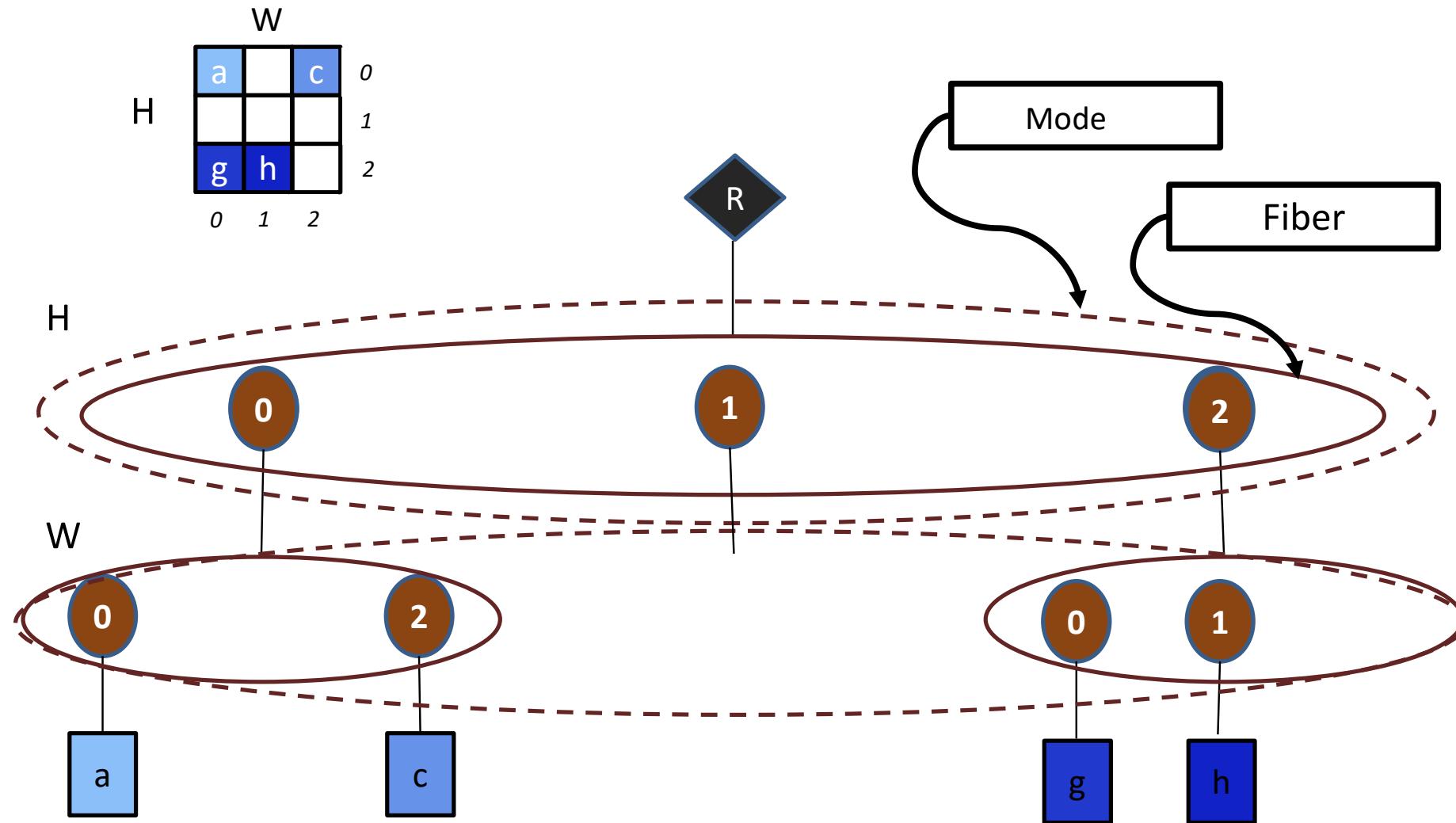
# Background: Fibertrees



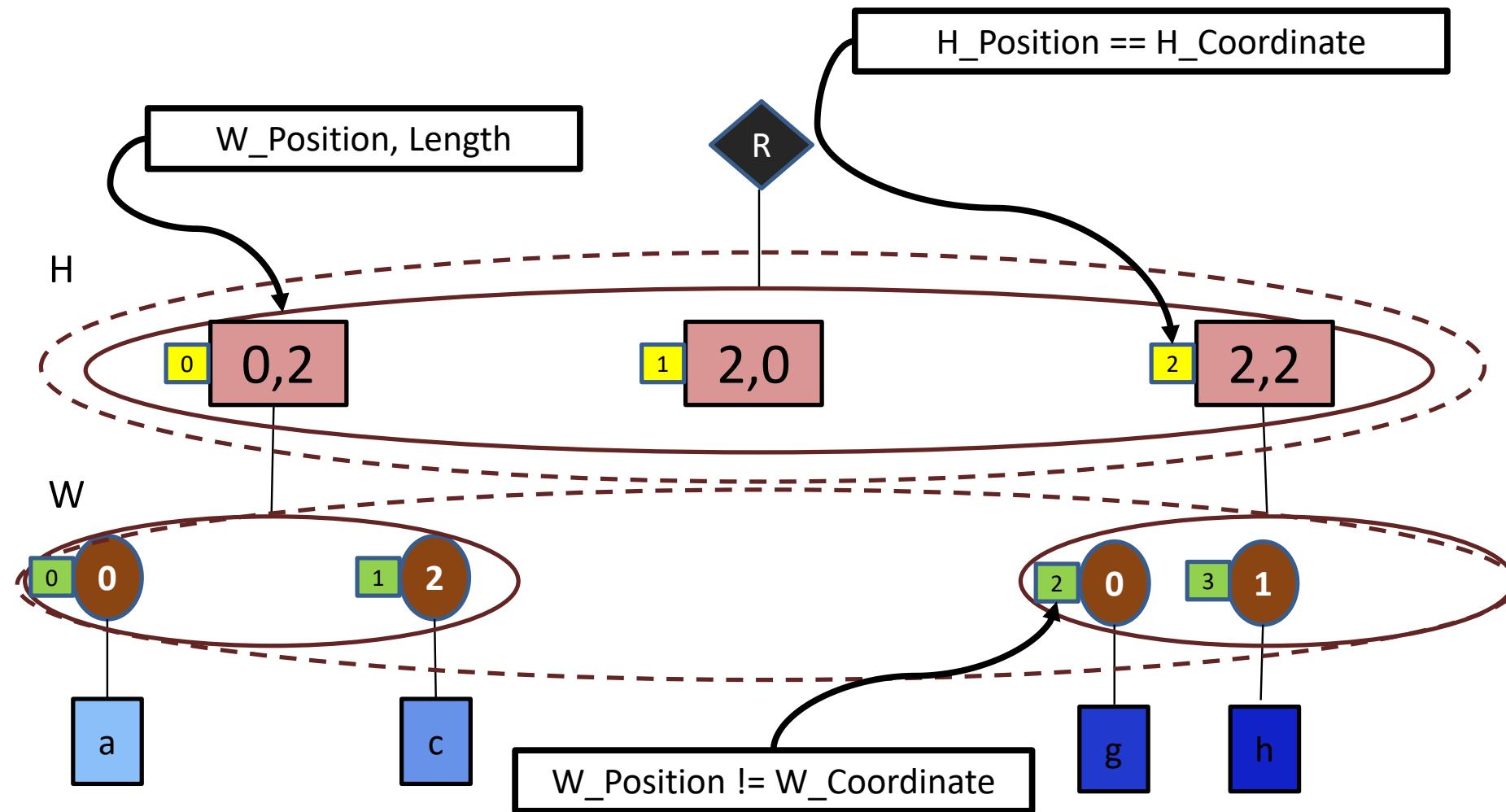
# Background: Fiber = (Coordinate, Payload)



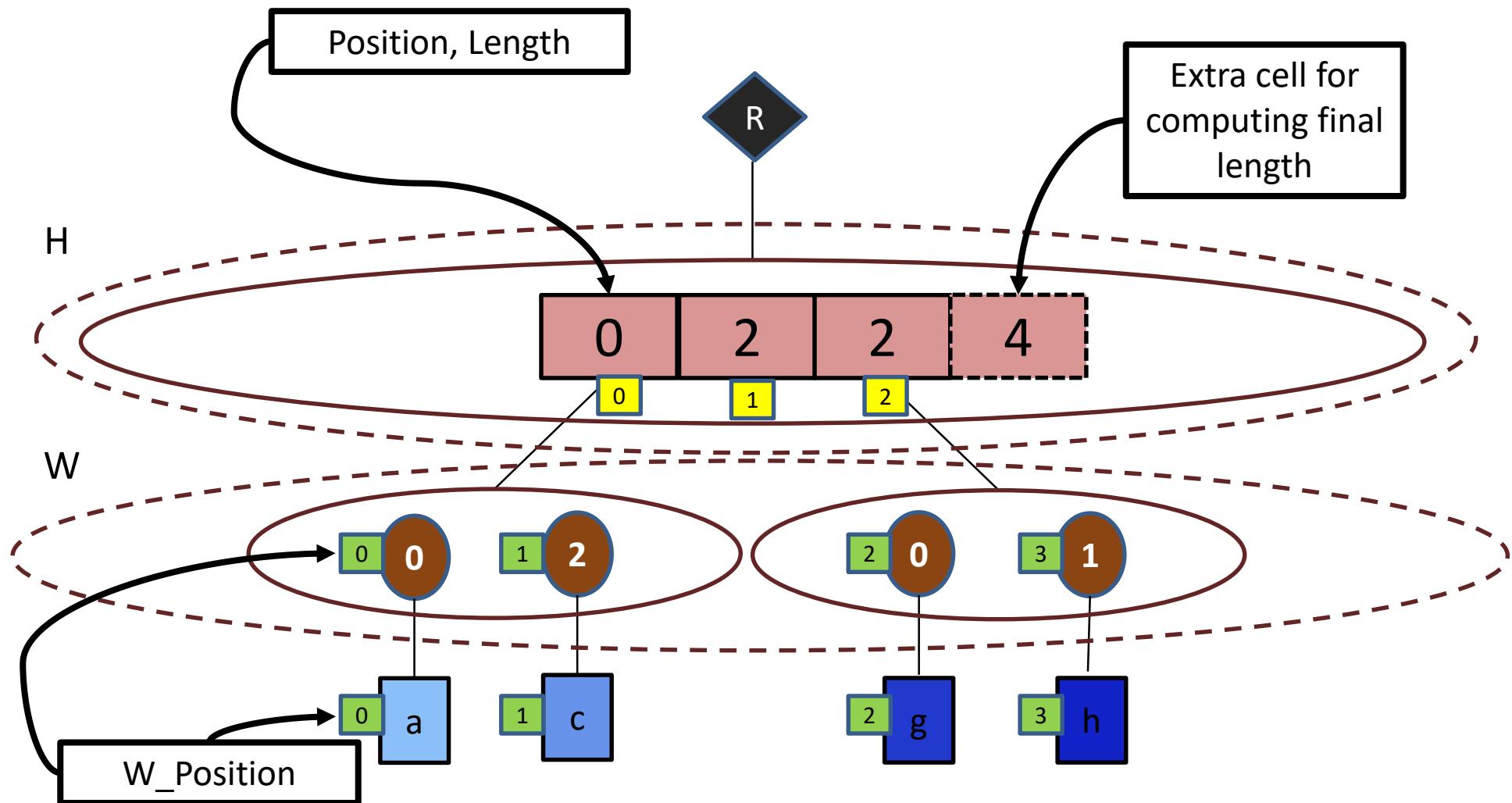
# Background: <Uncompressed, Compressed>



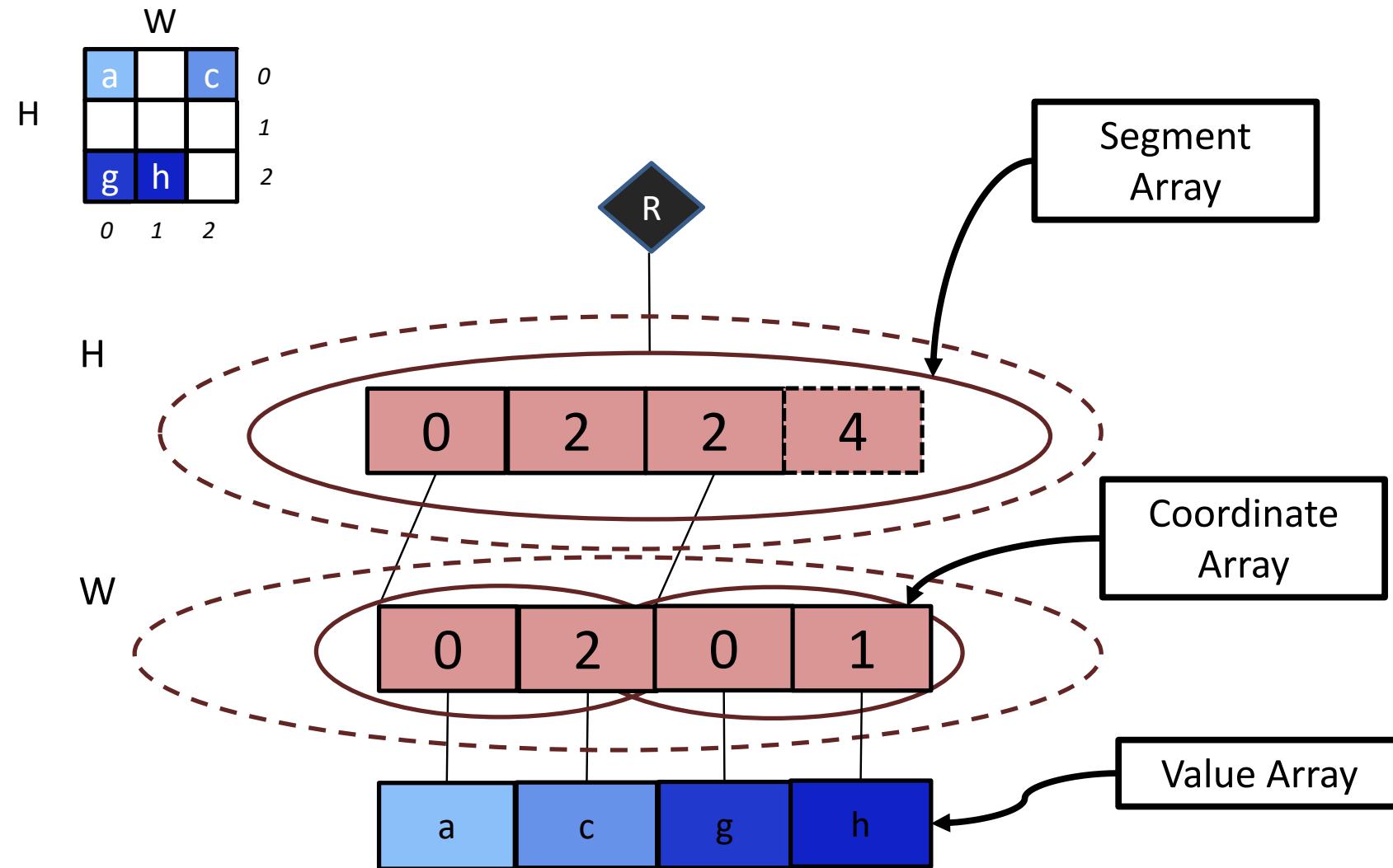
# Background: Uncompressed = (Position, Len)



# Background: <Position, Len> Array



# Background: Compressed Sparse Row



# Background: Compression Again

Compressed Sparse Row (CSR)										
Row Positions		<table border="1"><tr><td>0</td><td>1</td><td>3</td><td>4</td><td>5</td></tr></table>				0	1	3	4	5
0	1	3	4	5						
Col Coordinates		<table border="1"><tr><td>1</td><td>0</td><td>2</td><td>1</td><td>3</td></tr></table>				1	0	2	1	3
1	0	2	1	3						
Values		<table border="1"><tr><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td></tr></table>				1	2	3	4	5
1	2	3	4	5						
Fill Value		<table border="1"><tr><td>0</td></tr></table>				0				
0										

Compressed Sparse Row (CSR)										
Row Positions		<table border="1"><tr><td>0</td><td>1</td><td>3</td><td>4</td><td>5</td></tr></table>				0	1	3	4	5
0	1	3	4	5						
Col Coordinates		<table border="1"><tr><td>1</td><td>0</td><td>2</td><td>1</td><td>3</td></tr></table>				1	0	2	1	3
1	0	2	1	3						
Values		<table border="1"><tr><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td></tr></table>				1	2	3	4	5
1	2	3	4	5						
Fill Value		<table border="1"><tr><td>1</td></tr></table>				1				
1										

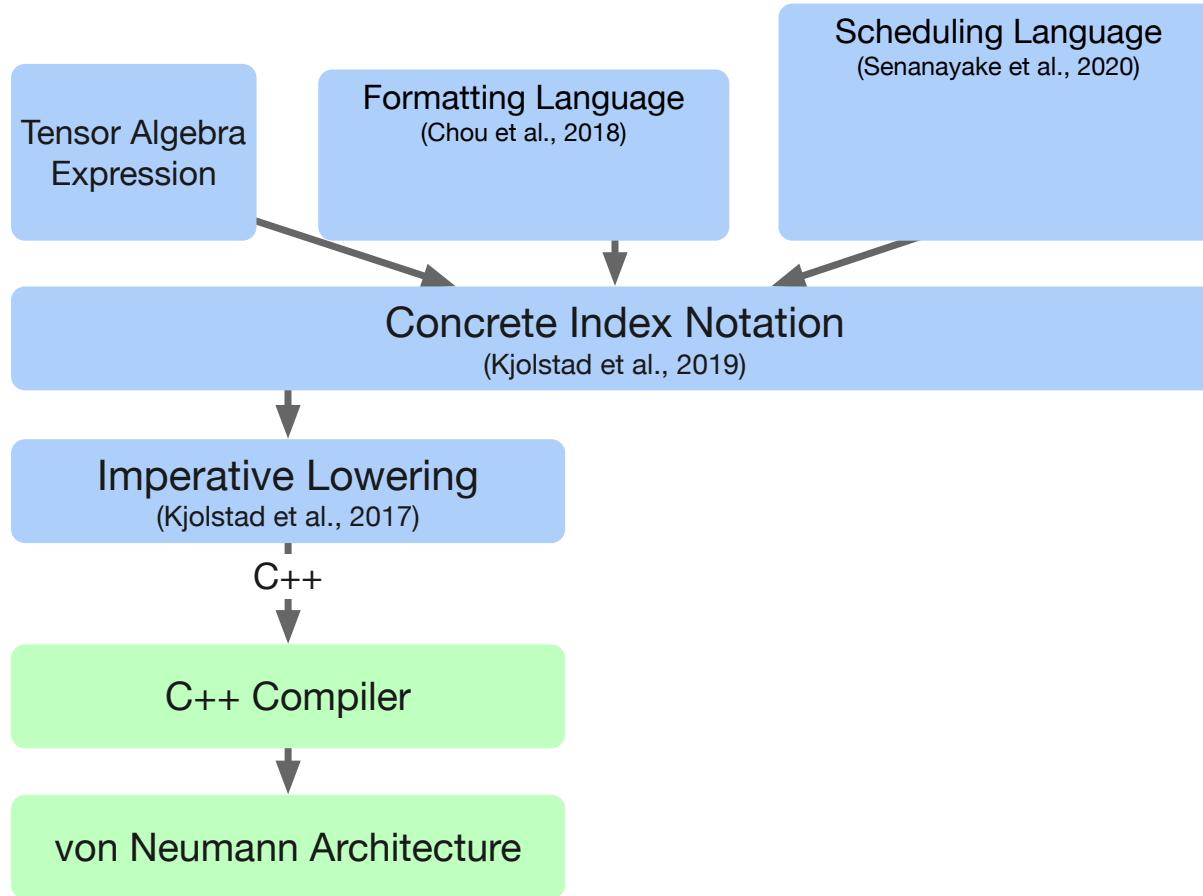
Coordinate (COO)										
Row Coordinates		<table border="1"><tr><td>0</td><td>1</td><td>1</td><td>2</td><td>3</td></tr></table>				0	1	1	2	3
0	1	1	2	3						
Col Coordinates		<table border="1"><tr><td>1</td><td>0</td><td>2</td><td>1</td><td>3</td></tr></table>				1	0	2	1	3
1	0	2	1	3						
Values		<table border="1"><tr><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td></tr></table>				1	2	3	4	5
1	2	3	4	5						
Fill Value		<table border="1"><tr><td>1</td></tr></table>				1				
1										

- Each rank (level) has a format attribute that are are ***data structure, representation, and format*** dependent: compressed, uncompressed, etc.
- Uncompressed means the occupancy is ***proportional*** to the max coordinate
- Compressed data has ***metadata overhead*** required to define the scheme
  - Data in a compression array is accessed by ***position***, the ***memory address space***
- ***Sparsity*** refers to the ***data pattern***

## Legend

- Other
- TACO Literature

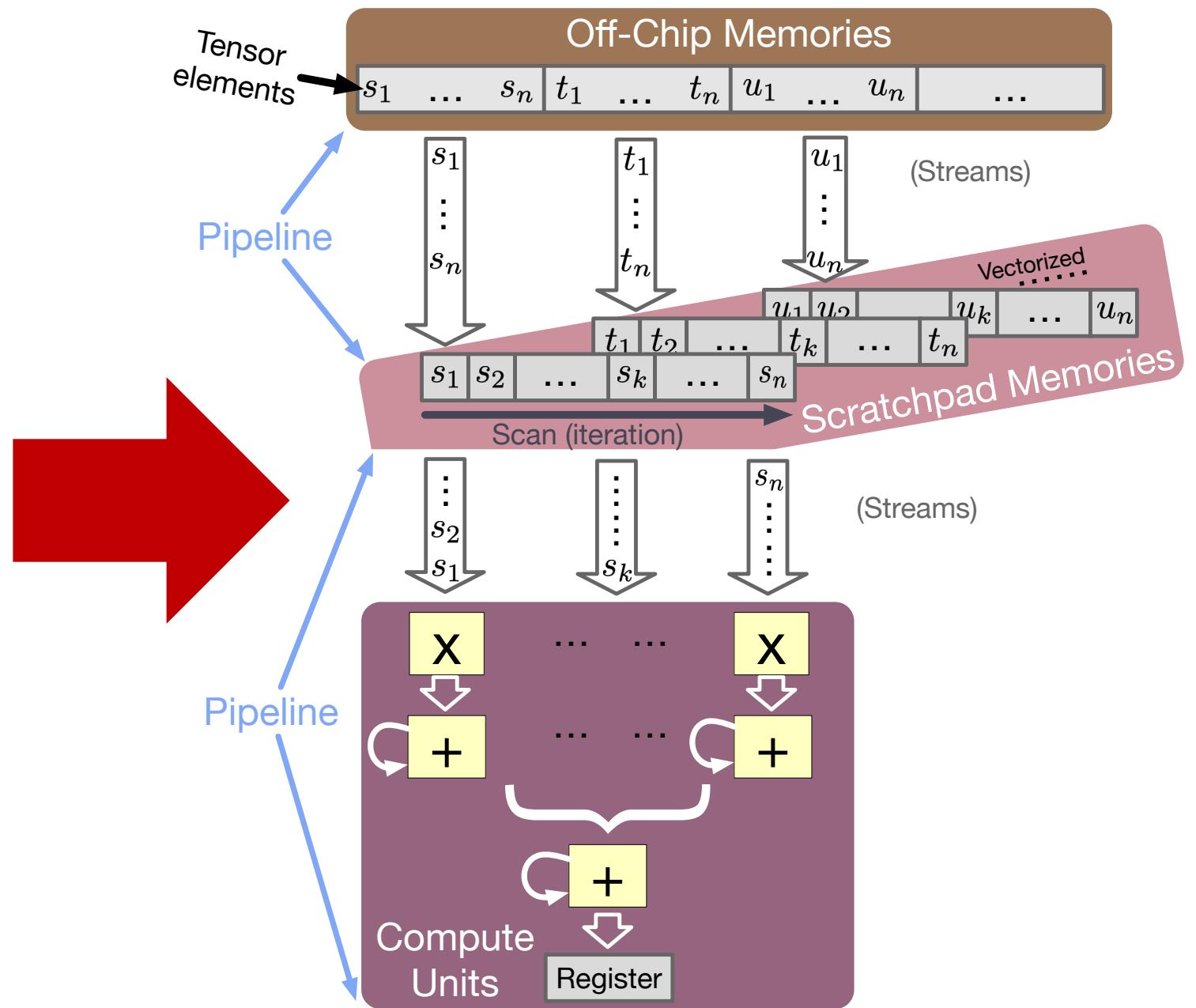
# TACO Overview



# Main Challenge

Imperative Representation  
 $\forall_i(x_i = b_i * c_i)$

Concrete Index Notation



# Goal

Create a *logical abstraction* that can represent *sparse dataflow* applications

Attempt to describe the *entire* space of sparse tensor algebra kernels  
spatially

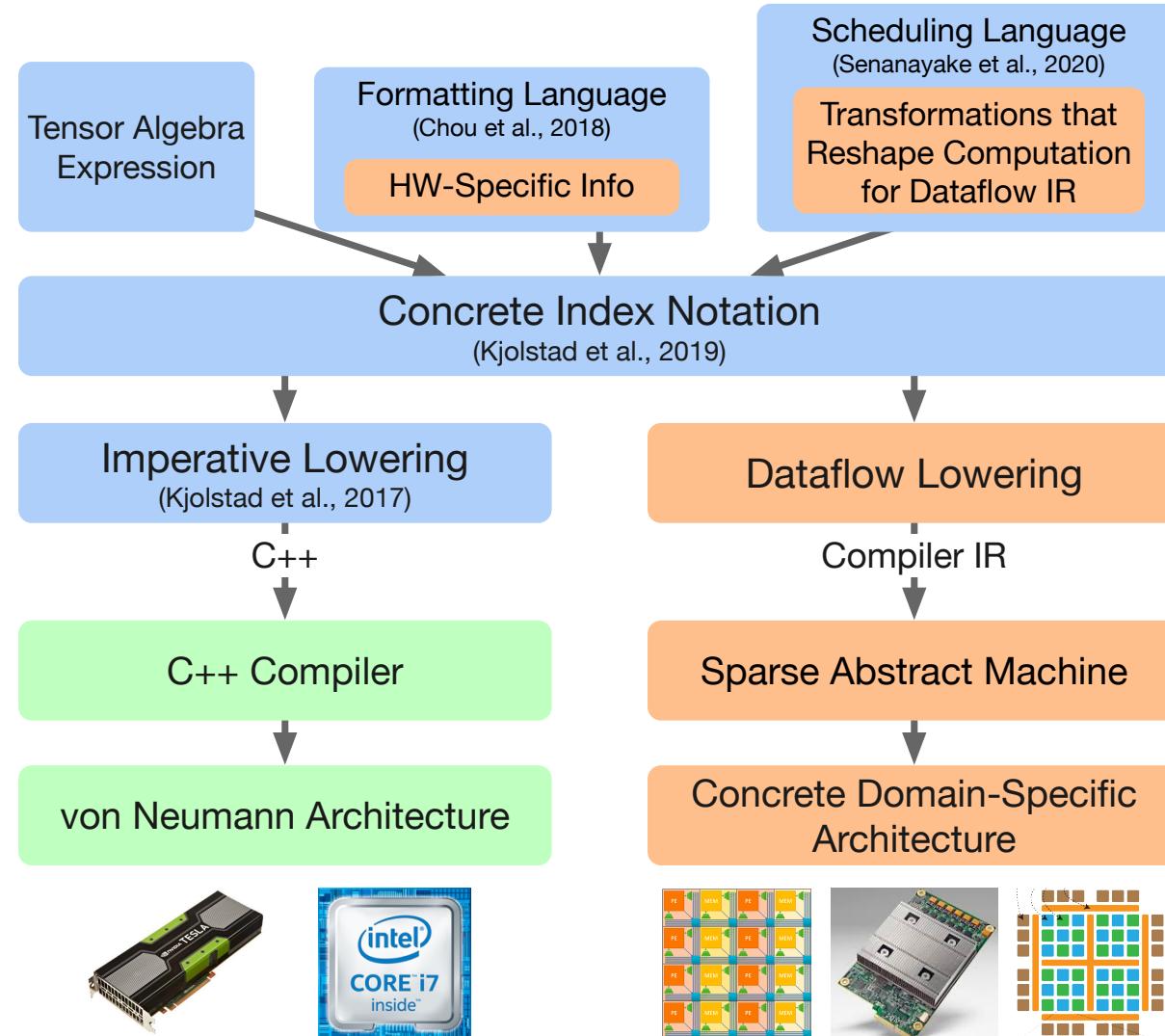
Ultimately, use this representation within the sparse compiler as an *IR*

Develop *efficient hardware implementations* of the logical blocks in  
conjunction

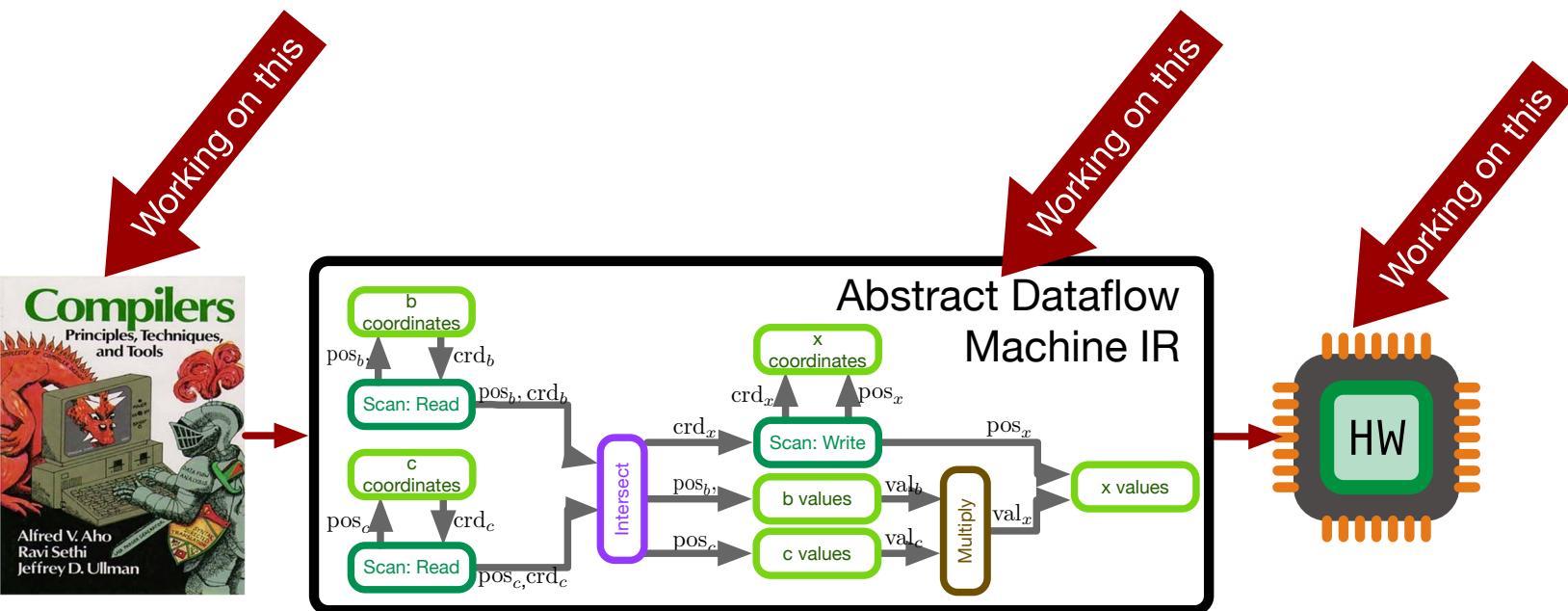
## Legend

- Sparse Dataflow Compiler
- Other
- TACO Literature

# Sparse Dataflow Compiler Overview



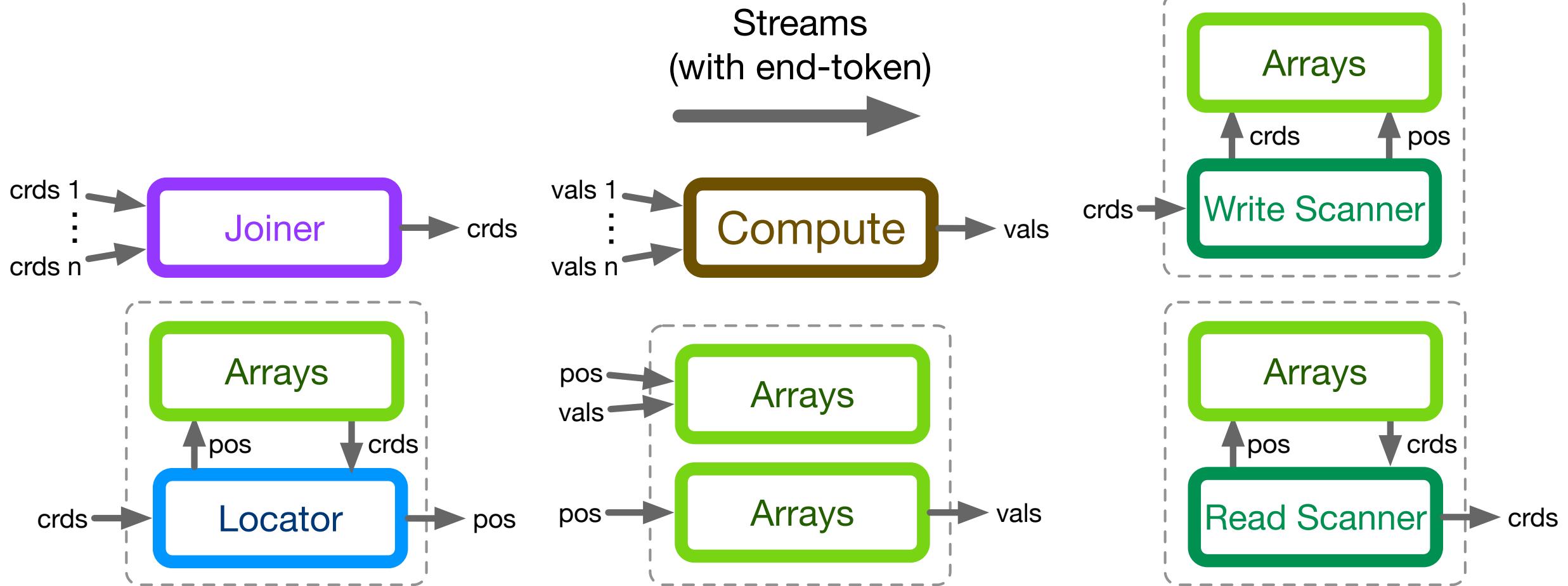
# Current Work



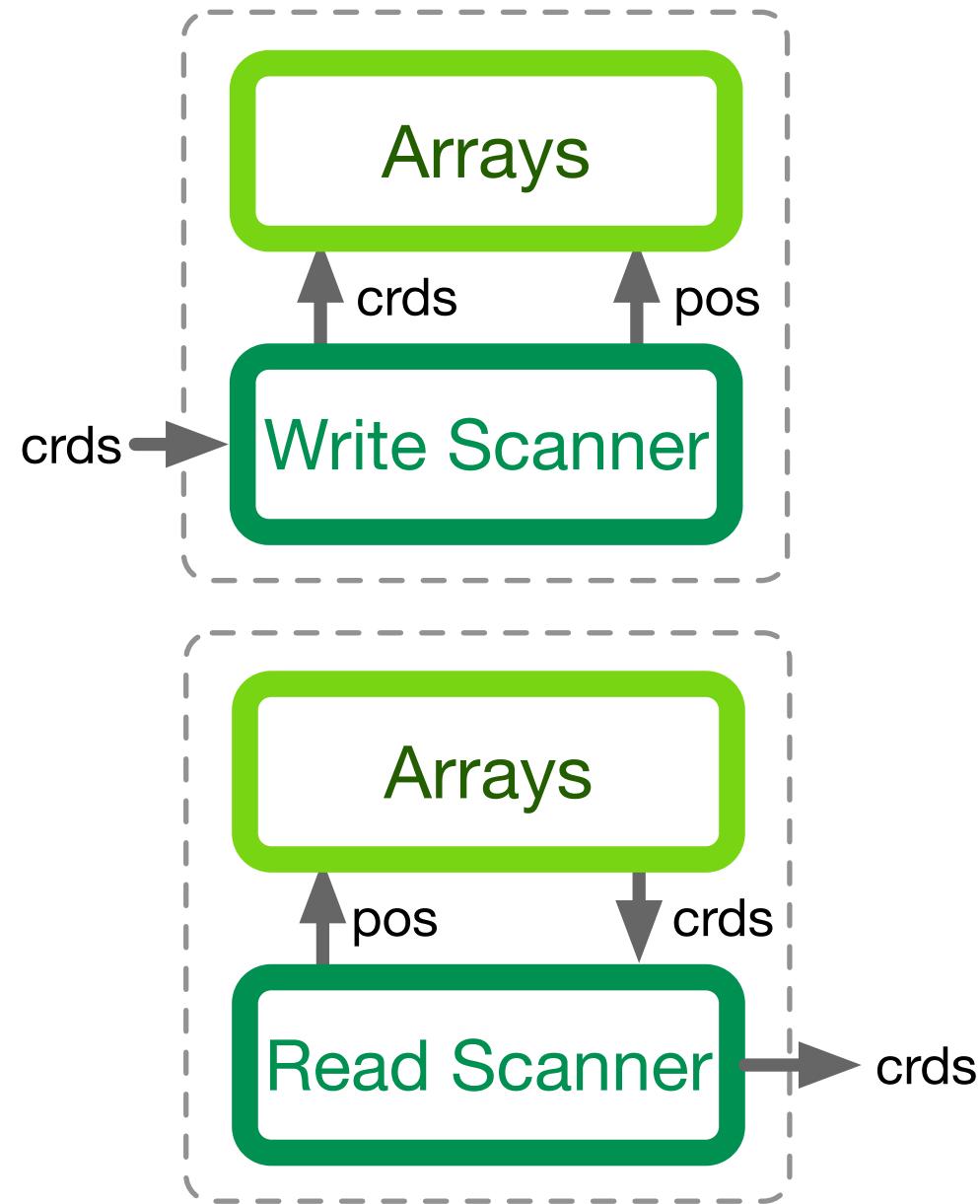
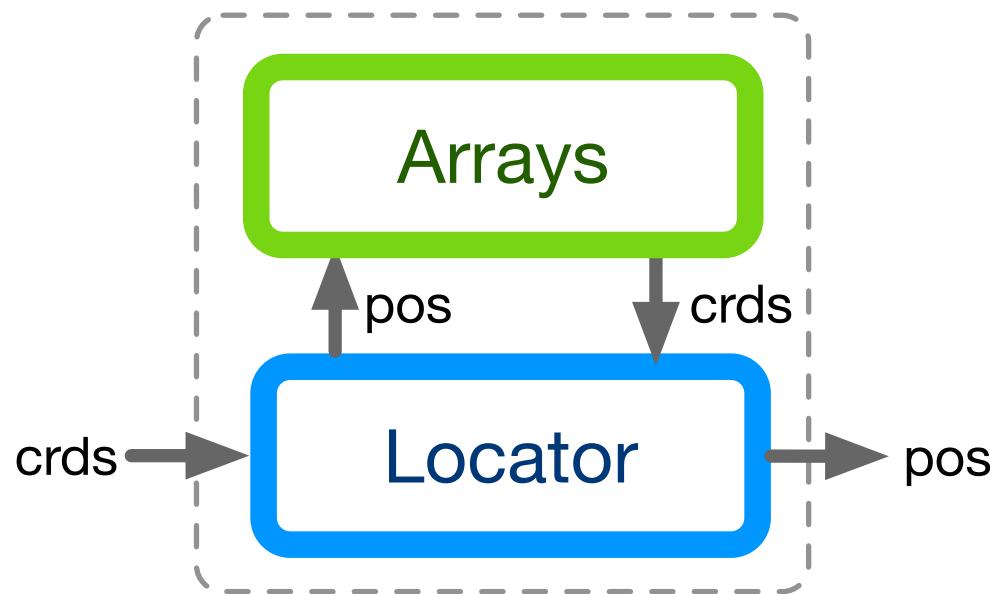
# Approach

- Draw sparse kernels spatially
- Start Simple, and increase in degrees of freedom
  - Started out with  $\text{SpV} * \text{V}$
  - Change vector compression formats:  $\text{SpV} * \text{SpV}$ ,  $\text{V} * \text{V}$
  - Change parallelization (x2):  $\text{SpV} * \text{V}$ ,  $\text{SpV} * \text{SpV}$
  - Change operator:  $\text{SpV} + \text{SpV}$
  - Increase dimensions and add reductions:  $\text{SpM} * \text{V}$
- Enumerate dataflow abstraction options to make granularity and abstraction decisions

# Primitives

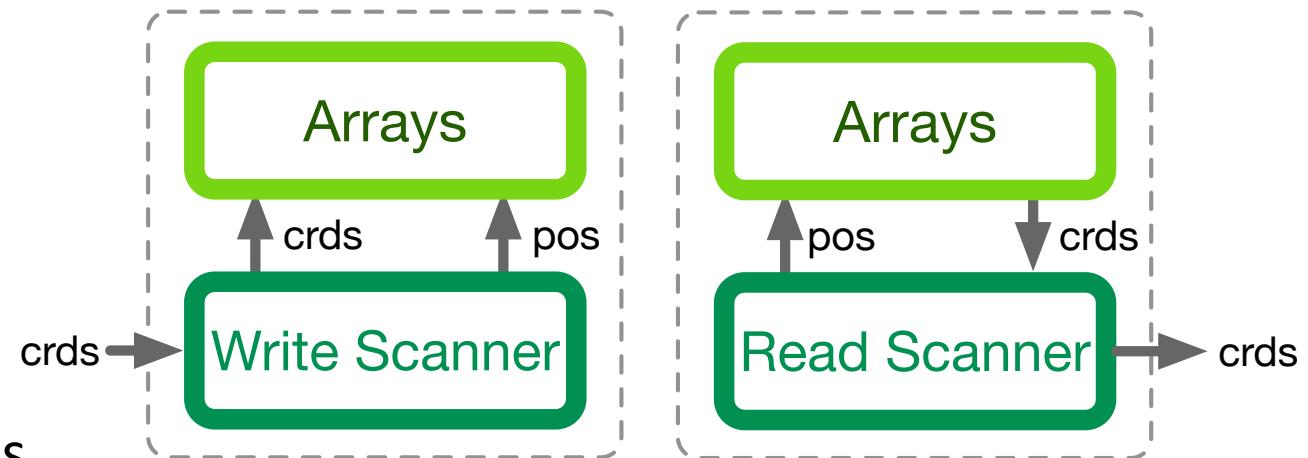


# Primitives



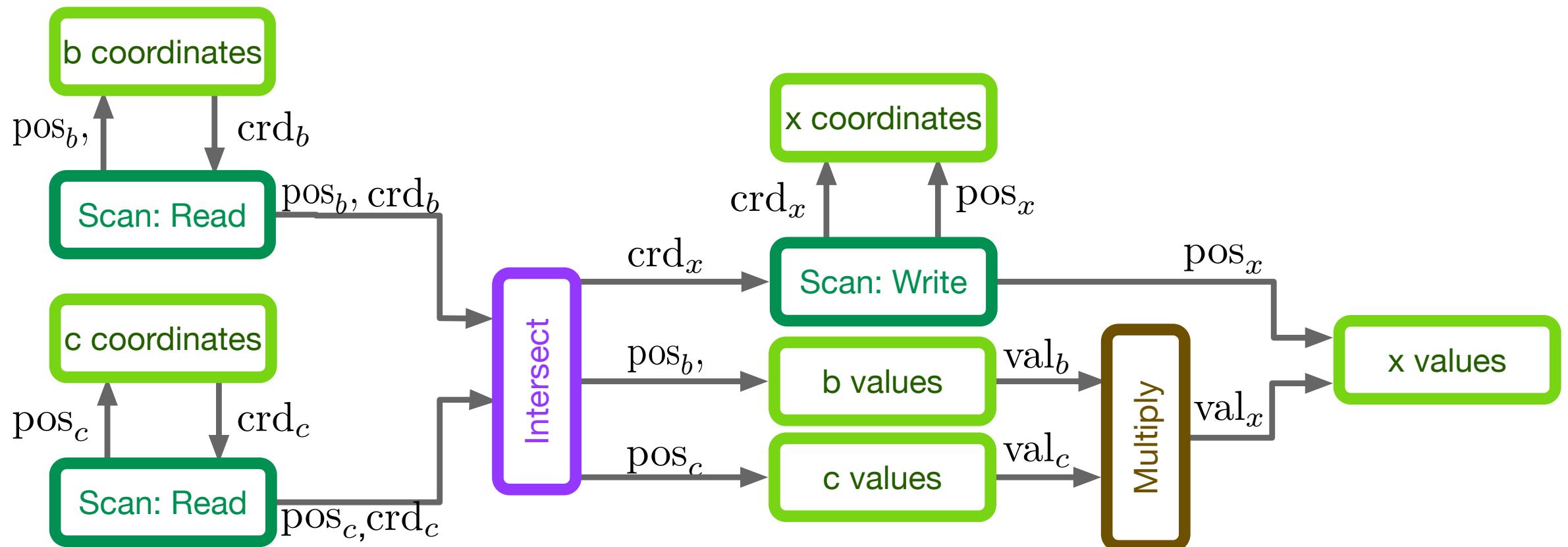
# Typing

- Define Interfaces with types
  - All streams are of valued elements
  - Positions and Coordinates are special types with properties
- Streams
  - Have implicit end-token that denotes the end of a rank
  - Could potentially generalize to streams of n-valued tuples (vectorization)
- Need a systematic way of checking type compatibility between modules



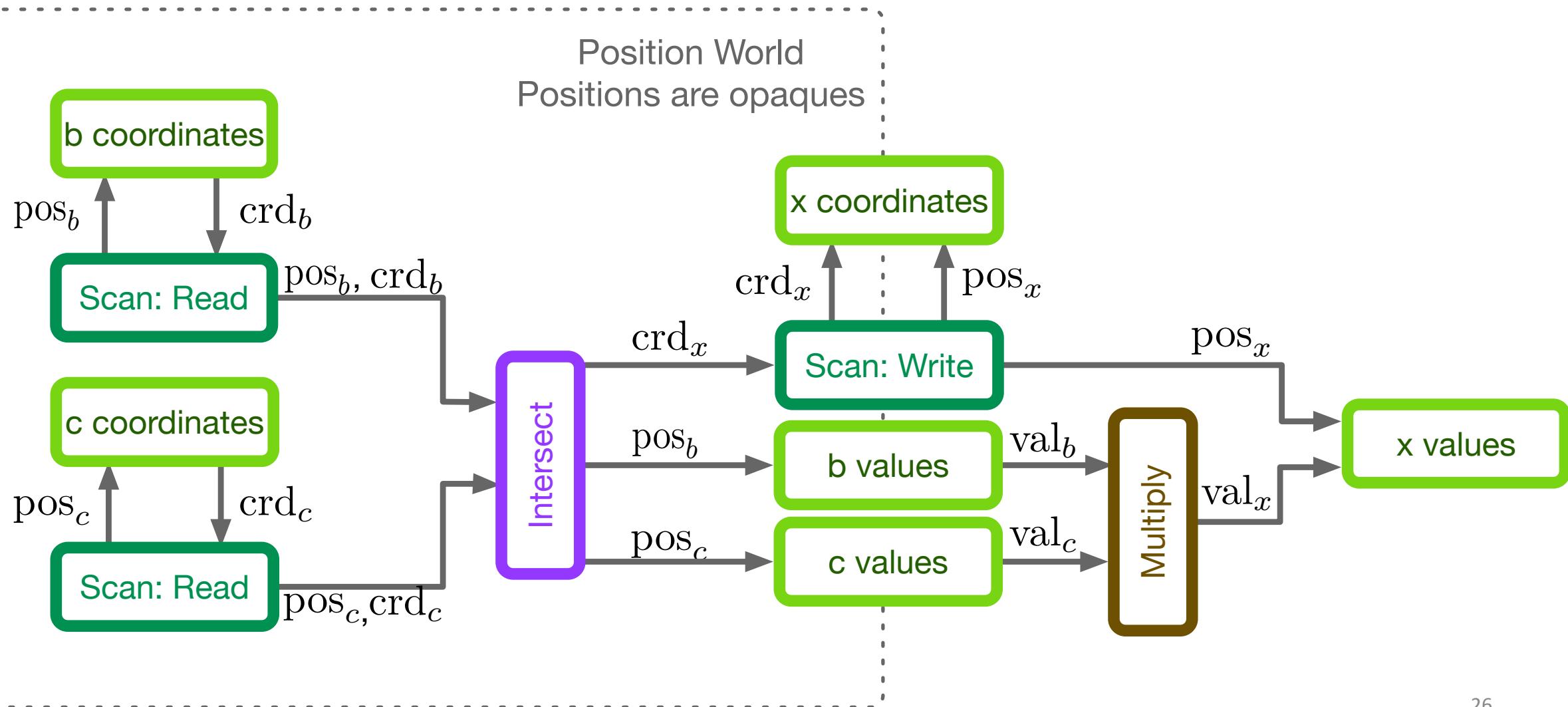
$$\text{SpV} = \text{SpV} * \text{SpV}$$

$$x_i = b_i * c_i$$



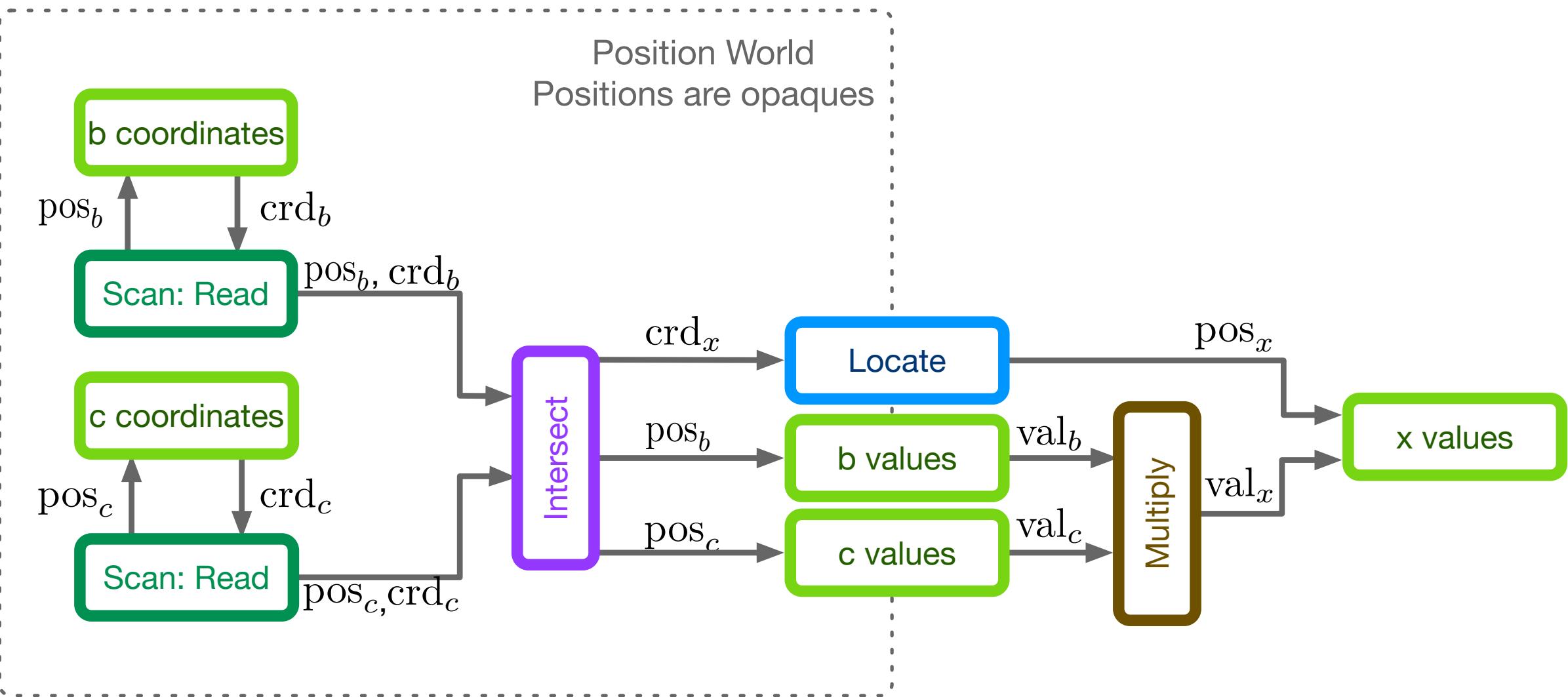
# SpV = SpV \* SpV Continued

$$x_i = b_i * c_i$$



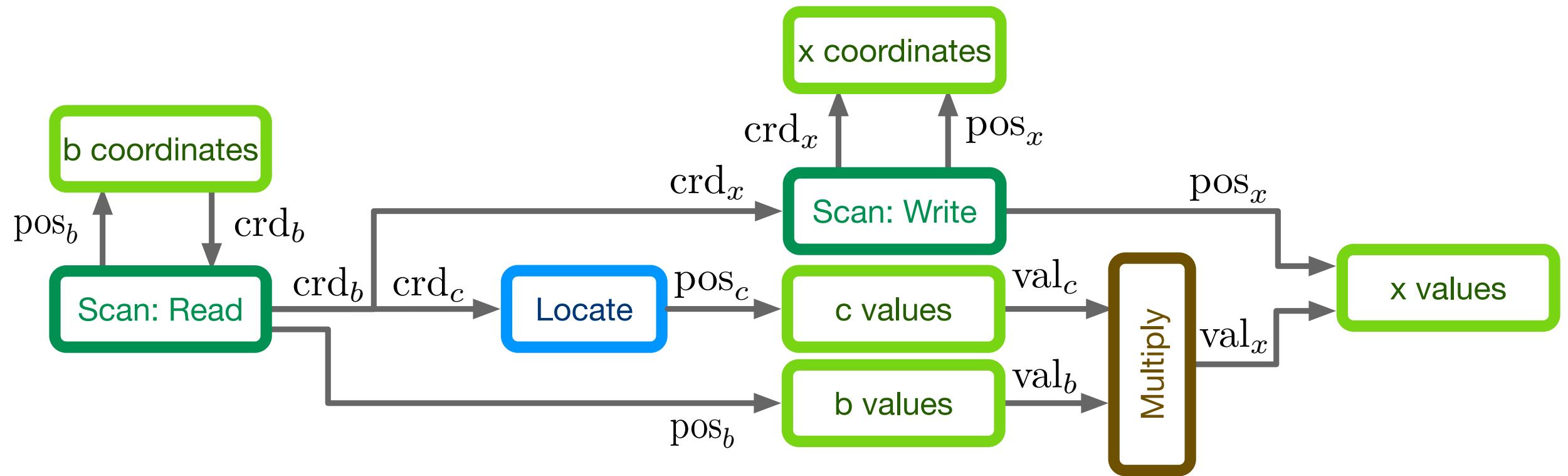
$$V = SpV * SpV$$

$$x_i = b_i * c_i$$



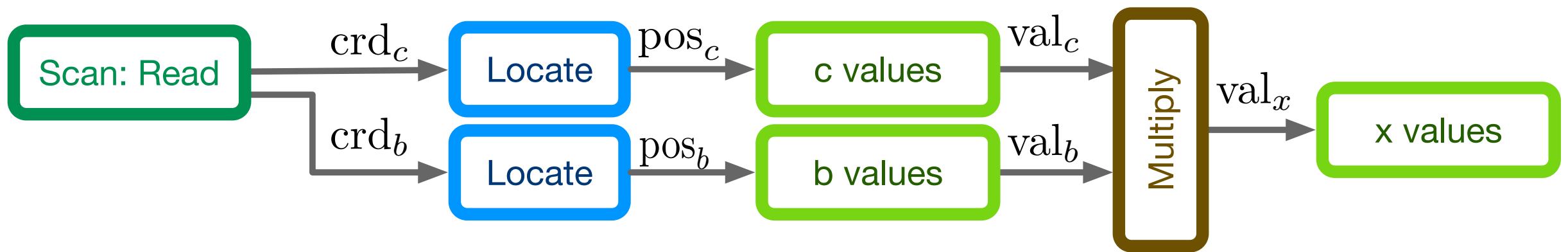
$$\text{SpV} = \text{SpV} * V$$

$$x_i = b_i * c_i$$



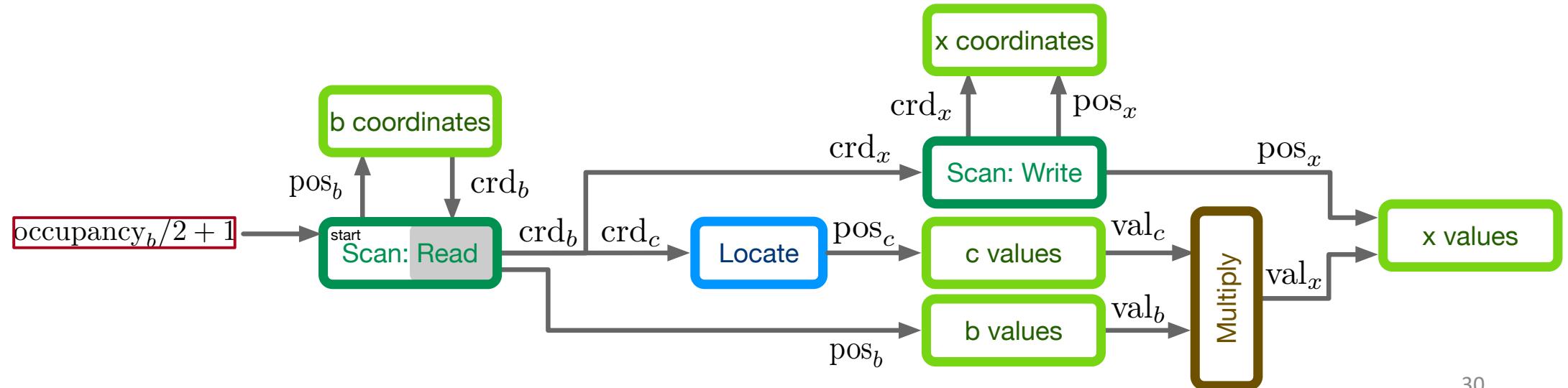
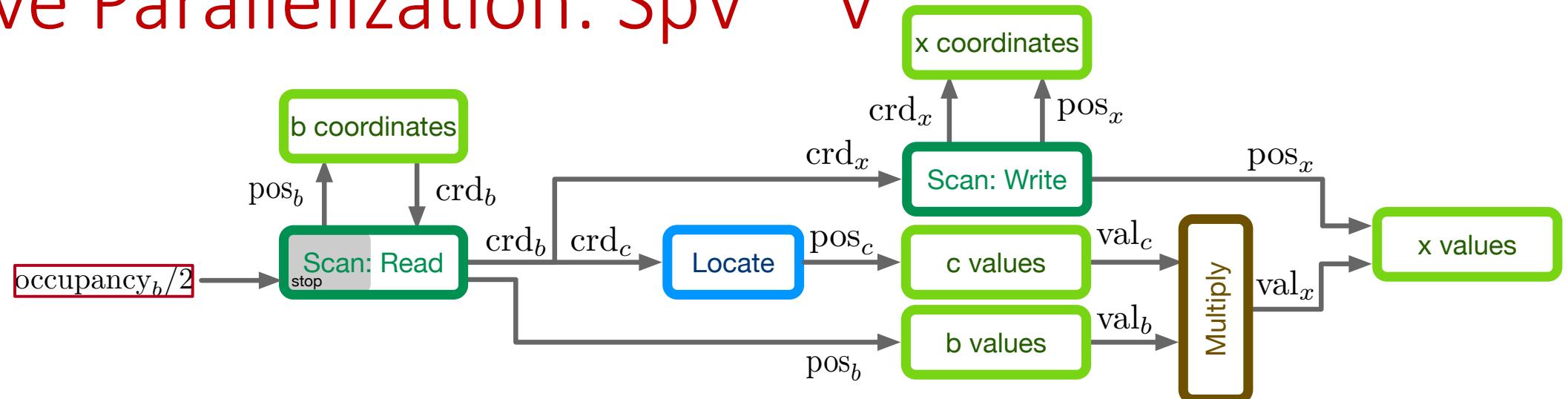
$$V = V * V$$

$$x_i = b_i * c_i$$



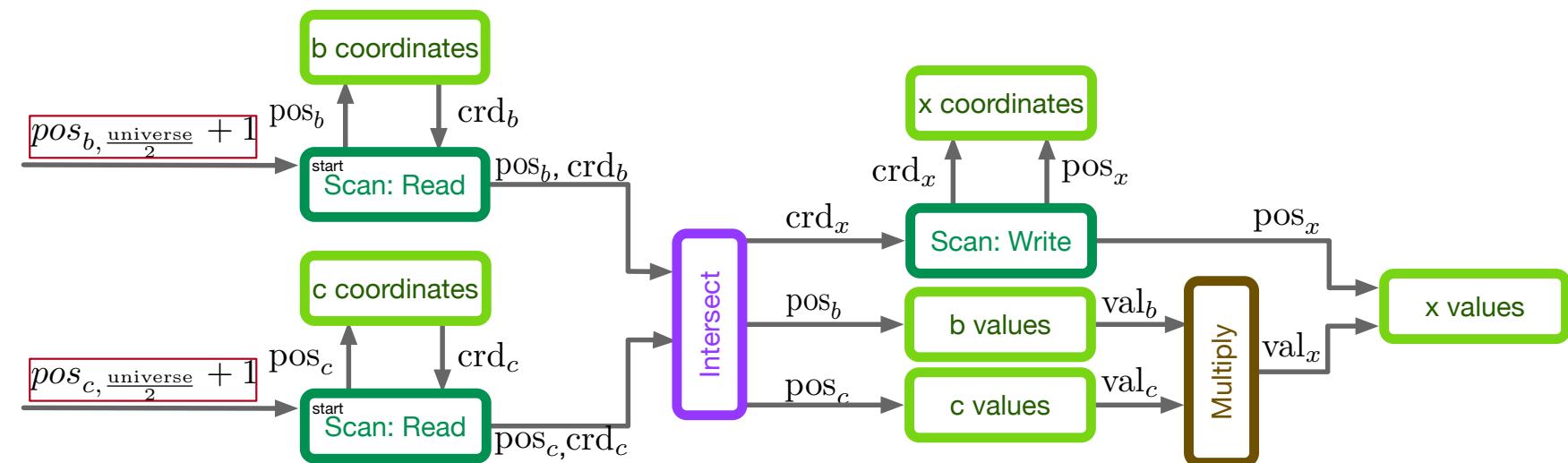
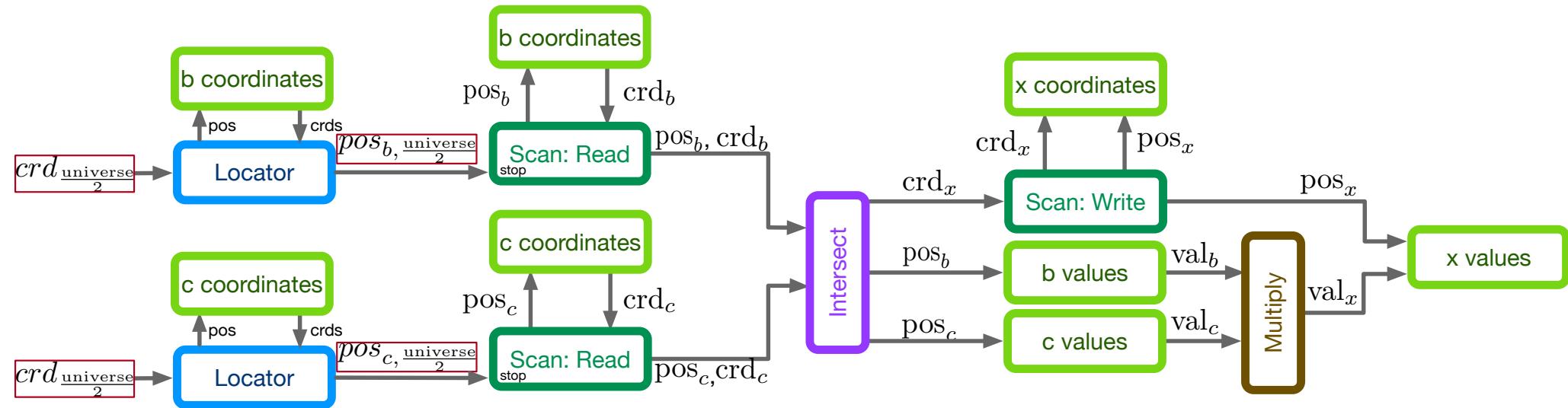
$$x_i = b_i * c_i$$

## Naïve Parallelization: SpV \* V



$$x_i = b_i * c_i$$

# Parallelization Coord Space: SpV \* SpV



# Parallelization Things to Think About

- Parallelize by other split methods (more complicated)
  - $\text{SpV} * \text{SpV} \rightarrow$  Cut one vector by occupancy/2 and locate into other input
  - Split both by occupancy/2 and then take the average of the 2 coordinates
  - Goal is to get the most balanced workload for both halves
  - Can always find sparsity pattern that is pathological for each method
- Parallelizing spatially is basically duplicating the same graph
- Missing compaction step of the resulting tensor at the end for readout/future use in cascaded kernels

# Next Steps

- Look into combining our current abstraction to include the fibertree abstraction in array storage
- Expanding kernels to other degrees of freedom: SpMV
- Analyzing more parallelization techniques (including vectorized streams)
- A compiler-prompted sparse hardware acceleration structure for coordinate skipping on intersections

# Preliminary Conclusions

- Working on a logical abstraction for sparse dataflow kernels
- Fibertree + fiber argument stream abstraction
- Parallelization is simple in hardware
- Typing and interfaces are important
- Expanding kernels one degree of freedom (complexity) at a time