# Stream-Dataflow Acceleration
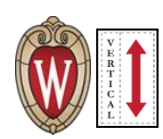
Tony Nowatzki[+], **Vinay Gangadhar***,

Newsha Ardalani*, Karu Sankaralingam*

**44th ISCA, Toronto, ON, Canada**

**Accelerator Session (6A-4)**

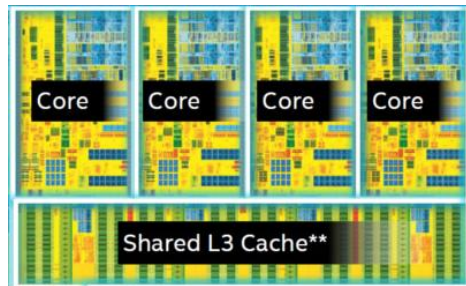**Tuesday June 27th, 2017**

*University of Wisconsin-Madison

[+]University of California, Los Angeles

# Era of Specialization

**Traditional Multicore**



*Application domain specialization*

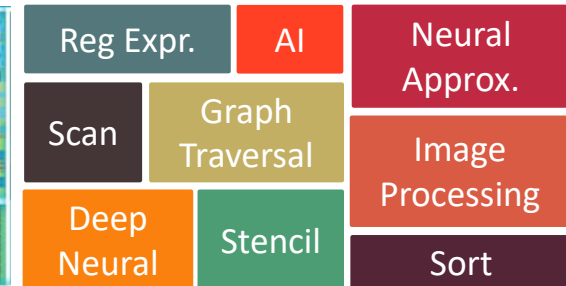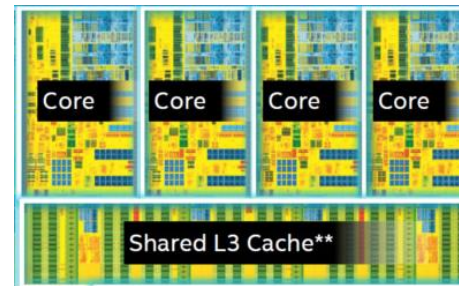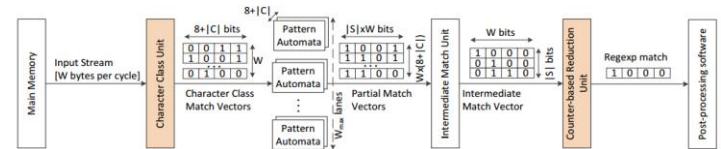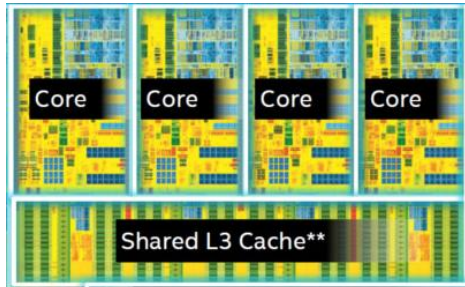| | | |
|---|---|---|
| Reg Expr. | AI | Neural Approx. |
| Scan | Graph Traversal | Image Processing |
| Deep Neural | Stencil | Sort |

# Era of Specialization

**Traditional Multicore**



*Application domain specialization*

# Era of Specialization

**Traditional Multicore**



*Application domain specialization*

Reg Expr. | AI | Neural Approx.
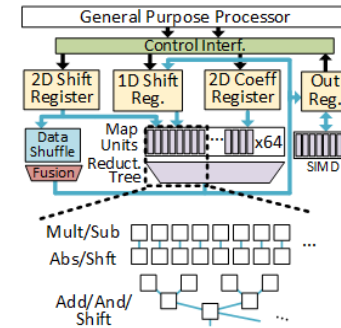Scan | Graph Traversal | Image Processing
Deep Neural | Stencil | Sort

# Era of Specialization

**Traditional Multicore**



Core | Core | Core | Core

Shared L3 Cache**

*Application domain specialization*

Core | Core | Core | Core

Shared L3 Cache**

| Reg Expr. | AI |
| Scan | Graph Traversal |
| Deep Neural | Stencil |

| Neural Approx. |
| Image Processing |
| Sort |



**Movidius Myriad VPU**



**NVIDIA DGX-1 AI Accelerator**



**Catapult FPGA Accelerator**



Deephi Tech CNN Processor

*Processor for convolution neural network*
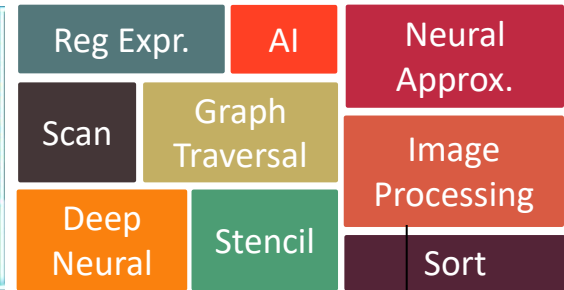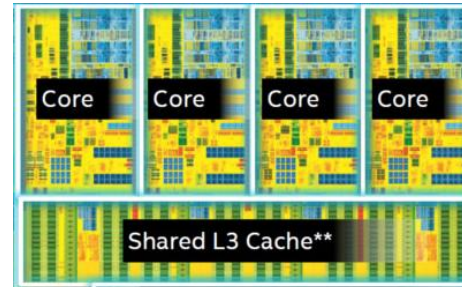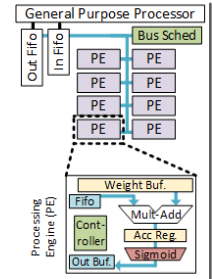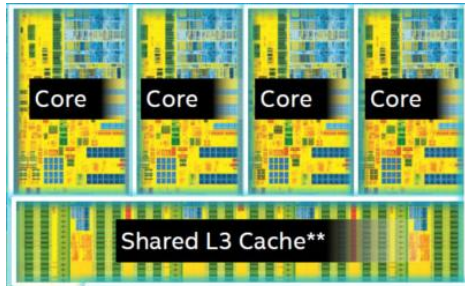
CNN processor principal architecture

# Era of Specialization

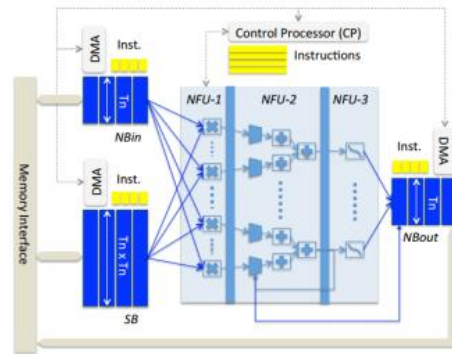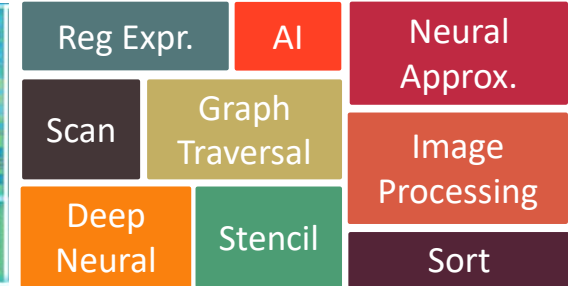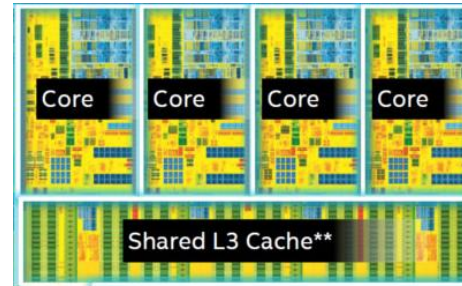**Traditional Multicore**



*Application domain specialization*

## Domain Specific Acceleration



| Reg Expr. | AI | Neural Approx. |
| Scan | Graph Traversal | Image Processing |
| Deep Neural | Stencil | Sort |

Fixed-function Accelerators for specific domain:
**Domain Specific Accelerators (DSAs)**

**+ High Efficiency**

10 – 100x
Performance/Power
or
Performance/Area

three orders of magnitude less energy than a state of the art software DBMS, while the performance-oriented design out-performs the same DBMS by **70X**

sor, the accelerator is **117X** faster, and it can reduce the total energy by **21X** The accelerator characteristics are obtained after layout at 65nm. Such a high throughput in
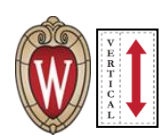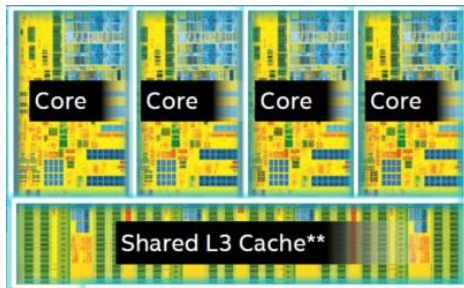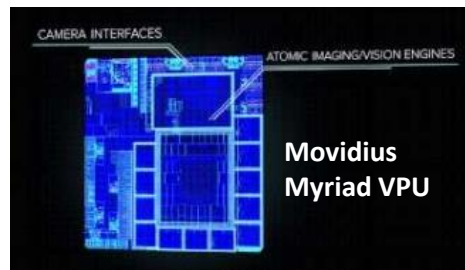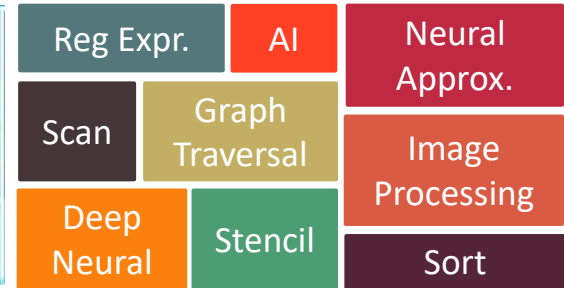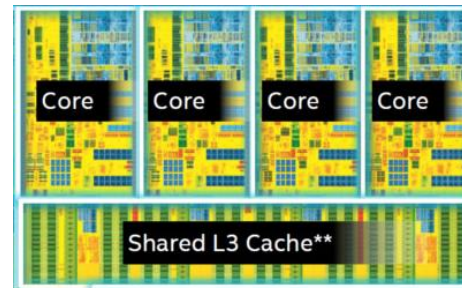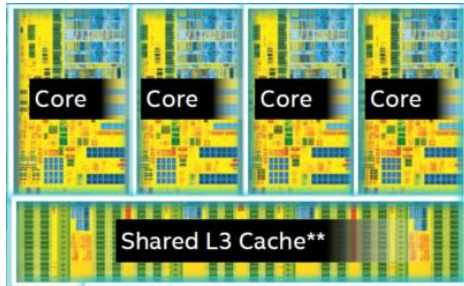
# Era of Specialization

**Traditional Multicore**



*Application domain specialization*

## Domain Specific Acceleration



| Reg Expr. | AI | Neural Approx. |
| Scan | Graph Traversal | Image Processing |
| Deep Neural | Stencil | Sort |

Fixed-function Accelerators for specific domain:
**Domain Specific Accelerators (DSAs)**

**+ High Efficiency**

10 – 100x
Performance/Power
or
Performance/Area

- Not programmable/re-configurable & Obsoletion prone

- Architecture, design, verification and fabrication cost

- Multi-DSA chip for "N" application domains → Area and cost inefficient

# The Universal Accelerator Dream...

Source:
Malitel Consulting

Deep Neural

Image Processing

Automated Driving

Compression

Regex Matching

Query Processing

Convert 100+ Accelerators

⬇

1 Programmable Accelerator Fabric

Standard programming and threading interface

**A generic programmable hardware accelerator matching the efficiency of Domain Specific Accelerators (DSAs) with an efficient hardware-software interface**

# Specialization Spectrum

*Generality*

→

**ASIC/ DSA**　　　　**GPGPU**　　**FPGA**　　**DSP**　　**SIMD**　　**GPP**

←

*Efficiency*

**(energy efficient computing)**

# Specialization Spectrum

*Generality*

*Efficiency*
**(energy efficient computing)**

**ASIC/ DSA**  **GPGPU**  **FPGA**  **DSP**  **SIMD**  **GPP**

Specialization Principles

*General Set of Micro-Architectural Mechanisms*

# Specialization Spectrum

*Generality*

**ASIC/ DSA**

**GPGPU    FPGA    DSP    SIMD    GPP**

*Efficiency*
**(energy efficient computing)**

Specialization Principles

Programmability / Re-configurability Features

*General Set of Micro-Architectural Mechanisms*

*Architecture with Flexible Hardware-Software Programming Interface*

# Specialization Spectrum

*Generality*

**ASIC/ DSA**

**GPGPU    FPGA    DSP    SIMD    GPP**

*Efficiency*
**(energy efficient computing)**

Specialization Principles

Programmability / Re-configurability Features

*General Set of Micro-Architectural Mechanisms*

**+**

*Architecture with Flexible Hardware-Software Programming Interface*

**Programmable Hardware Accelerator**

# Specialization Spectrum

*Generality*

**ASIC/ DSA**

**GPGPU    FPGA    DSP    SIMD    GPP**

*Efficiency*
**(energy efficient computing)**

Specialization Principles

Programmability / Re-configurability Features

*General Set of Micro-Architectural Mechanisms*

**+**

*Architecture with Flexible Hardware-Software Programming Interface*

**Programmable Hardware Accelerator**

Efficiency close to DSAs/ASICs

Trivial adaptation of new algorithms/applications

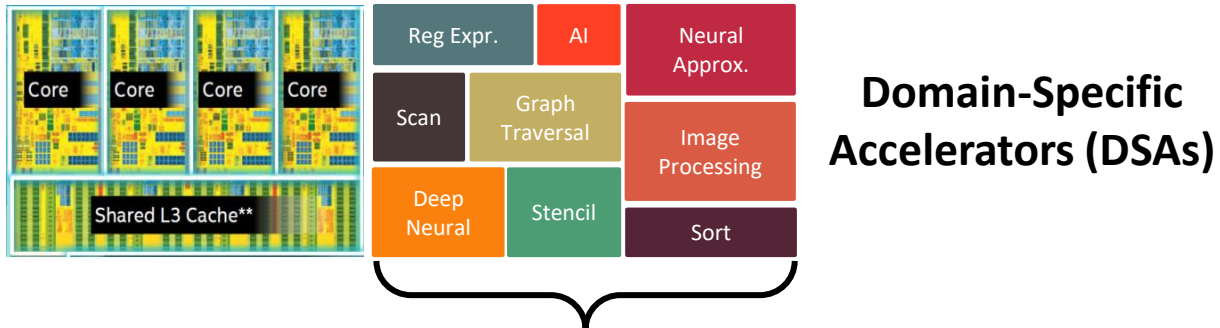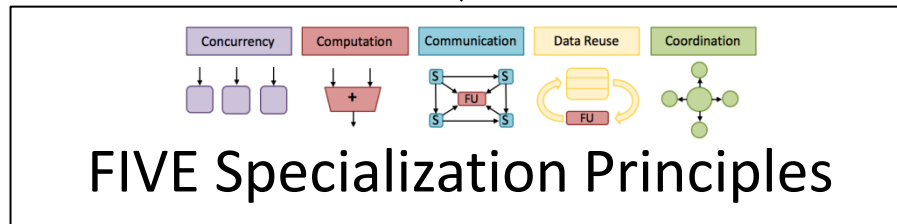Retain programmability

# Background Work*

*IEEE Micro Top-Picks 2017: *Domain Specialization is Generally Unnecessary for Accelerators*



**Domain-Specific Accelerators (DSAs)**
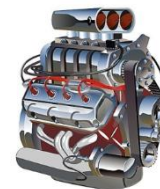
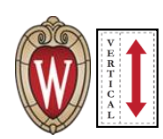Commonality in DSAs ?

FIVE Specialization Principles
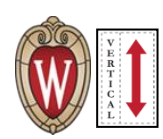
Micro-Architectural Mechanisms

*Programmable Hardware Accelerator Architecture*

# Our Work:
# Stream-Dataflow Acceleration

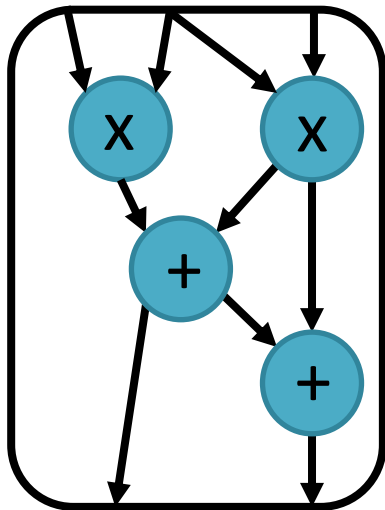Exploit common accelerator application behavior:

# Our Work: Stream-Dataflow Acceleration

Exploit common accelerator application behavior:

## Dataflow Computation

- Stream-Dataflow **Execution model**
  - Abstracts typical accelerator computation phases



*Dataflow Graph*

ISCA 2017 Stream-Dataflow Acceleration Talk

# Our Work:
# Stream-Dataflow Acceleration

**From Memory**

Local storage

**Reuse Stream**

**Memory Stream**

**Recurrence Stream**

X  X
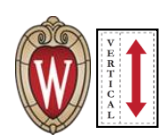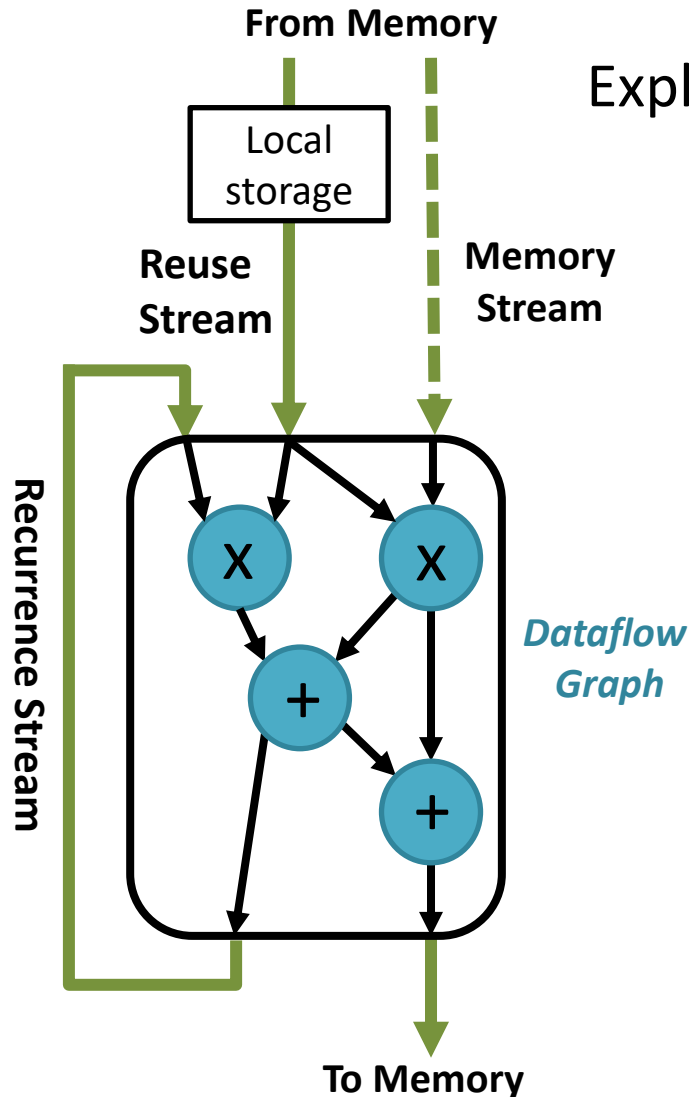
+

+

*Dataflow Graph*

**To Memory**

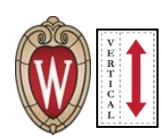Exploit common accelerator application behavior:

## Dataflow Computation

- Stream-Dataflow **Execution model**
  – Abstracts typical accelerator computation phases
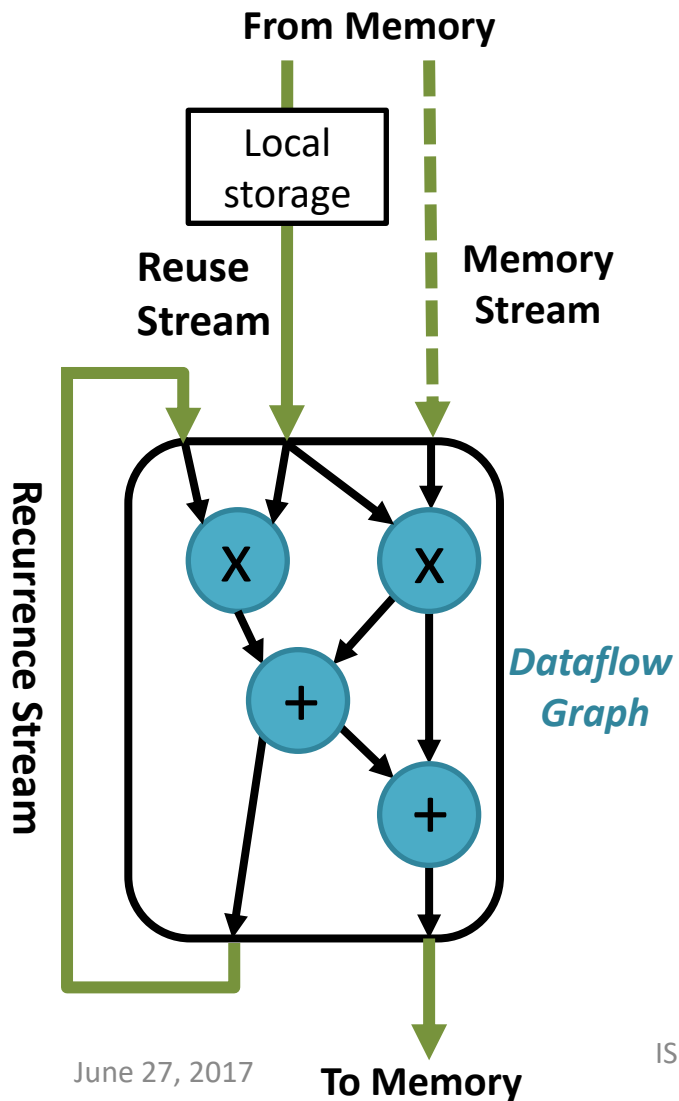
## Stream Patterns and Interface

- Stream-Dataflow **ISA encoding** and **Hardware-Software interface**
  – Exposes parallelism available in these phases

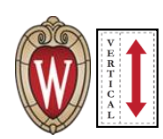UCLA

# Stream-Dataflow Acceleration



## Stream-Dataflow Model

**From Memory**

Local storage

**Reuse Stream**

**Memory Stream**

**Recurrence Stream**

X    X

+

+

*Dataflow Graph*
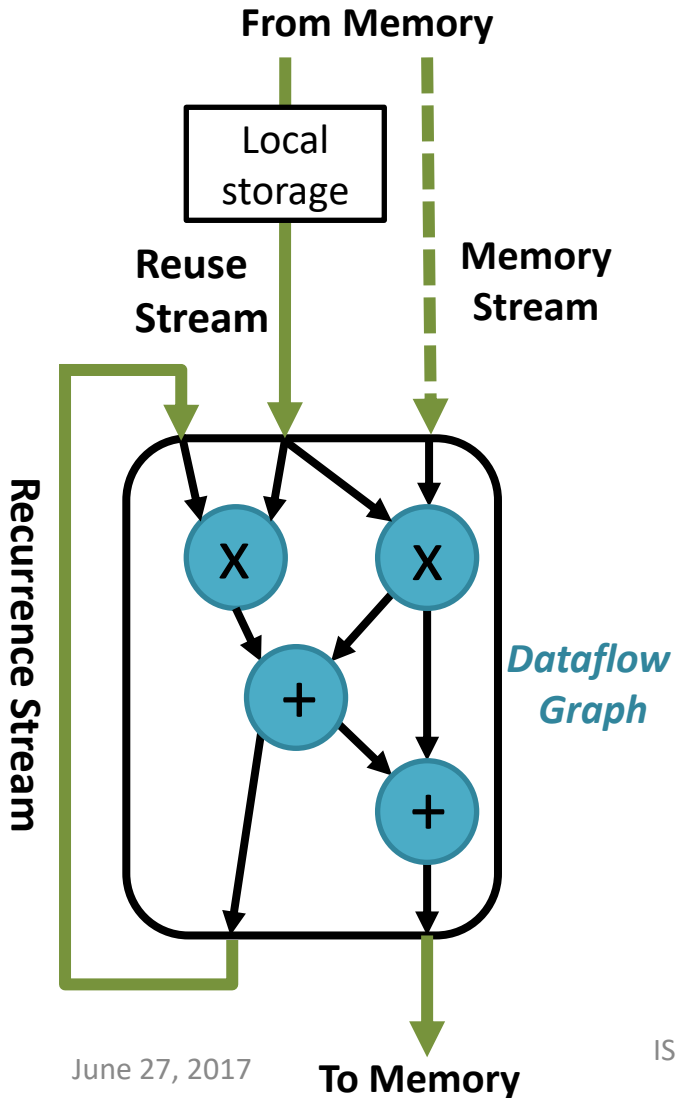
**To Memory**

## Programmable Stream-Dataflow Accelerator

# Stream-Dataflow Acceleration

## Stream-Dataflow Model

**From Memory**

Local storage

**Reuse Stream**

**Memory Stream**

**Recurrence Stream**

X    X

+

+

*Dataflow Graph*

**To Memory**

## Programmable Stream-Dataflow Accelerator

Memory/Cache Hierarchy

Memory Interface

Programmable Scratchpad

Reconfigurable Fabric

# Stream-Dataflow Acceleration

## Stream-Dataflow Model

**From Memory**



Local storage

**Reuse Stream**

**Memory Stream**

**Recurrence Stream**

X  X

+

+

*Dataflow Graph*

**To Memory**

## Programmable Stream-Dataflow Accelerator



Memory/Cache Hierarchy

Memory Interface

Input Data Streams

... Input Data Streams

Programmable Scratchpad

Reconfigurable Fabric

- Data-parallel program kernels streaming data from memory

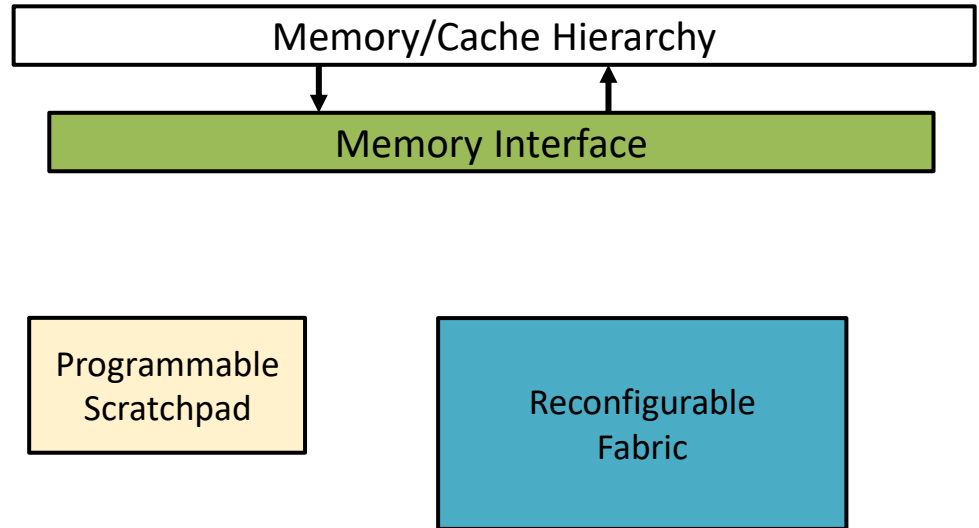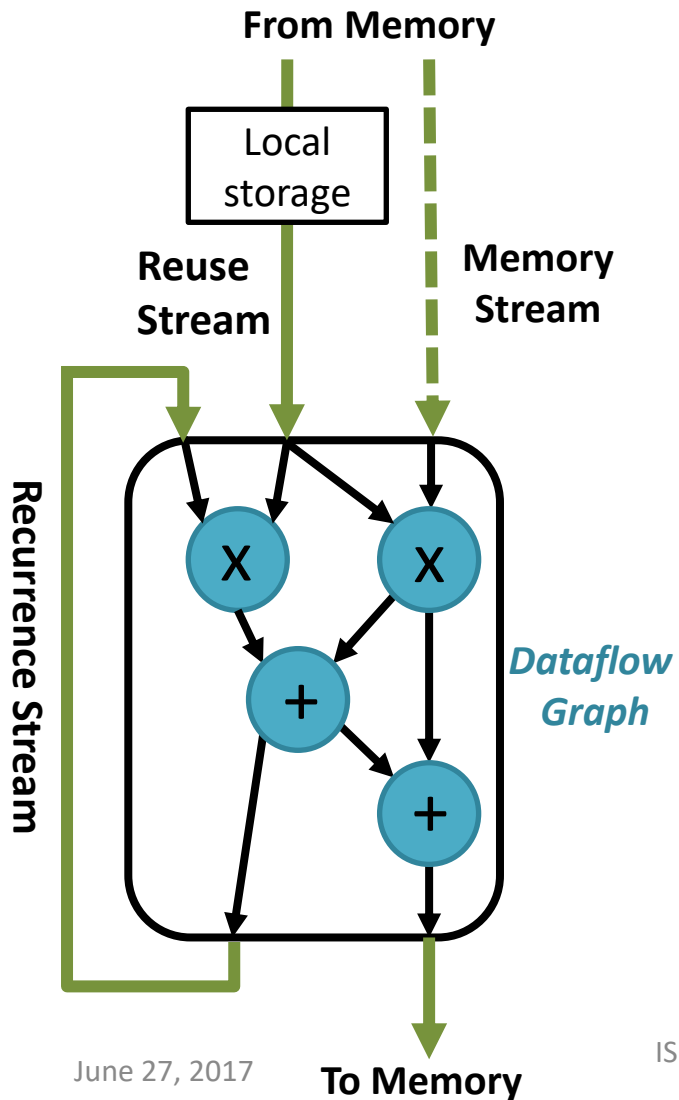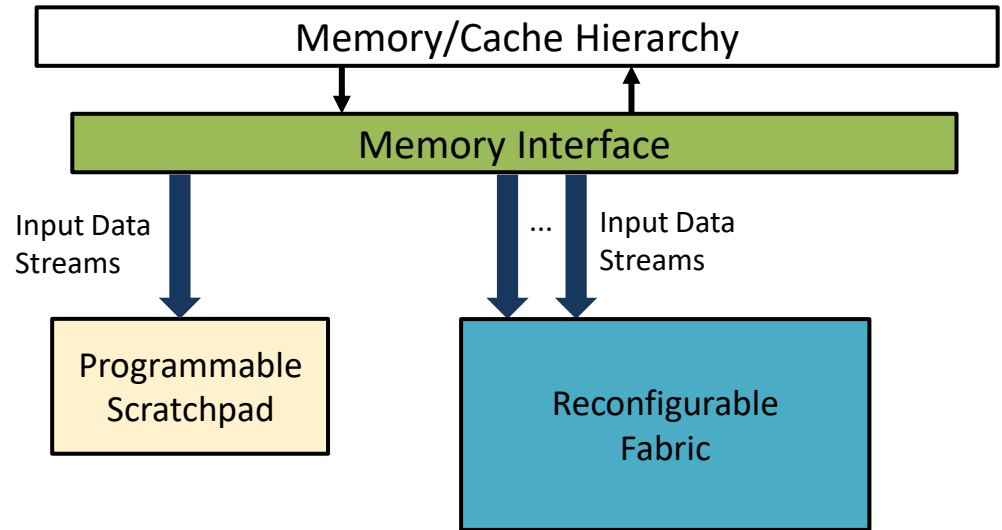ISCA 2017 Stream-Dataflow Acceleration Talk
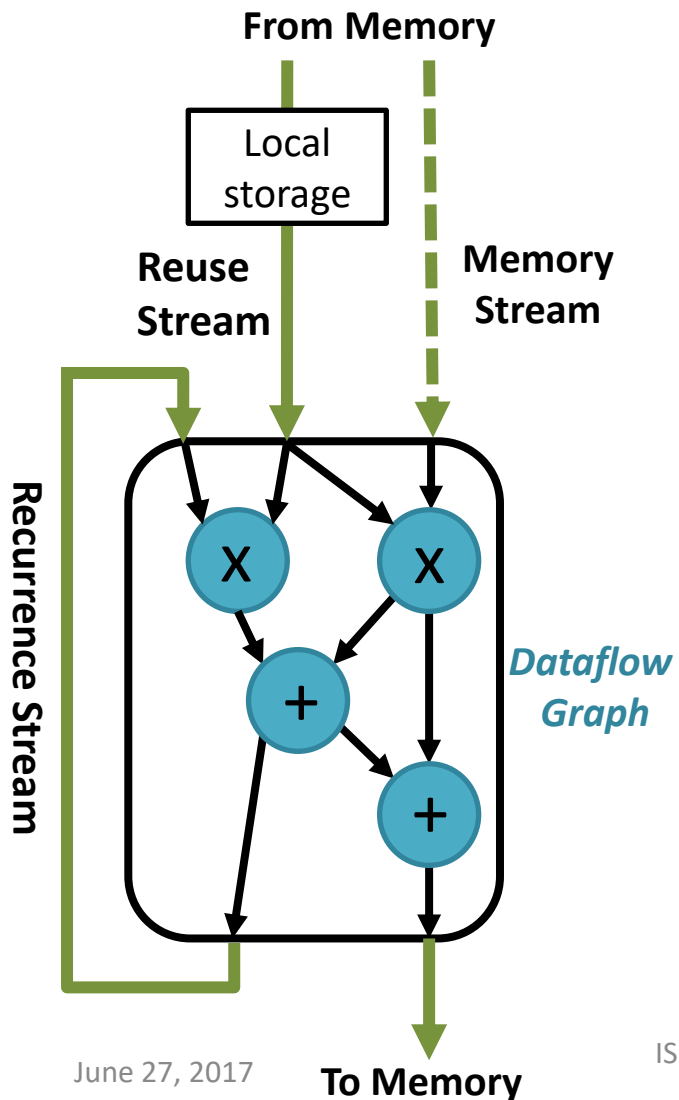
# Stream-Dataflow Acceleration
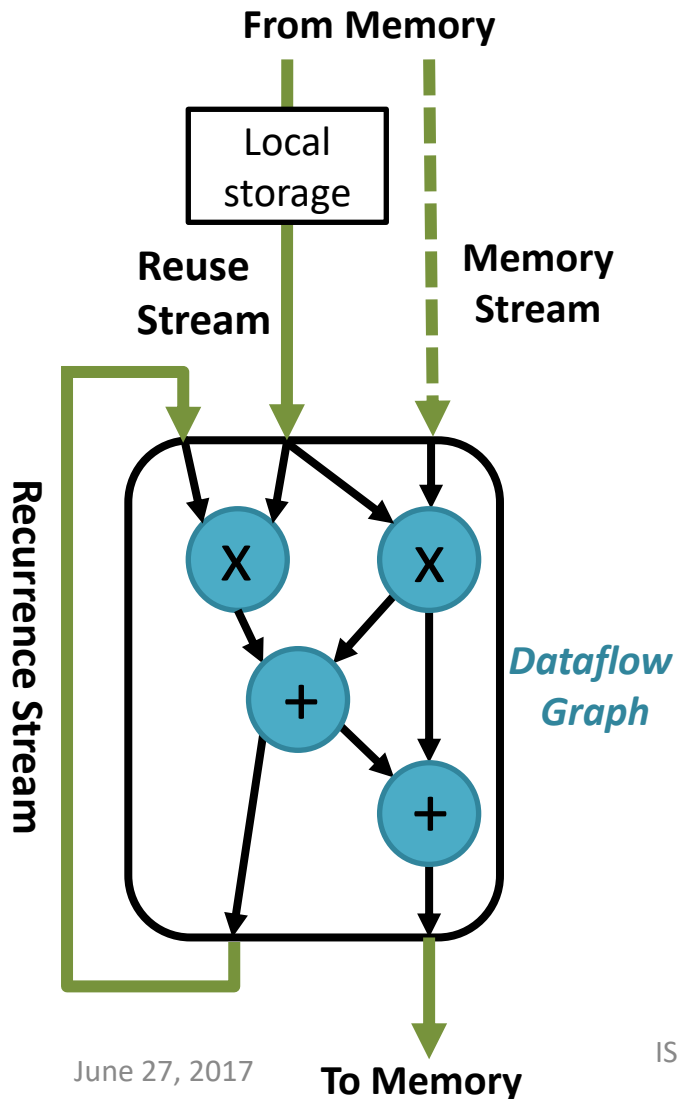
## Stream-Dataflow Model



## Programmable Stream-Dataflow Accelerator



- Data-parallel program kernels streaming data from memory
- Dataflow computation fabric operates on data streams iteratively

ISCA 2017 Stream-Dataflow Acceleration Talk

# Stream-Dataflow Acceleration

## Stream-Dataflow Model

**From Memory**

Local storage

**Reuse Stream**  **Memory Stream**

**Recurrence Stream**

*Dataflow Graph*

X  X

+

+

**To Memory**

## Programmable Stream-Dataflow Accelerator

Memory/Cache Hierarchy

Memory Interface

Input Data Streams    Output Data Streams    ...    Input Data Streams    ...    Output Data Streams

Programmable Scratchpad

Reconfigurable Fabric

Reuse streams

Recurring Data Streams

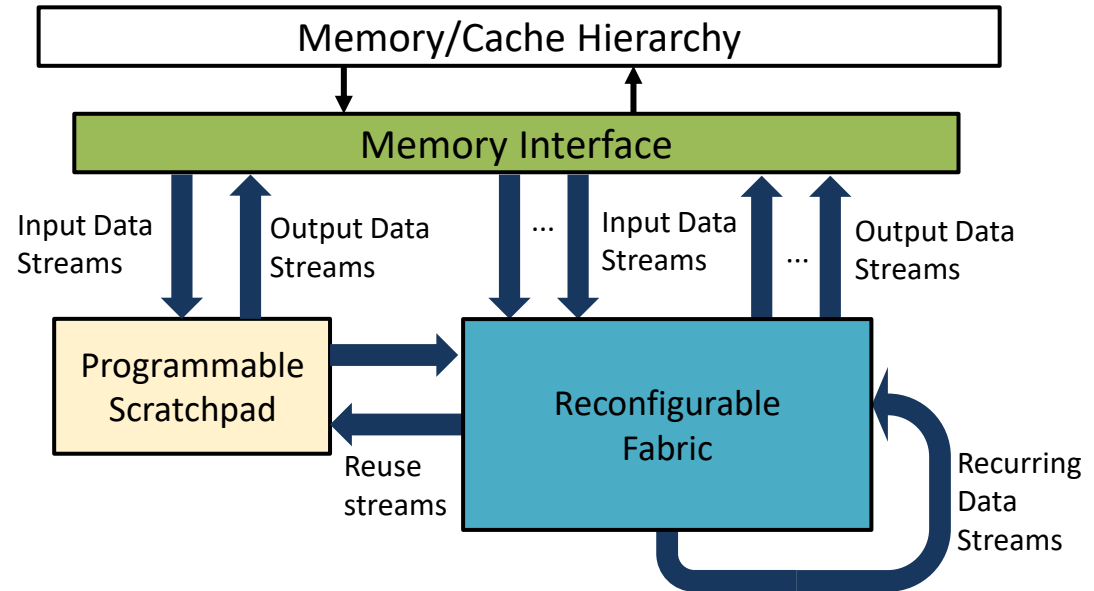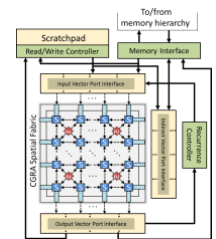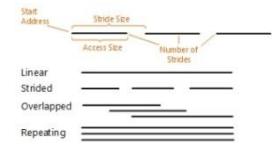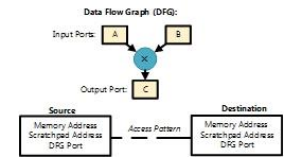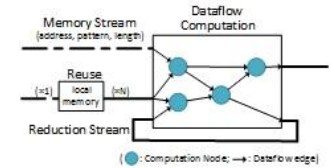- Data-parallel program kernels streaming data from memory
- Dataflow computation fabric operates on data streams iteratively
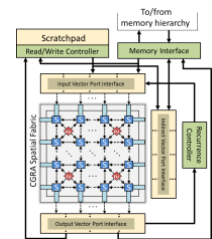- Computed output streams stored back to memory
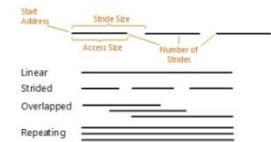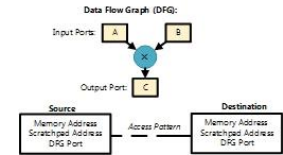
ISCA 2017 Stream-Dataflow Acceleration Talk

# Outline

- Motivation and Overview

- Stream-Dataflow Execution Model

- Hardware-Software Interface and Example program

- Stream-Dataflow Accelerator Architecture

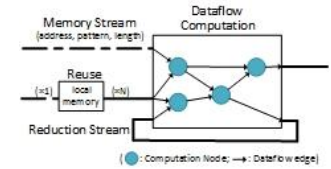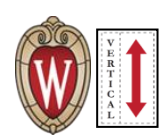- Evaluation and Results

# Outline

- Motivation and Overview

- **Stream-Dataflow Execution Model**

- Hardware-Software Interface and Example program

- Stream-Dataflow Accelerator Architecture

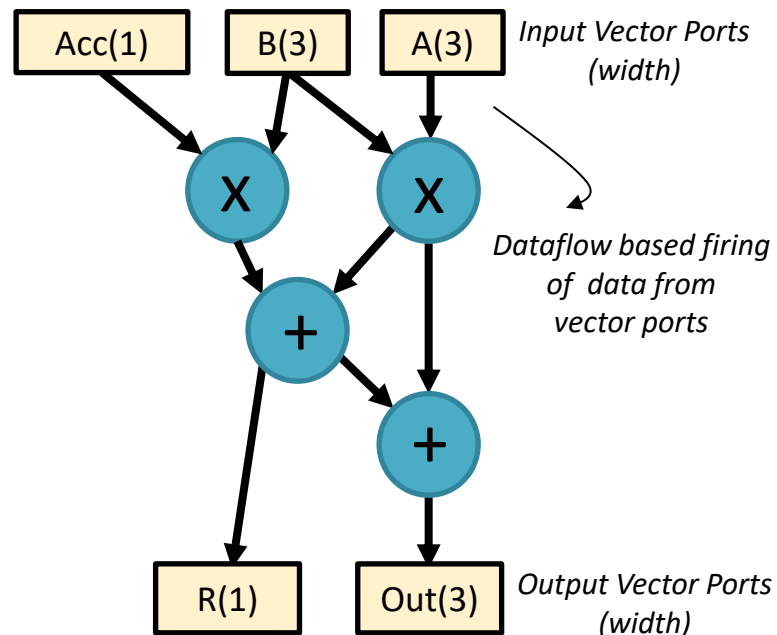- Evaluation and Results

# Stream-Dataflow Execution Model
## *Programmer Abstractions for Stream-Dataflow Model*

# Stream-Dataflow Execution Model

*Programmer Abstractions for Stream-Dataflow Model*

- ***Computation abstraction*** – Dataflow Graph (DFG) with input/output vector ports



Input Vector Ports (width)

Dataflow based firing of data from vector ports

Output Vector Ports (width)

# Stream-Dataflow Execution Model
*Programmer Abstractions for Stream-Dataflow Model*

- ***Computation abstraction*** – Dataflow Graph (DFG) with input/output vector ports



*Dataflow Graph*

# Stream-Dataflow Execution Model

*Programmer Abstractions for Stream-Dataflow Model*

- *Computation abstraction* – Dataflow Graph (DFG) with input/output vector ports

- *Data abstraction* – Streams of data fetched from memory and stored back to memory

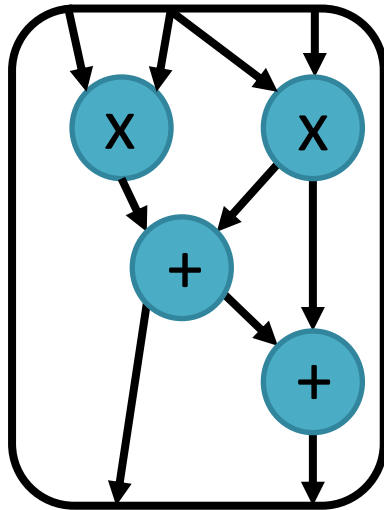**From Memory**

**Memory Stream**

*Dataflow Graph*

# Stream-Dataflow Execution Model

*Programmer Abstractions for Stream-Dataflow Model*

- *Computation abstraction* – Dataflow Graph (DFG) with input/output vector ports

- *Data abstraction* – Streams of data fetched from memory and stored back to memory

**From Memory**

**Memory Stream**

X    X

+

+

*Dataflow Graph*

**To Memory**

# Stream-Dataflow Execution Model

*Programmer Abstractions for Stream-Dataflow Model*

- *Computation abstraction* – Dataflow Graph (DFG) with input/output vector ports
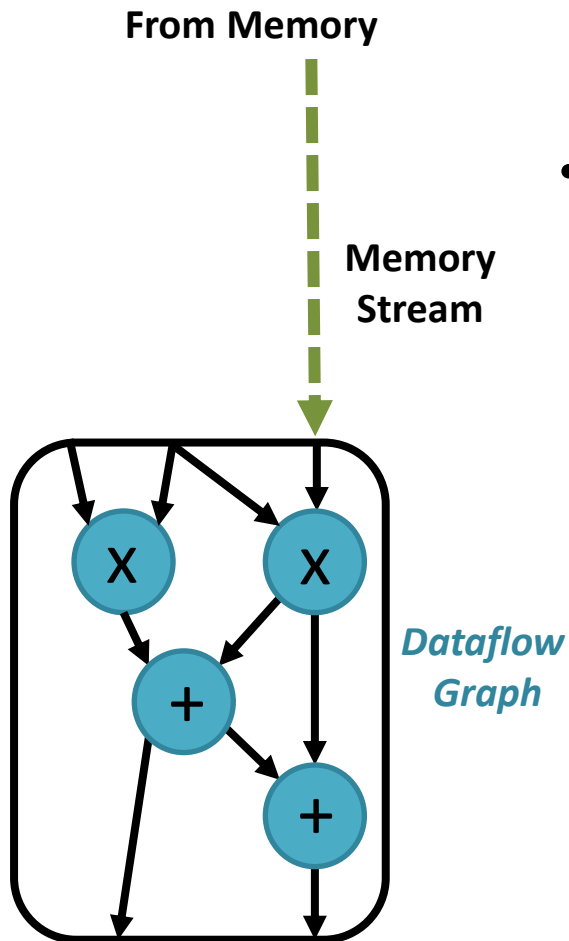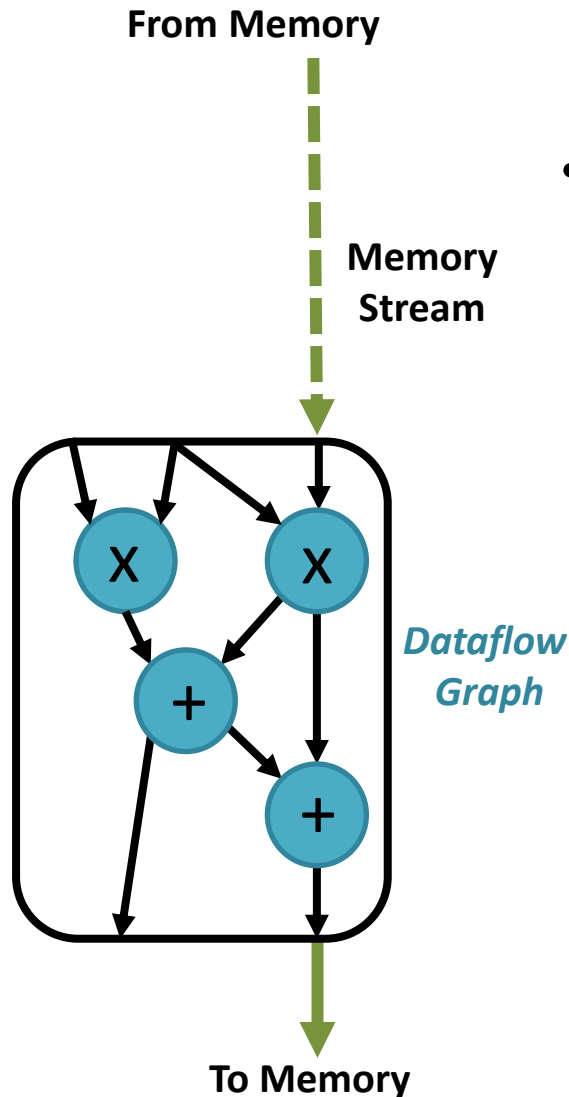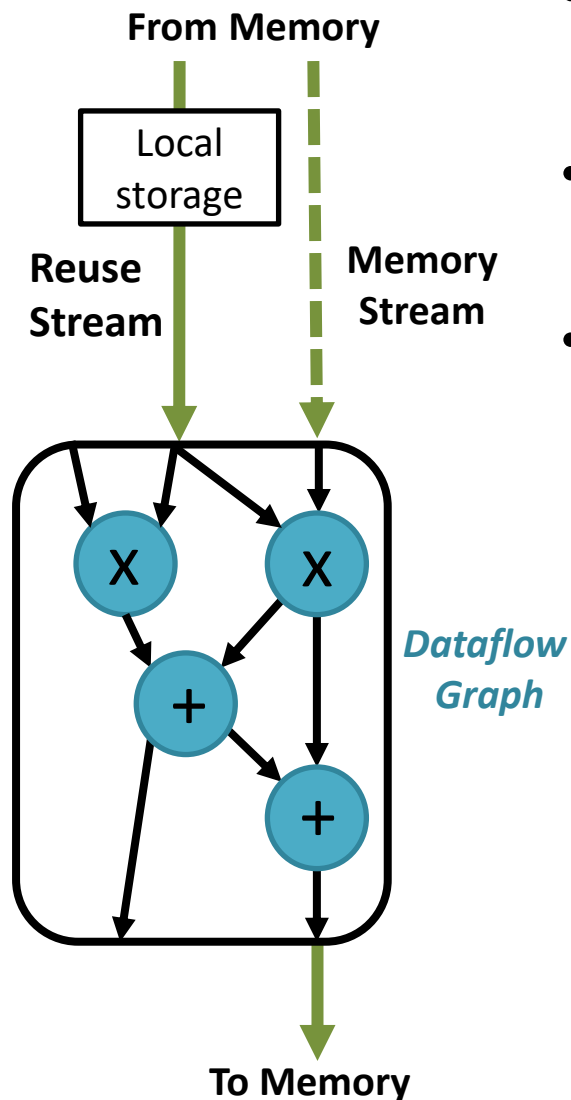
- *Data abstraction* – Streams of data fetched from memory and stored back to memory

- *Reuse abstraction* – Streams of data fetched once from memory, stored in local storage (programmable scratchpad) and reused again

**From Memory**

Local storage

**Reuse Stream**    **Memory Stream**

X    X

+

+

*Dataflow Graph*

**To Memory**

# Stream-Dataflow Execution Model
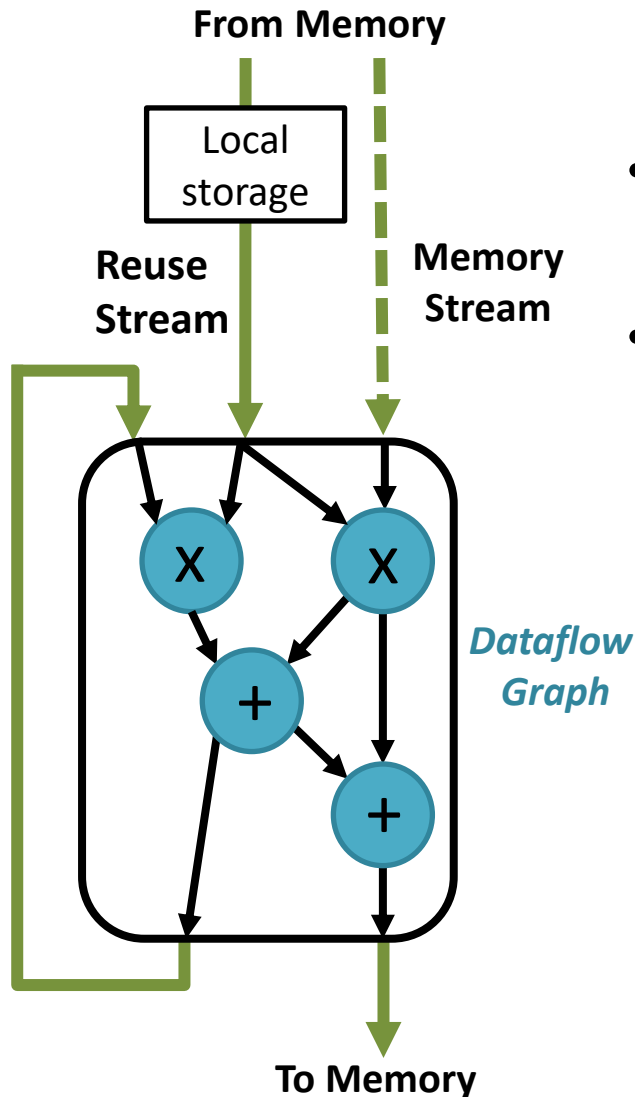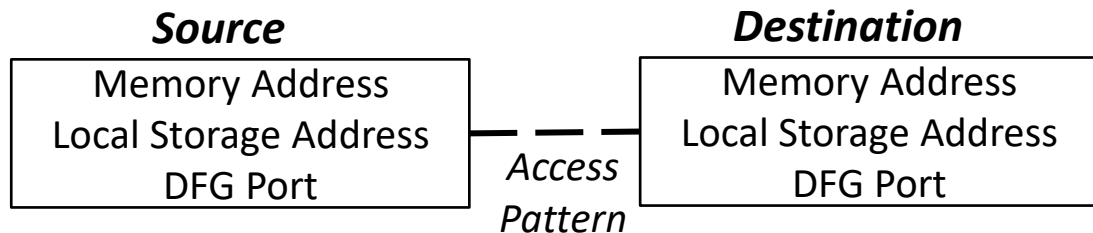
*Programmer Abstractions for Stream-Dataflow Model*

- *Computation abstraction* – Dataflow Graph (DFG) with input/output vector ports

- *Data abstraction* – Streams of data fetched from memory and stored back to memory

- *Reuse abstraction* – Streams of data fetched once from memory, stored in local storage (programmable scratchpad) and reused again



**From Memory**

Local storage

**Reuse Stream**

**Memory Stream**

**Recurrence Stream**

*Dataflow Graph*

**To Memory**

# Stream-Dataflow Execution Model

*Programmer Abstractions for Stream-Dataflow Model*



**From Memory**

Local storage

**Reuse Stream**

**Memory Stream**
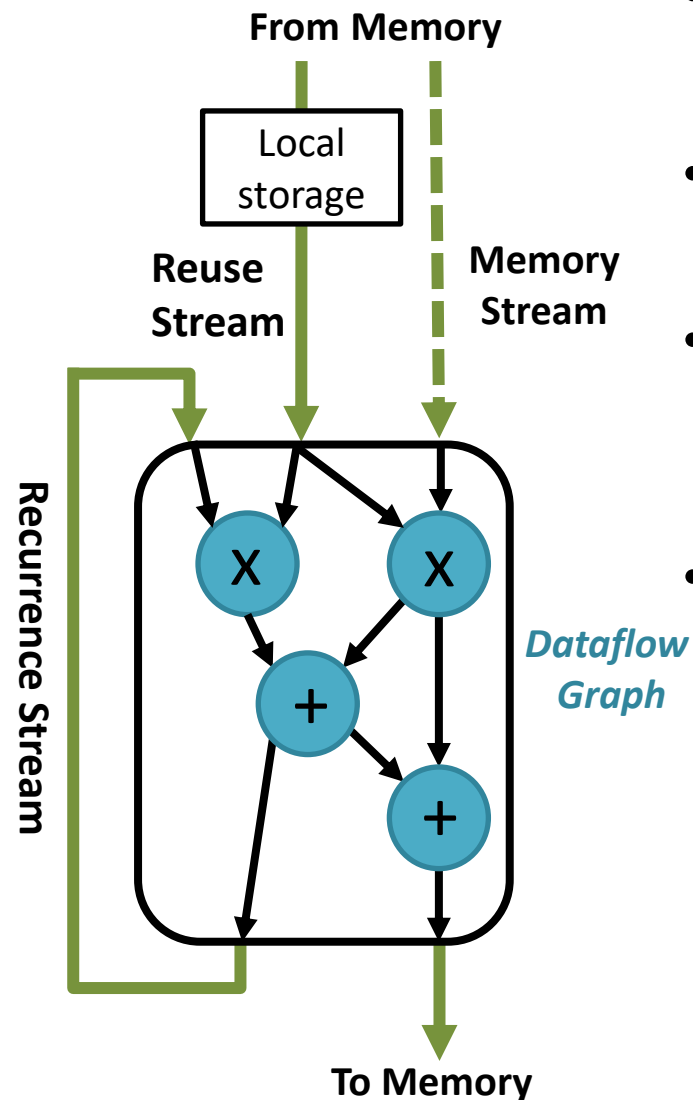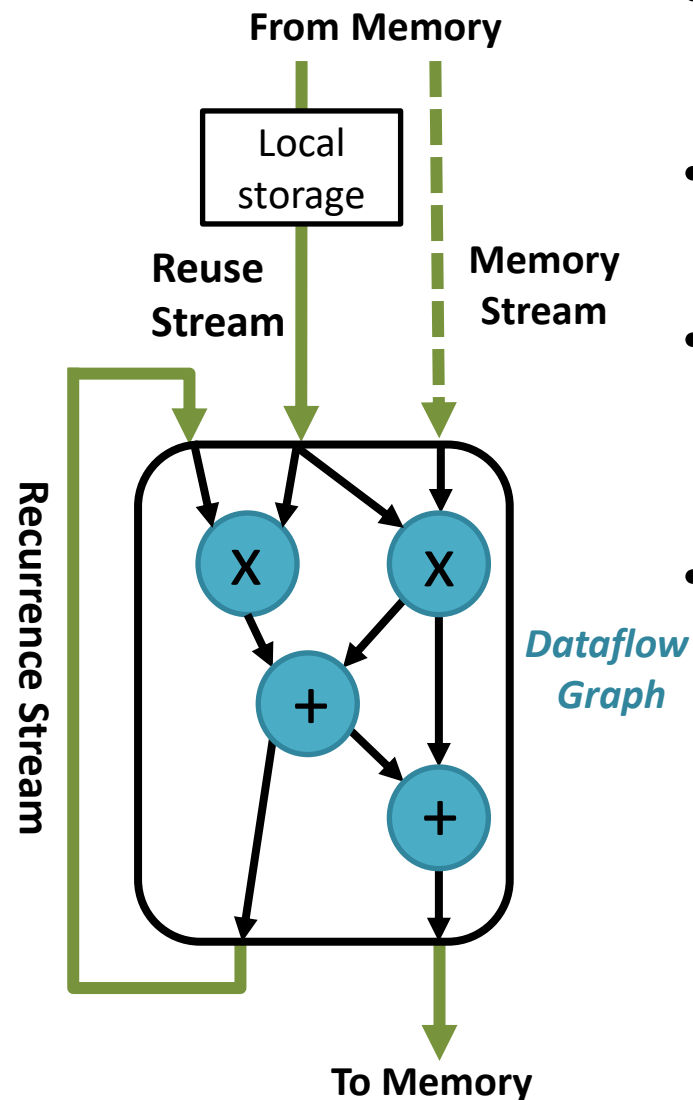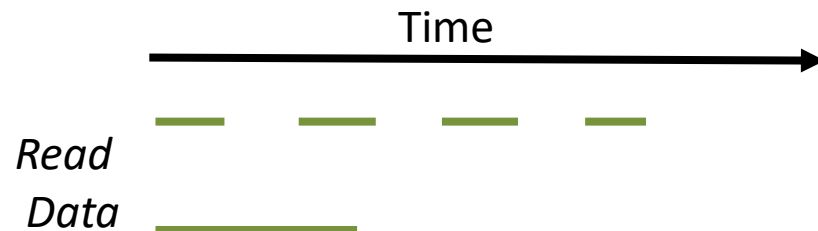
*Dataflow Graph*

**Recurrence Stream**

**To Memory**

- *Computation abstraction* – Dataflow Graph (DFG) with input/output vector ports

- *Data abstraction* – Streams of data fetched from memory and stored back to memory

- *Reuse abstraction* – Streams of data fetched once from memory, stored in local storage (programmable scratchpad) and reused again

- *Communication abstraction* – Stream-Dataflow data movement commands and barriers

| *Source* | *Destination* |
|---|---|
| Memory Address<br>Local Storage Address<br>DFG Port | Memory Address<br>Local Storage Address<br>DFG Port |

*Access Pattern*

# Stream-Dataflow Execution Model

*Programmer Abstractions for Stream-Dataflow Model*

**From Memory**

Local storage

**Reuse Stream**   **Memory Stream**

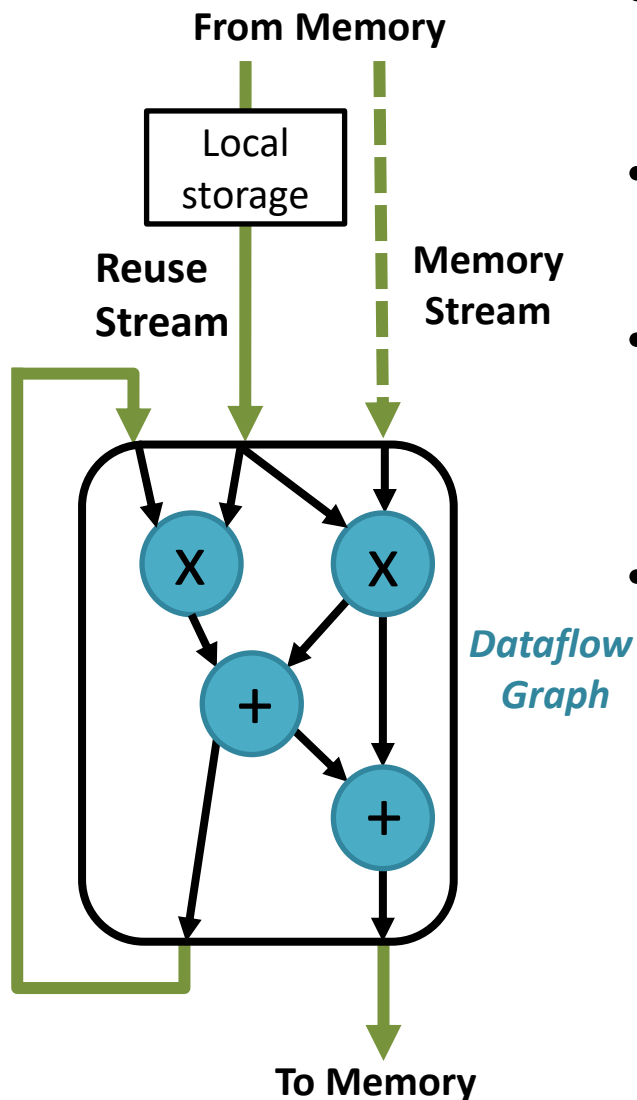**Recurrence Stream**

X   X

+

+

*Dataflow Graph*

**To Memory**
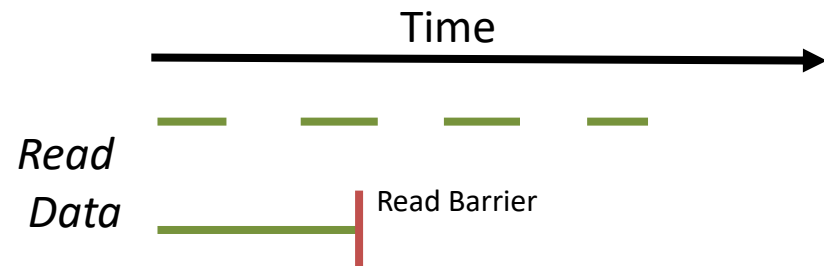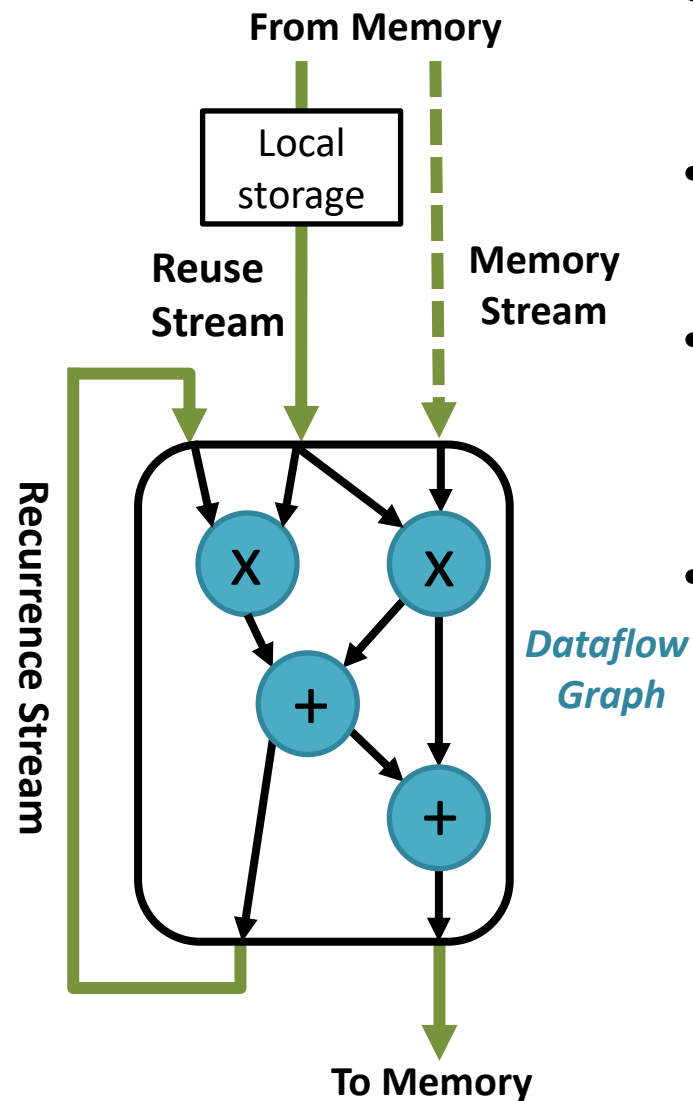
- *Computation abstraction* – Dataflow Graph (DFG) with input/output vector ports

- *Data abstraction* – Streams of data fetched from memory and stored back to memory

- *Reuse abstraction* – Streams of data fetched once from memory, stored in local storage (programmable scratchpad) and reused again

- *Communication abstraction* – Stream-Dataflow data movement commands and barriers

Time ⟶

# Stream-Dataflow Execution Model

*Programmer Abstractions for Stream-Dataflow Model*



*Dataflow Graph*

**From Memory**

Local storage

**Reuse Stream**
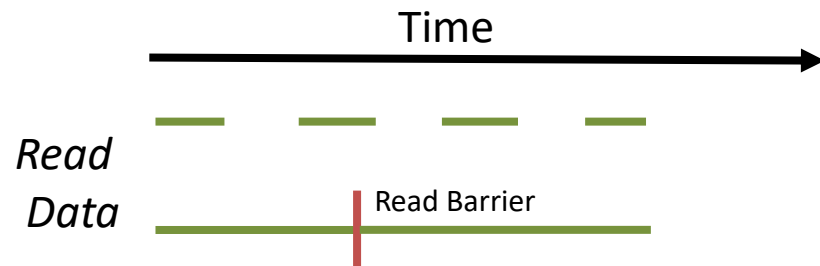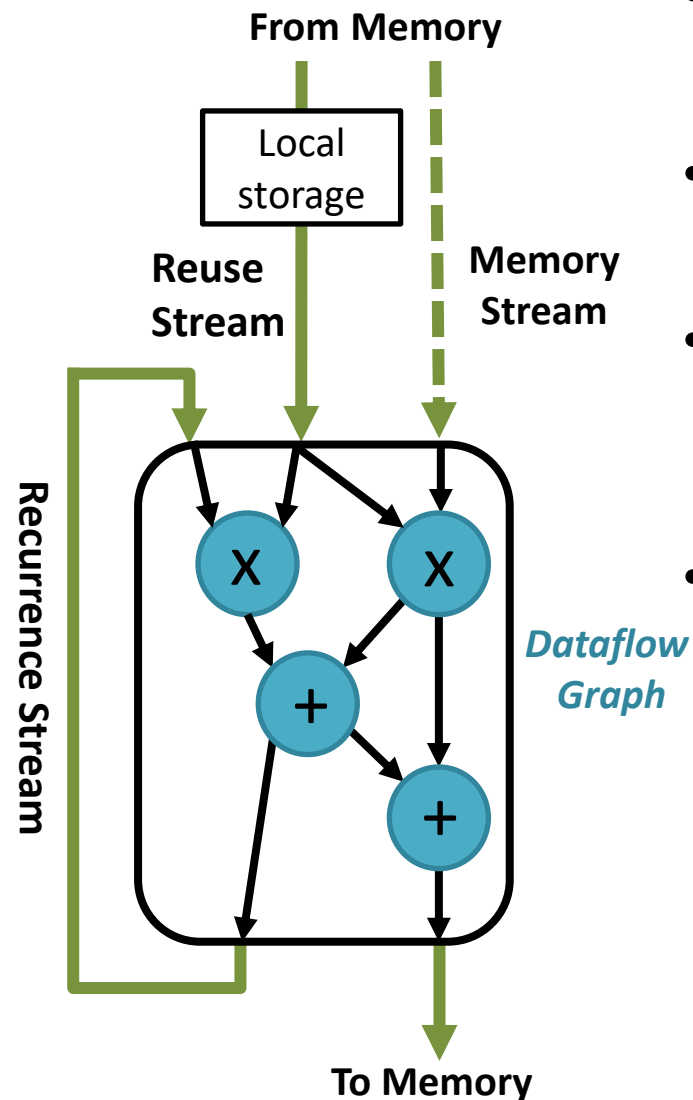
**Memory Stream**

**Recurrence Stream**

**To Memory**

- *Computation abstraction* – Dataflow Graph (DFG) with input/output vector ports

- *Data abstraction* – Streams of data fetched from memory and stored back to memory

- *Reuse abstraction* – Streams of data fetched once from memory, stored in local storage (programmable scratchpad) and reused again

- *Communication abstraction* – Stream-Dataflow data movement commands and barriers

Time

*Read Data*

# Stream-Dataflow Execution Model

*Programmer Abstractions for Stream-Dataflow Model*



**From Memory**

Local storage

**Reuse Stream**   **Memory Stream**

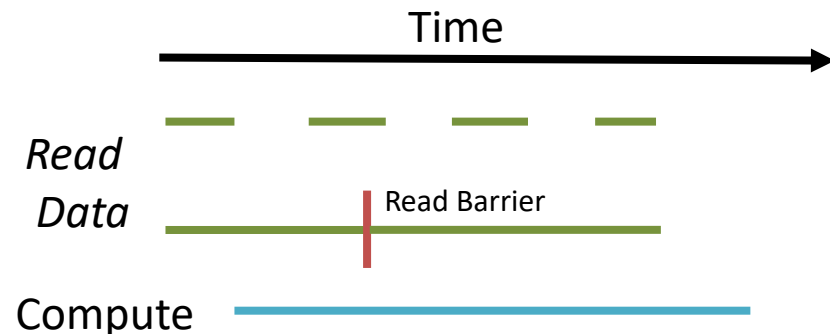**Recurrence Stream**

X  X

+

+

*Dataflow Graph*

**To Memory**
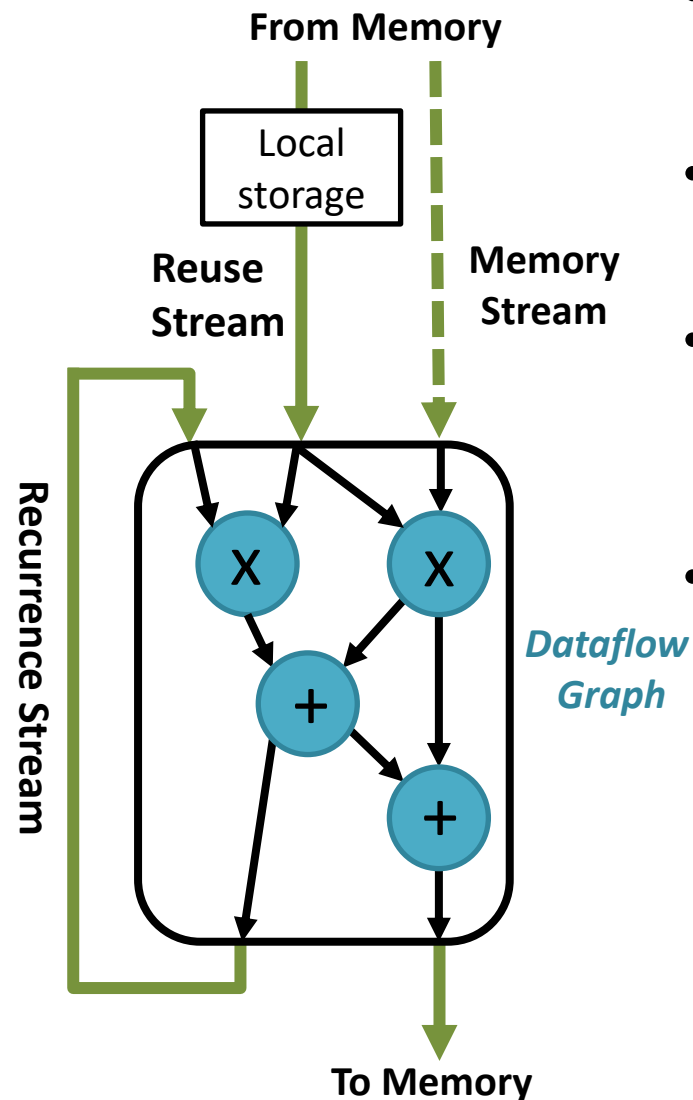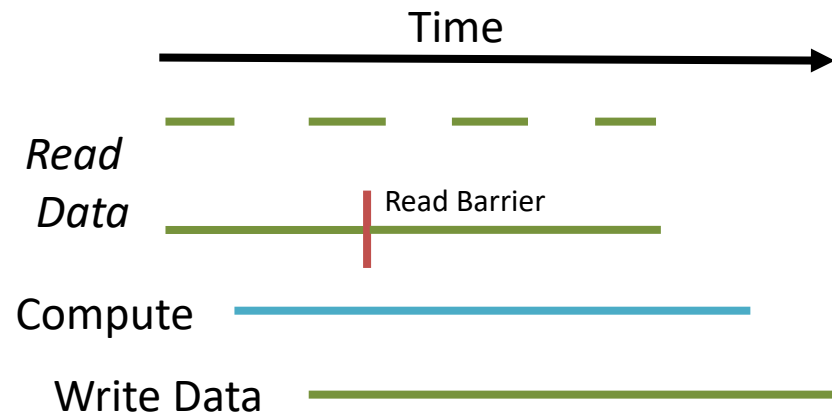
- *Computation abstraction* – Dataflow Graph (DFG) with input/output vector ports

- *Data abstraction* – Streams of data fetched from memory and stored back to memory

- *Reuse abstraction* – Streams of data fetched once from memory, stored in local storage (programmable scratchpad) and reused again

- *Communication abstraction* – Stream-Dataflow data movement commands and barriers

Time

*Read Data*

Read Barrier

# Stream-Dataflow Execution Model
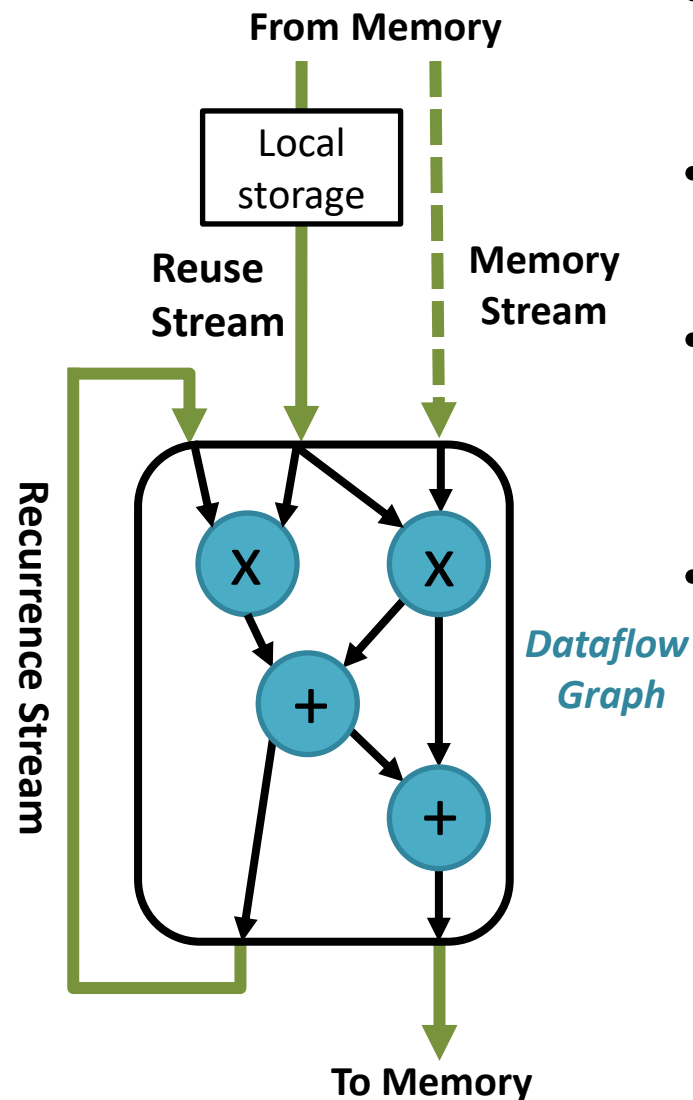
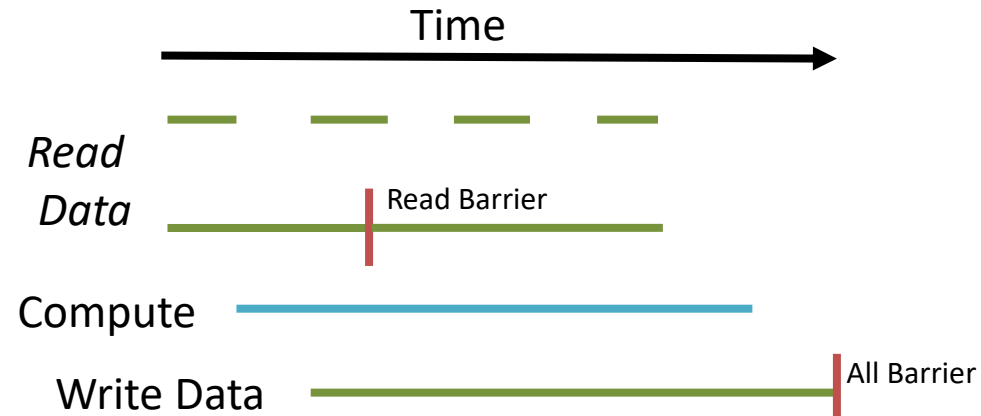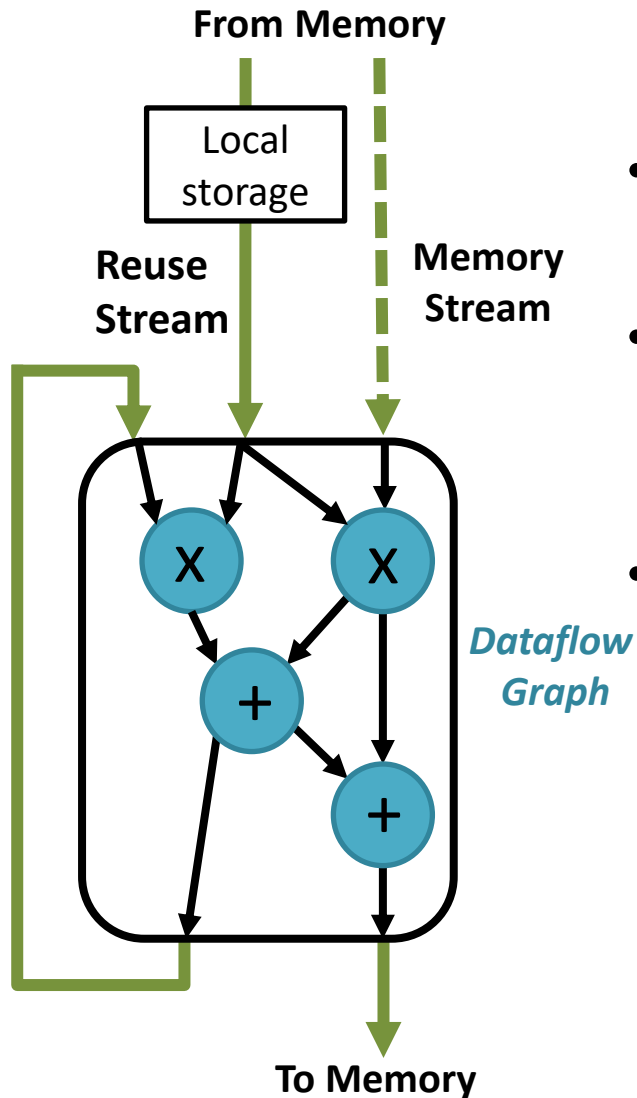*Programmer Abstractions for Stream-Dataflow Model*



- **Computation abstraction** – Dataflow Graph (DFG) with input/output vector ports

- **Data abstraction** – Streams of data fetched from memory and stored back to memory

- **Reuse abstraction** – Streams of data fetched once from memory, stored in local storage (programmable scratchpad) and reused again

- **Communication abstraction** – Stream-Dataflow data movement commands and barriers
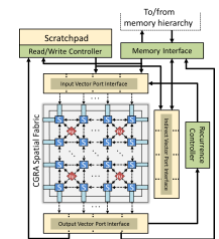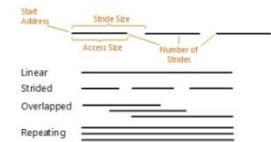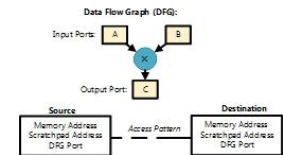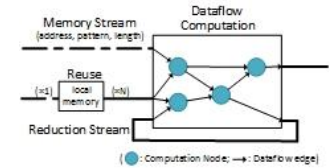
# Stream-Dataflow Execution Model

*Programmer Abstractions for Stream-Dataflow Model*



**From Memory**

Local storage

**Reuse Stream**

**Memory Stream**

**Recurrence Stream**

X   X

+

+

*Dataflow Graph*

**To Memory**

- *Computation abstraction* – Dataflow Graph (DFG) with input/output vector ports

- *Data abstraction* – Streams of data fetched from memory and stored back to memory

- *Reuse abstraction* – Streams of data fetched once from memory, stored in local storage (programmable scratchpad) and reused again

- *Communication abstraction* – Stream-Dataflow data movement commands and barriers

Time

*Read Data*

Read Barrier

Compute

UCLA

# Stream-Dataflow Execution Model

*Programmer Abstractions for Stream-Dataflow Model*



**From Memory**

Local storage

**Reuse Stream**     **Memory Stream**

**Recurrence Stream**

X    X

+

+

*Dataflow Graph*

**To Memory**

- *Computation abstraction* – Dataflow Graph (DFG) with input/output vector ports

- *Data abstraction* – Streams of data fetched from memory and stored back to memory

- *Reuse abstraction* – Streams of data fetched once from memory, stored in local storage (programmable scratchpad) and reused again

- *Communication abstraction* – Stream-Dataflow data movement commands and barriers

Time

*Read Data*     Read Barrier

Compute

Write Data

# Stream-Dataflow Execution Model

*Programmer Abstractions for Stream-Dataflow Model*

- *Computation abstraction* – Dataflow Graph (DFG) with input/output vector ports

- *Data abstraction* – Streams of data fetched from memory and stored back to memory

- *Reuse abstraction* – Streams of data fetched once from memory, stored in local storage (programmable scratchpad) and reused again

- *Communication abstraction* – Stream-Dataflow data movement commands and barriers



**From Memory**

Local storage

**Reuse Stream**    **Memory Stream**

**Recurrence Stream**

X    X

+

+

*Dataflow Graph*

**To Memory**

Time

*Read Data*    Read Barrier

Compute

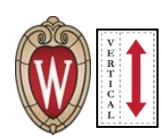Write Data    All Barrier

# Outline

- Motivation and Overview

- Stream-Dataflow Execution Model

- Hardware-Software Interface and Example program
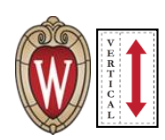
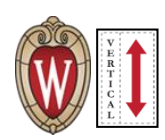- Stream-Dataflow Accelerator Architecture

- Evaluation and Results

# Outline

- Motivation and Overview

- Stream-Dataflow Execution Model

- **Hardware-Software Interface and Example program**

- Stream-Dataflow Accelerator Architecture

- Evaluation and Results

**Stream-Dataflow ISA Interface**

Express any data-stream pattern of accelerator applications using simple, flexible and yet efficient encoding

# Stream-Dataflow ISA

# Stream-Dataflow ISA

- Set-up Interface:
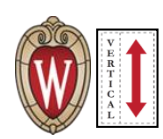    - **SD_Config** – Configuration data stream for dataflow computation fabric (CGRA)
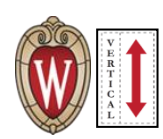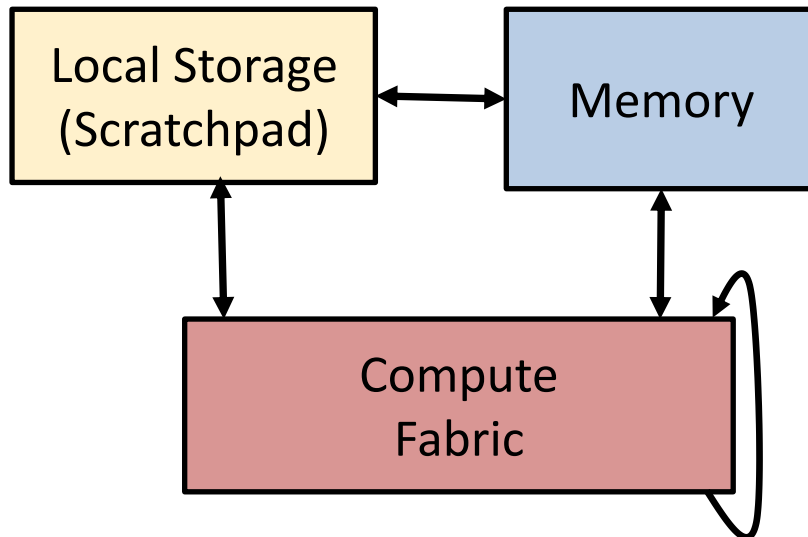
# Stream-Dataflow ISA

- Set-up Interface:
    - **SD_Config** – Configuration data stream for dataflow computation fabric (CGRA)

- Control Interface:
    **SD_Barrier_Scratch_Rd**, **SD_Barrier_Scratch_Wr, SD_Barrier_All**

# Stream-Dataflow ISA

- Set-up Interface:
  - ❑ **SD_Config** – Configuration data stream for dataflow computation fabric (CGRA)

- Control Interface:
  **SD_Barrier_Scratch_Rd, SD_Barrier_Scratch_Wr, SD_Barrier_All**

- Stream Interface → **SD_[source]_[dest]**
  Source/Dest Parameters: *Address (memory or local_storage), DFG Port number*
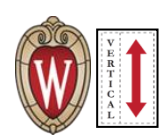  Pattern Parameters: *access_size, stride_size, num_strides*

# Stream-Dataflow ISA
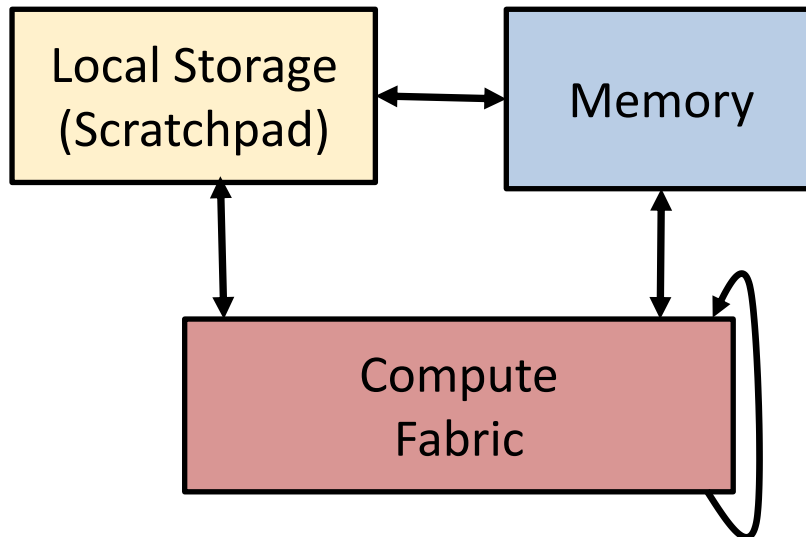
- Set-up Interface:
  - ❑ **SD_Config** – Configuration data stream for dataflow computation fabric (CGRA)

- Control Interface:
  **SD_Barrier_Scratch_Rd**, **SD_Barrier_Scratch_Wr**, **SD_Barrier_All**

- Stream Interface → **SD_[source]_[dest]**
  Source/Dest Parameters: *Address (memory or local_storage), DFG Port number*
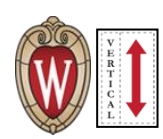  Pattern Parameters: *access_size, stride_size, num_strides*

| Command Name | Parameters | Description |
|---|---|---|
| SD_Config | Address, Size | Stream CGRA configuration from given address |
| SD_Mem_Scratch | Source Mem Address, Stride, Access Size, Num Strides, Dest. Scratch Address | Read from memory with pattern to scratchpad |
| SD_Scratch_Port | Source Scratch Address, Stride, Access Size, Strides, Input Port # | Read from scratchpad with pattern to input port |
| SD_Mem_Port | Source Mem Address, Stride, Access Size, Num Strides, Input Port # | Read from memory with pattern to input port |
| SD_Const_Port | Constant Value, Num Elements, Input Port # | Send constant value to input port |
| SD_Clean_Port | Num Elements, Output Port # | Throw away some elements from output port |
| SD_Port_Port | Output Port #, Num Elements, Input Port # | Issue recurrence between input-output port pairs |
| SD_Port_Scratch | Output Port #, Num Elements, Scratch Address | Write from port to scratchpad |
| SD_Port_Mem | Output Port #, Stride, Access Size, Num Strides, Dest. Mem Address | Write from port to memory with pattern |
| SD_Mem_IndPort | Source Mem Address, Stride, Access Size, Num Strides, Indirect Port # | Read the addresses from memory with pattern to indirect port |
| SD_IndPort_Port | Indirect Port #, Offset Address, Input Port # | Indirect load from addresses present in indirect port |
| SD_IndPort_Mem | Indirect Port #, Output Port #, Dest. Offset Address | Indirect store to addresses present in indirect port |
| SD_Barrier_Scratch_Rd | - | Barrier for scratchpad reads |
| SD_Barrier_Scratch_Wr | - | Barrier for scratchpad writes |
| SD_Barrier_All | - | Barrier to wait for all commands completion |

# Stream-Dataflow ISA

- **Set-up Interface:**
  - ❑ **SD_Config** – Configuration data stream for dataflow computation fabric (CGRA)

- **Control Interface:**
  **SD_Barrier_Scratch_Rd**, **SD_Barrier_Scratch_Wr**, **SD_Barrier_All**

- **Stream Interface → SD_[source]_[dest]**
  Source/Dest Parameters: *Address (memory or local_storage), DFG Port number*
  Pattern Parameters: *access_size, stride_size, num_strides*
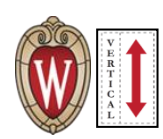
# Stream-Dataflow ISA

- Set-up Interface:
  - ❑ **SD_Config** – Configuration data stream for dataflow computation fabric (CGRA)

- Control Interface:
  **SD_Barrier_Scratch_Rd**, **SD_Barrier_Scratch_Wr**, **SD_Barrier_All**

- Stream Interface → **SD_[source]_[dest]**
  Source/Dest Parameters: *Address (memory or local_storage), DFG Port number*
  Pattern Parameters: *access_size, stride_size, num_strides*

# Stream-Dataflow
# Hardware-Software Interface

**Access Pattern**

**Source**
Memory,
Local Storage,
DFG Port

**Destination**
Memory,
Local Storage,
DFG Port

# Stream-Dataflow
# Hardware-Software Interface

**Source**

Memory,
Local Storage,
DFG Port

**Access Pattern**

Start Address

Stride

Access Size

Number of Strides

**Destination**

Memory,
Local Storage,
DFG Port

# Stream-Dataflow
# Hardware-Software Interface

## Access Pattern

**Source**
Memory,
Local Storage,
DFG Port

**Start Address**

Stride

Access Size

**Number of Strides**

**Destination**
Memory,
Local Storage,
DFG Port

access_size = 4

mem_addr
= 0xA

num_strides
= 2

memory_stride = 8

# Stream-Dataflow Hardware-Software Interface

## Access Pattern

**Source**
Memory,
Local Storage,
DFG Port

**Start Address**

Stride

Access Size

**Number of Strides**

**Destination**
Memory,
Local Storage,
DFG Port

**Example Access Patterns**

Linear

Strided

Overlapped

Repeating

Offset-Indirect

# Stream-Dataflow Hardware-Software Interface

## Access Pattern

**Source**
Memory, Local Storage, DFG Port

**Start Address**

Stride

Access Size

**Number of Strides**

**Destination**
Memory, Local Storage, DFG Port

**Example Access Patterns**

Linear

Strided

Overlapped

Repeating

Offset-Indirect

2D Direct Streams

# Stream-Dataflow Hardware-Software Interface

UCLA

**Access Pattern**

**Source**
Memory, Local Storage, DFG Port

Start Address

Stride

Access Size

Number of Strides

**Destination**
Memory, Local Storage, DFG Port

**Example Access Patterns**

Linear

Strided

Overlapped

Repeating

2D Direct Streams

Offset-Indirect

2D Indirect Streams

# Stream-Dataflow ISA Encoding

Stream:



Dataflow:

# Stream-Dataflow ISA Encoding

## Stream:



## Dataflow:



*Dataflow Graph*

Vector A[0:2]   Vector B[0:2]

C

*Specified in a Domain Specific Language (DSL)*

# Stream-Dataflow ISA Encoding

## Stream:

*Stream Encoding*
**<address, access_size, stride_size, length>**

## Dataflow:



*Specified in a Domain Specific Language (DSL)*

# Stream-Dataflow ISA Encoding

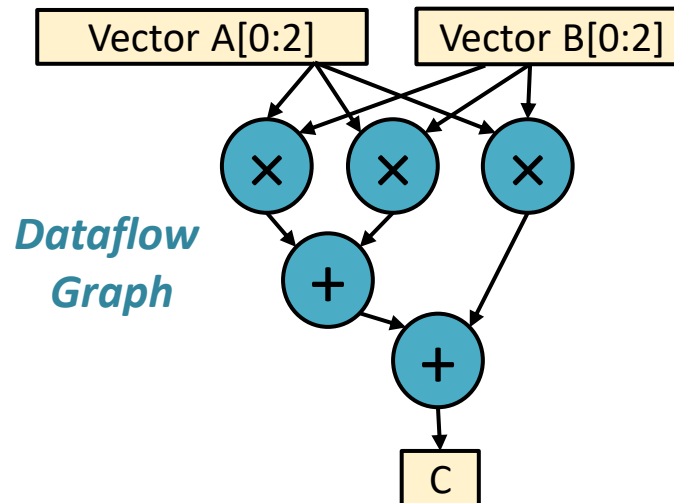## Stream:

Time →

```
for i = 1 to 100:
  ... = a[2*i];
```

*Stream Encoding*
**<address, access_size, stride_size, length>**
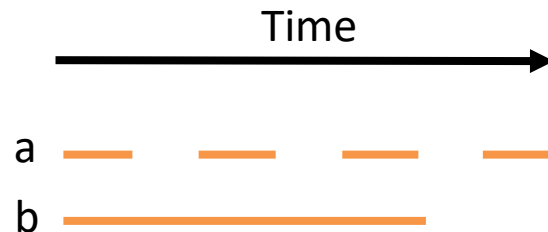
## Dataflow:

| Vector A[0:2] | Vector B[0:2] |



*Dataflow Graph*

*Specified in a Domain Specific Language (DSL)*

# Stream-Dataflow ISA Encoding

## Stream:

```
for i = 1 to 100:
    ... = a[2*i];
```

Time →

a

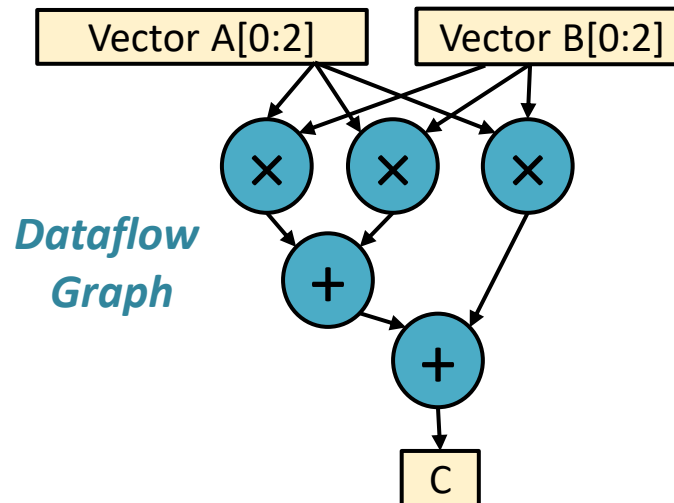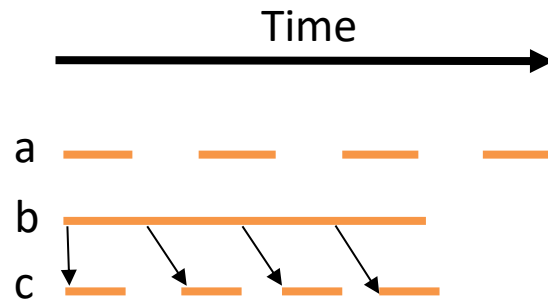*Stream Encoding*
**<address, access_size, stride_size, length>**

*Eg: <a, 1, 2, 100>*

## Dataflow:

Vector A[0:2]   Vector B[0:2]

*Dataflow Graph*

× × ×

+

+

C

*Specified in a Domain Specific Language (DSL)*

# Stream-Dataflow ISA Encoding

## Stream:

```
for i = 1 to 100:
    ... = a[2*i];
    ... = b[i];
```
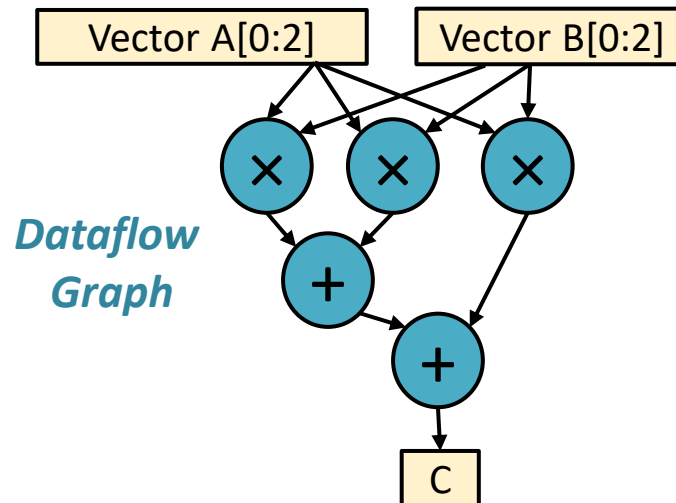
Time →

a

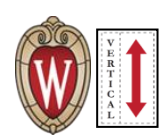*Stream Encoding*

**<address, access_size, stride_size, length>**

*Eg: <a, 1, 2, 100>*

## Dataflow:

Vector A[0:2]    Vector B[0:2]

×  ×  ×

+

+

C

*Dataflow Graph*

*Specified in a Domain Specific Language (DSL)*

# Stream-Dataflow ISA Encoding

## Stream:

```
for i = 1 to 100:
    ... = a[2*i];
    ... = b[i];
```

Time

a

b

*Stream Encoding*

**<address, access_size, stride_size, length>**

*Eg:  <a, 1,  2, 100>*

*<b, 1, 1, 100>*

## Dataflow:

Vector A[0:2]     Vector B[0:2]

*Dataflow Graph*

×  ×  ×

+

+

C

*Specified in a Domain Specific Language (DSL)*

# Stream-Dataflow ISA Encoding

## Stream:

```
for i = 1 to 100:
    ... = a[2*i];
    ... = b[i];
    c[b[i]] = ...
```

Time →

a

b

*Stream Encoding*

**<address, access_size, stride_size, length>**

*Eg: <a, 1, 2, 100>*

*<b, 1, 1, 100>*

## Dataflow:

Vector A[0:2]    Vector B[0:2]

*Dataflow Graph*

×  ×  ×

+

+

C

*Specified in a Domain Specific Language (DSL)*

# Stream-Dataflow ISA Encoding

## Stream:

```
for i = 1 to 100:
    ... = a[2*i];
    ... = b[i];
    c[b[i]] = ...
```

Time →

a
b
c

*Stream Encoding*

**<address, access_size, stride_size, length>**

*Eg: <a, 1, 2, 100>*

*<b, 1, 1, 100>*

**<stream, start, offset_address>**

*IND<[prev], c, 100>*

## Dataflow:



Vector A[0:2]    Vector B[0:2]

×  ×  ×

+

+

C

*Dataflow Graph*

*Specified in a Domain Specific Language (DSL)*

# Example Code: Dot Product

## Original Program

```
for(int i = 0 to N) {
  c += a[i] * b[i];
}
```

# Example Code: Dot Product

## Original Program

```
for(int i = 0 to N) {
  c += a[i] * b[i];
}
```
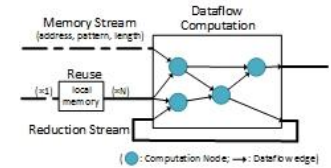
**Dataflow Encoding**

# Example Code: Dot Product

**Stream ISA Encoding**

```
Send a[0: N] → P1
Send b[0: N] → P2
Get P3 → c
```

**Original Program**

```
for(int i = 0 to N) {
  c += a[i] * b[i];
}
```
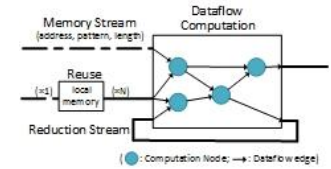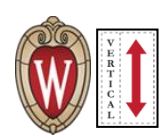
**Dataflow Encoding**

# Outline

- Motivation and Overview

- Stream-Dataflow Execution Model

- Hardware-Software Interface and Example program

- **Stream-Dataflow Accelerator Architecture**

- **Evaluation and Results**

# Outline

- Motivation and Overview

- Stream-Dataflow Execution Model

- Hardware-Software Interface and Example program

- **Stream-Dataflow Accelerator Architecture**
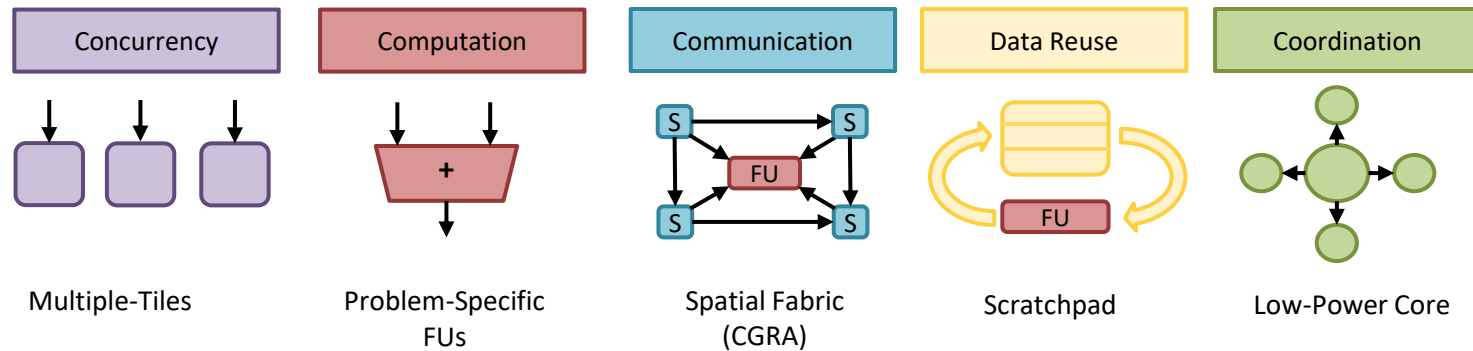
- Evaluation and Results

# Requirements for Stream-Dataflow Accelerator Architecture

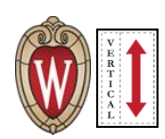1. Should employ the common specialization principles and hardware mechanisms

   (*IEEE Micro Top-Picks 2017: *Domain Specialization is Generally Unnecessary for Accelerators)*



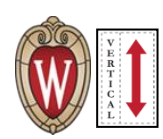| Concurrency | Computation | Communication | Data Reuse | Coordination |
| Multiple-Tiles | Problem-Specific FUs | Spatial Fabric (CGRA) | Scratchpad | Low-Power Core |

2. Programmability features without the inefficiencies of existing data-parallel architectures* (with less power, area and control overheads)

*More detailed analysis contrasting data-parallel architectures and stream-dataflow architecture in paper*

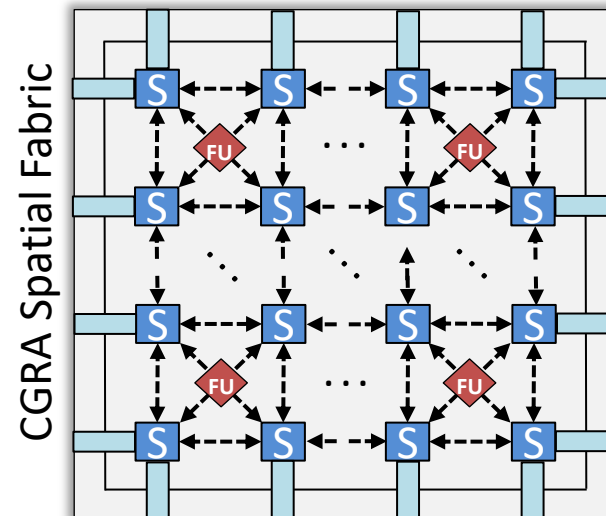# Stream-Dataflow Accelerator Architecture

—— 512b   ‐‐‐‐ 64b

# Stream-Dataflow Accelerator Architecture

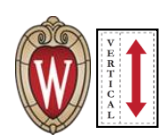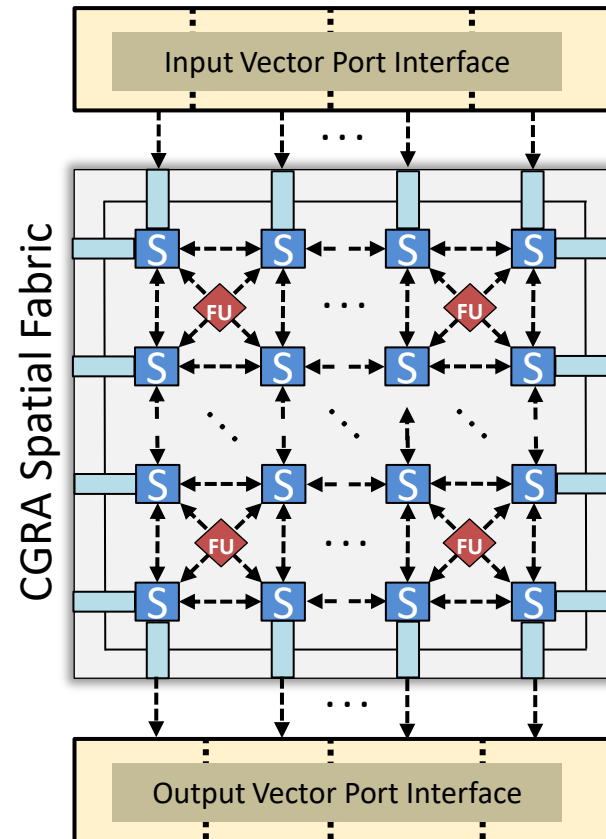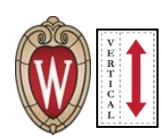—— 512b  ---- 64b

## Dataflow:

- Coarse grained reconfigurable architecture (CGRA) for data parallel execution

# Stream-Dataflow Accelerator Architecture

—— 512b ---- 64b

## Dataflow:

- Coarse grained reconfigurable architecture (CGRA) for data parallel execution

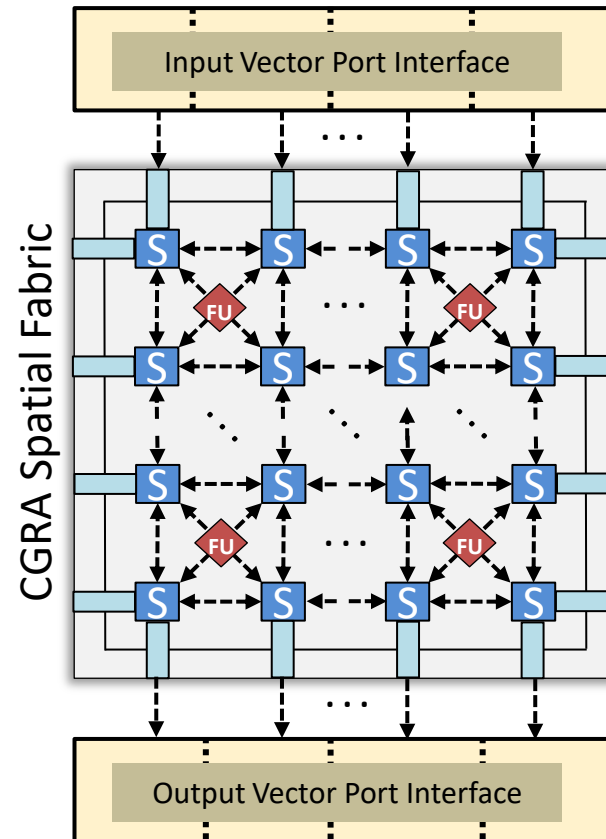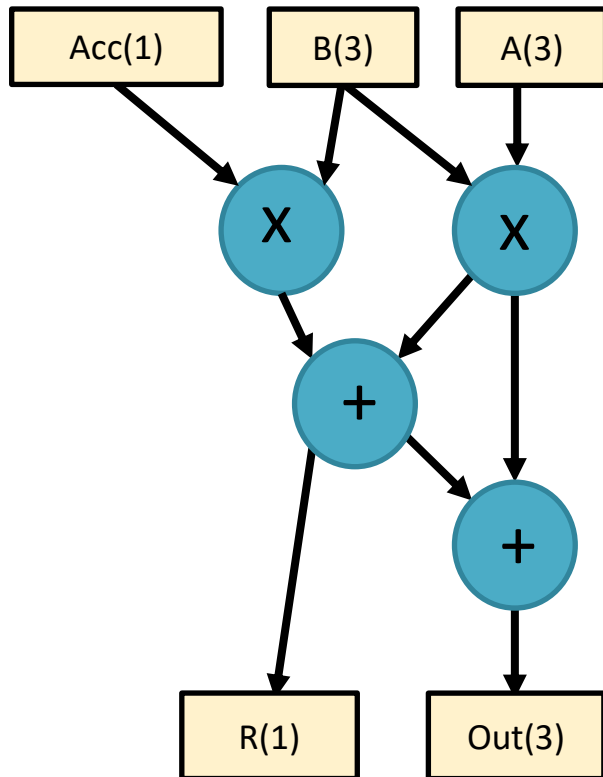- Direct vector port interface into and out of CGRA for vector execution

# Stream-Dataflow Accelerator Architecture

— 512b  ---- 64b

## Dataflow:

- Coarse grained reconfigurable architecture (CGRA) for data parallel execution

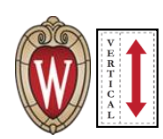- Direct vector port interface into and out of CGRA for vector execution

# Stream-Dataflow Accelerator Architecture

— 512b ---- 64b

## Dataflow:

- Coarse grained reconfigurable architecture (CGRA) for data parallel execution

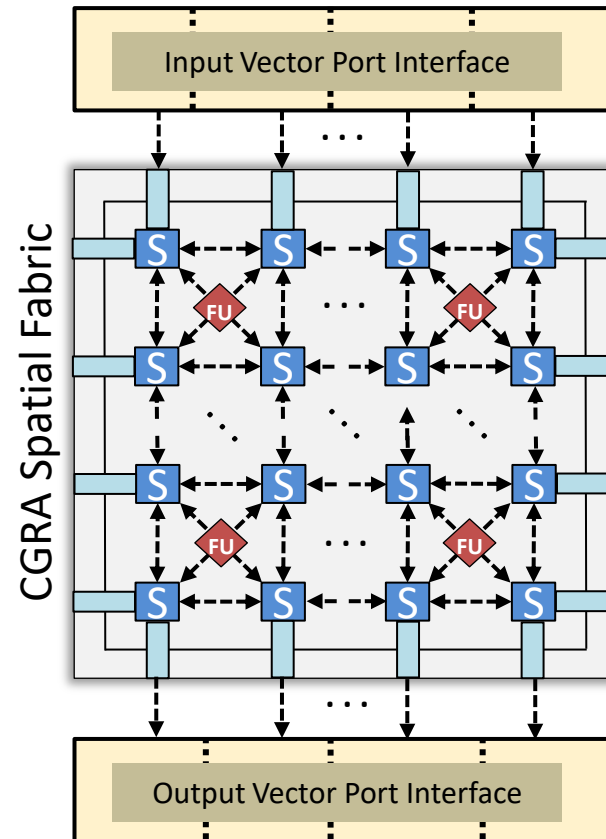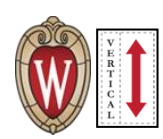- Direct vector port interface into and out of CGRA for vector execution

# Stream-Dataflow Accelerator Architecture

—— 512b  ---- 64b

## Dataflow:

- Coarse grained reconfigurable architecture (CGRA) for data parallel execution

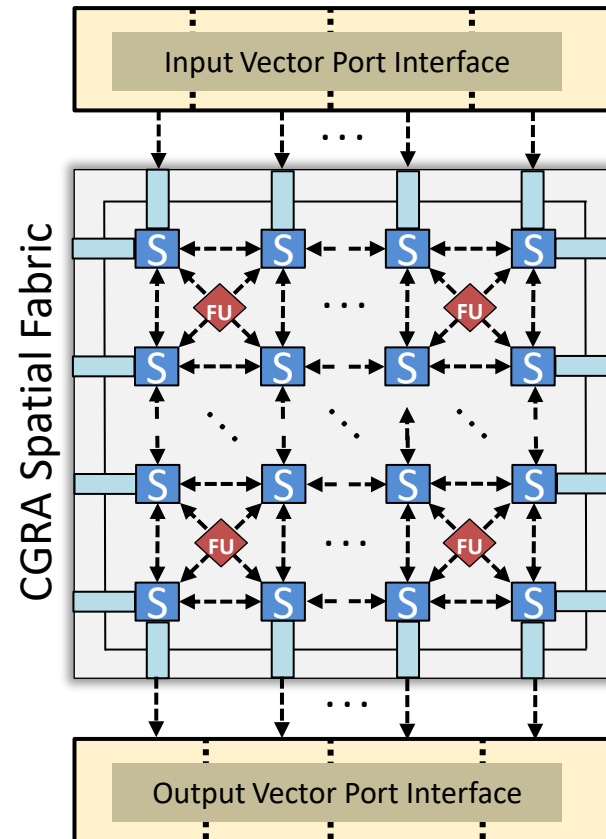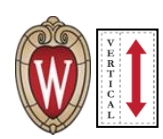- Direct vector port interface into and out of CGRA for vector execution

## Stream Interface:
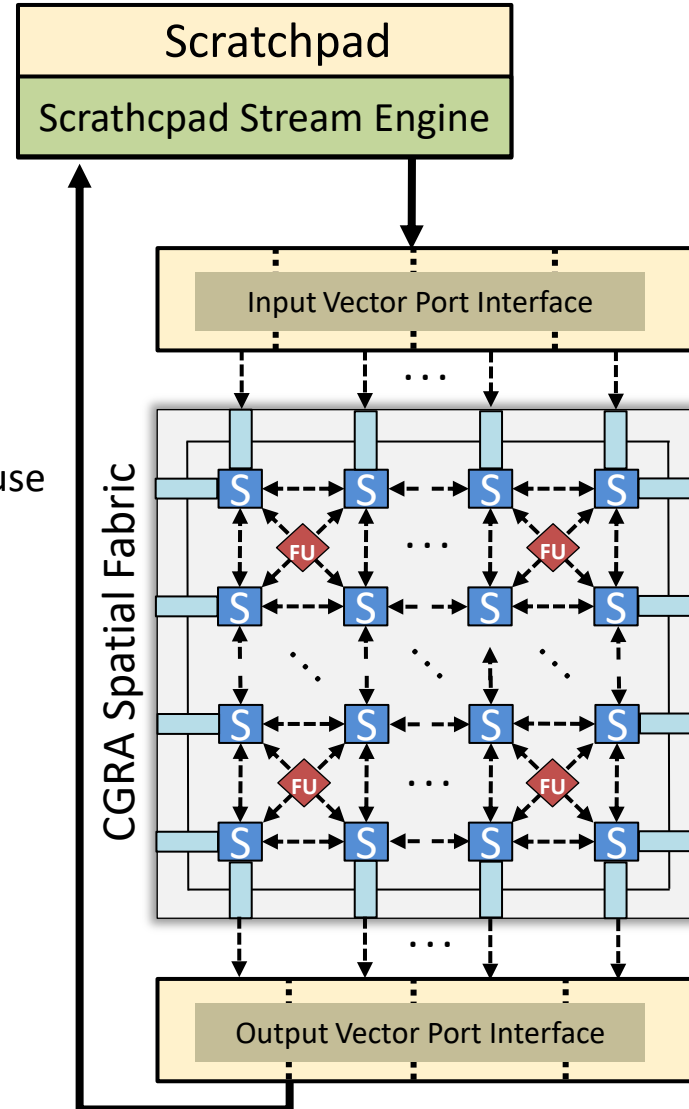
# Stream-Dataflow Accelerator Architecture

—— 512b    ---- 64b

## Dataflow:

- Coarse grained reconfigurable architecture (CGRA) for data parallel execution

- Direct vector port interface into and out of CGRA for vector execution

## Stream Interface:

- Programmable scratchpad and supporting stream-engine for data-locality and data-reuse

# Stream-Dataflow Accelerator Architecture
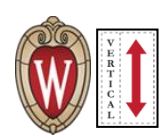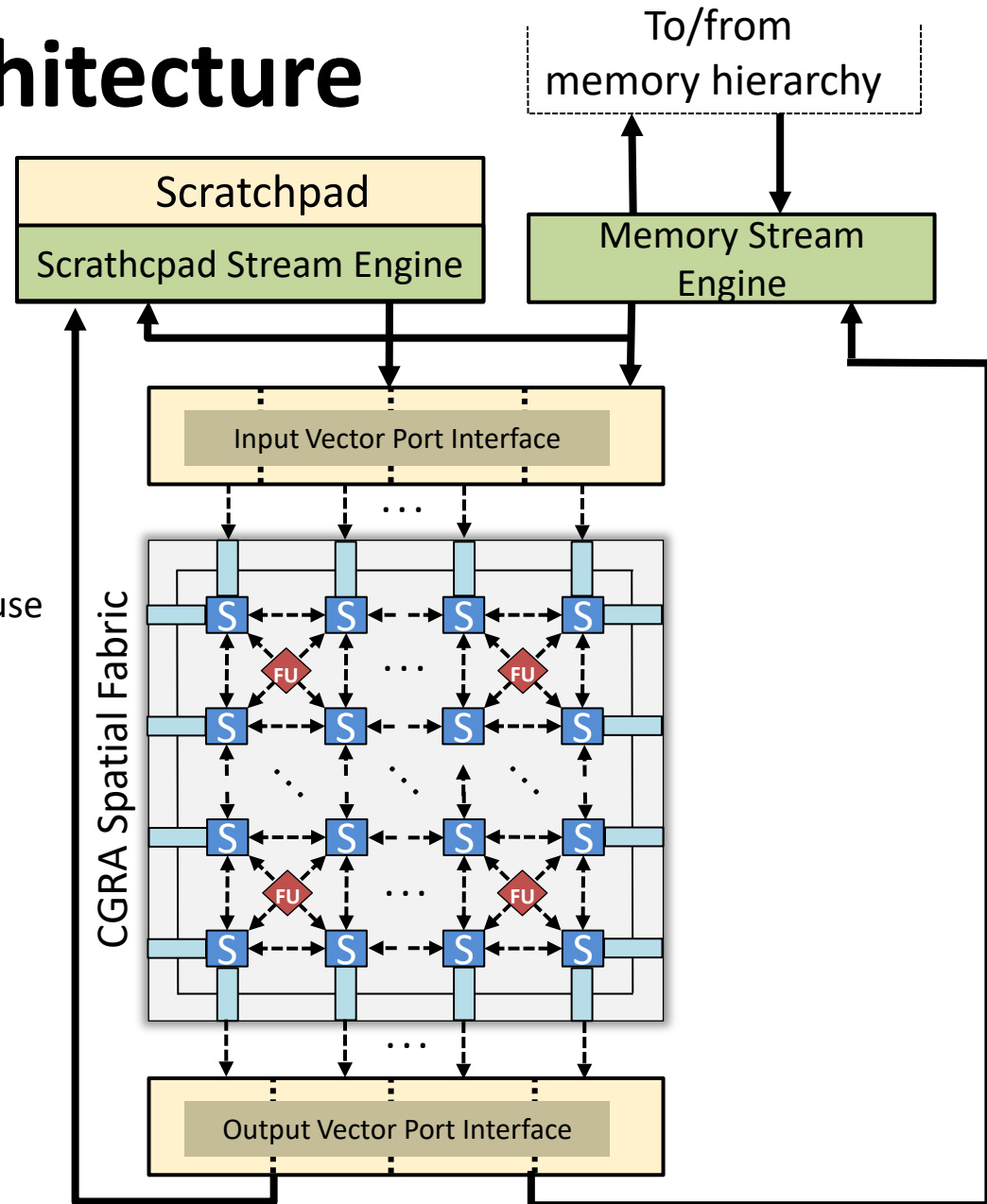
— 512b ---- 64b

## Dataflow:

- Coarse grained reconfigurable architecture (CGRA) for data parallel execution

- Direct vector port interface into and out of CGRA for vector execution

## Stream Interface:

- Programmable scratchpad and supporting stream-engine for data-locality and data-reuse

- Memory stream-engine to facilitate data streaming in and out of the accelerator



To/from memory hierarchy

Scratchpad

Scrathcpad Stream Engine

Memory Stream Engine

Input Vector Port Interface

CGRA Spatial Fabric

Output Vector Port Interface

# Stream-Dataflow Accelerator Architecture
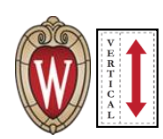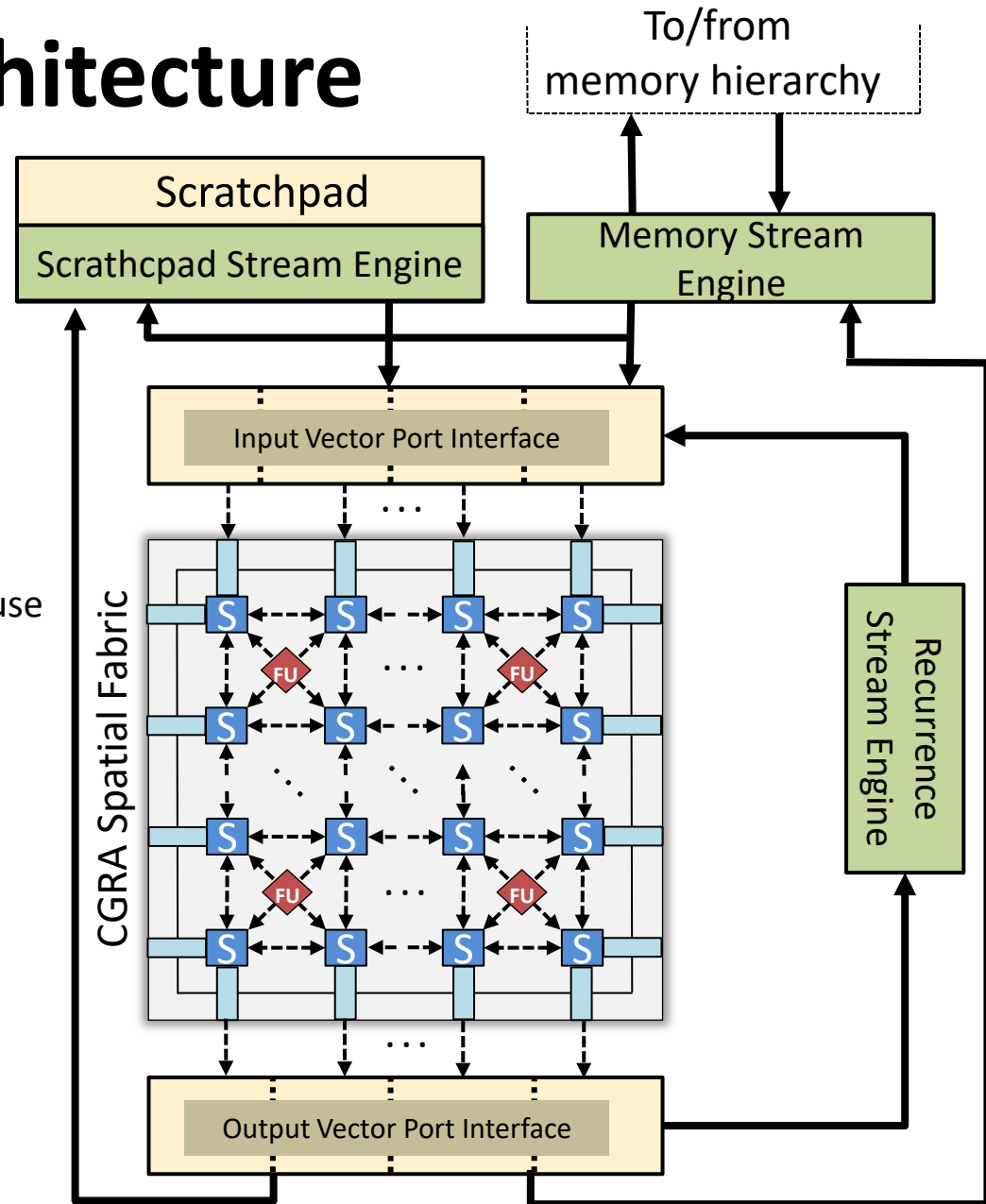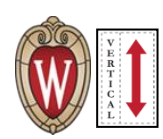
—— 512b ---- 64b

## Dataflow:

- Coarse grained reconfigurable architecture (CGRA) for data parallel execution

- Direct vector port interface into and out of CGRA for vector execution

## Stream Interface:

- Programmable scratchpad and supporting stream-engine for data-locality and data-reuse

- Memory stream-engine to facilitate data streaming in and out of the accelerator

- Recurrence stream-engine to support recurrent data stream



To/from memory hierarchy

Scratchpad

Scrathcpad Stream Engine

Memory Stream Engine

Input Vector Port Interface

CGRA Spatial Fabric

Recurrence Stream Engine

Output Vector Port Interface

# Stream-Dataflow Accelerator Architecture

— 512b ---- 64b

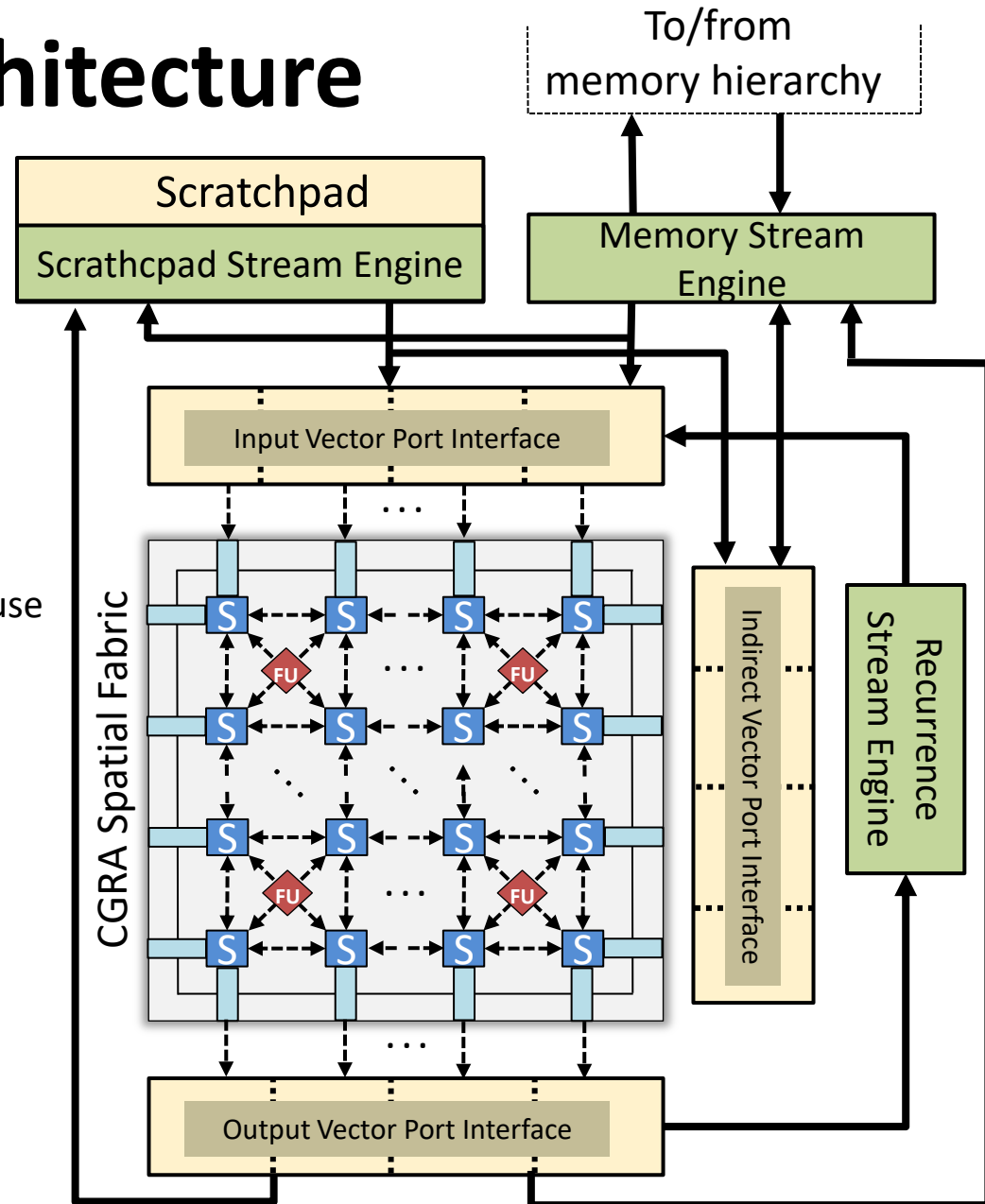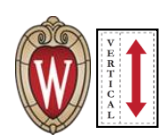## Dataflow:

- Coarse grained reconfigurable architecture (CGRA) for data parallel execution

- Direct vector port interface into and out of CGRA for vector execution

## Stream Interface:

- Programmable scratchpad and supporting stream-engine for data-locality and data-reuse

- Memory stream-engine to facilitate data streaming in and out of the accelerator

- Recurrence stream-engine to support recurrent data stream

- Indirect vector port interface for streaming addresses (indirect load/stores)
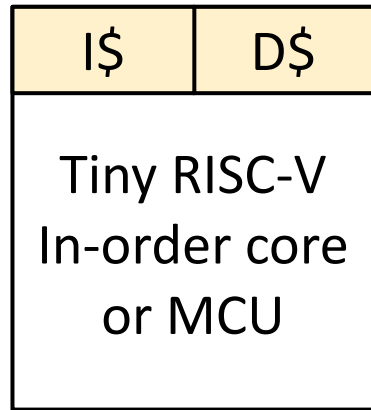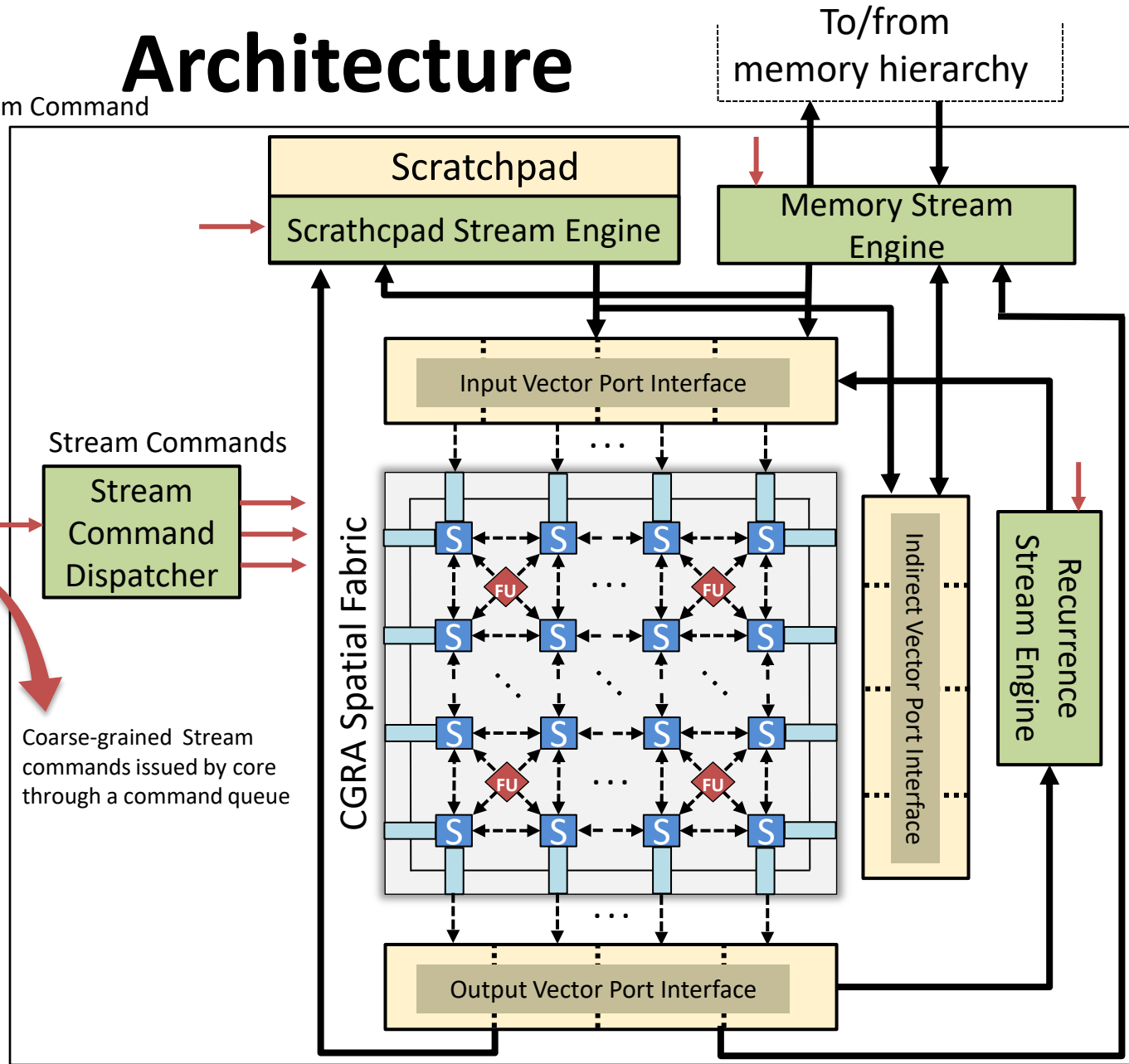


To/from memory hierarchy
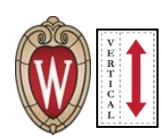
Scratchpad

Scratchcpad Stream Engine

Memory Stream Engine

Input Vector Port Interface

CGRA Spatial Fabric

Indirect Vector Port Interface

Recurrence Stream Engine

Output Vector Port Interface

# Stream-Dataflow Accelerator Architecture

To/from memory hierarchy

— 512b  ----- 64b  — Stream Command



Scratchpad

Scrathcpad Stream Engine

Memory Stream Engine

Input Vector Port Interface

I$  D$

Tiny RISC-V In-order core or MCU

Stream Commands

Stream Command Dispatcher

CGRA Spatial Fabric

Indirect Vector Port Interface

Recurrence Stream Engine

Coarse-grained Stream commands issued by core through a command queue

Output Vector Port Interface

- Stream command interface exposed to a general purpose programmable core

- Non-intrusive accelerator design
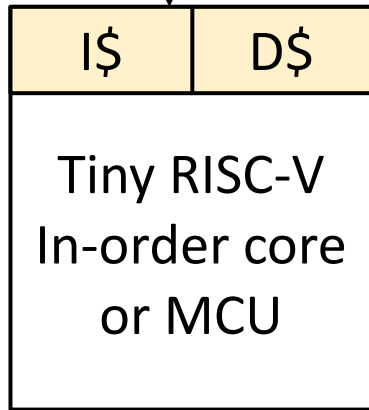
# Stream-Dataflow Accelerator Architecture

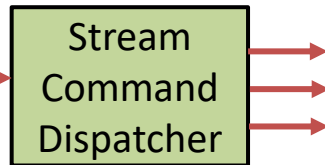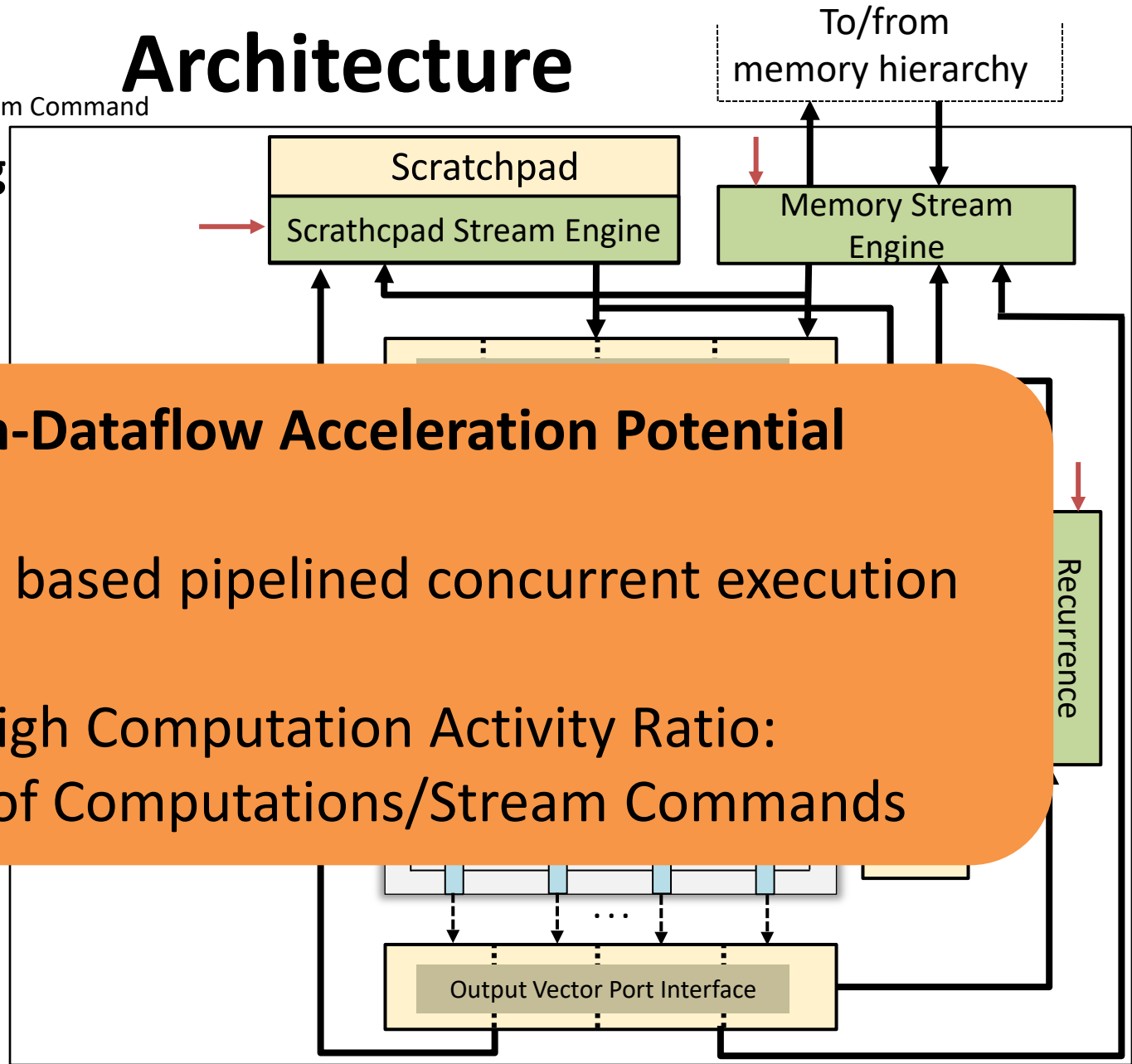To/from memory hierarchy

—— 512b  ----- 64b  —— Stream Command

## Stream ISA Encoding

```
Send a[0: N] → P1
Send b[0: N] → P2
Get P3 → c
```

I$   D$

Tiny RISC-V In-order core or MCU

**Scratchpad**

Scratchpad Stream Engine

**Memory Stream Engine**

Stream Commands

**Stream Command Dispatcher**

Input Vector Port Interface

CGRA Spatial Fabric

S S S S
FU      FU
S S S S
S S S S
FU      FU
S S S S

Indirect Vector Port Interface

Recurrence Stream Engine

Output Vector Port Interface

Coarse-grained Stream commands issued by core through a command queue

- Stream command interface exposed to a general purpose programmable core

- Non-intrusive accelerator design

# Stream-Dataflow Accelerator Architecture



**Stream ISA Encoding**

```
Send a[0: N] → P1
Send b[0: N] → P2
Get P3 → c
```

To/from memory hierarchy

Scratchpad

Scrathcpad Stream Engine

Memory Stream Engine

Recurrence

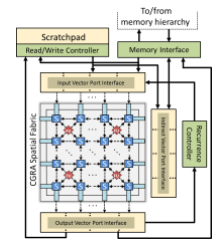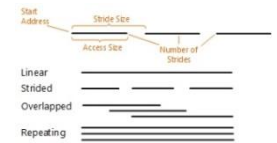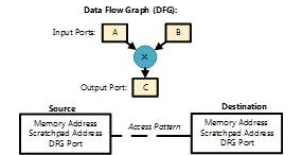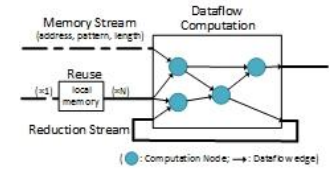Output Vector Port Interface

512b ---- 64b ── Stream Command

## Stream-Dataflow Acceleration Potential

*1.* Dataflow based pipelined concurrent execution

2. High Computation Activity Ratio:
Number of Computations/Stream Commands

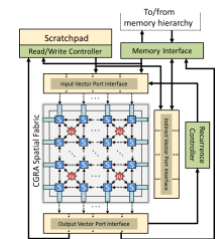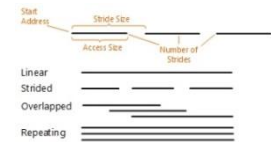- ...general purpose programmable core

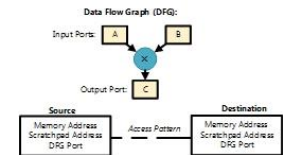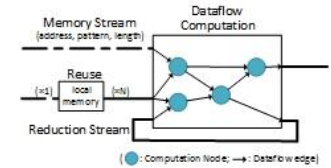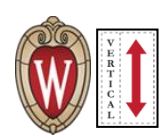- Non-intrusive accelerator design

# Outline

- Motivation and Overview

- Stream-Dataflow Execution Model

- Hardware-Software Interface and Example program

- Stream-Dataflow Accelerator Architecture

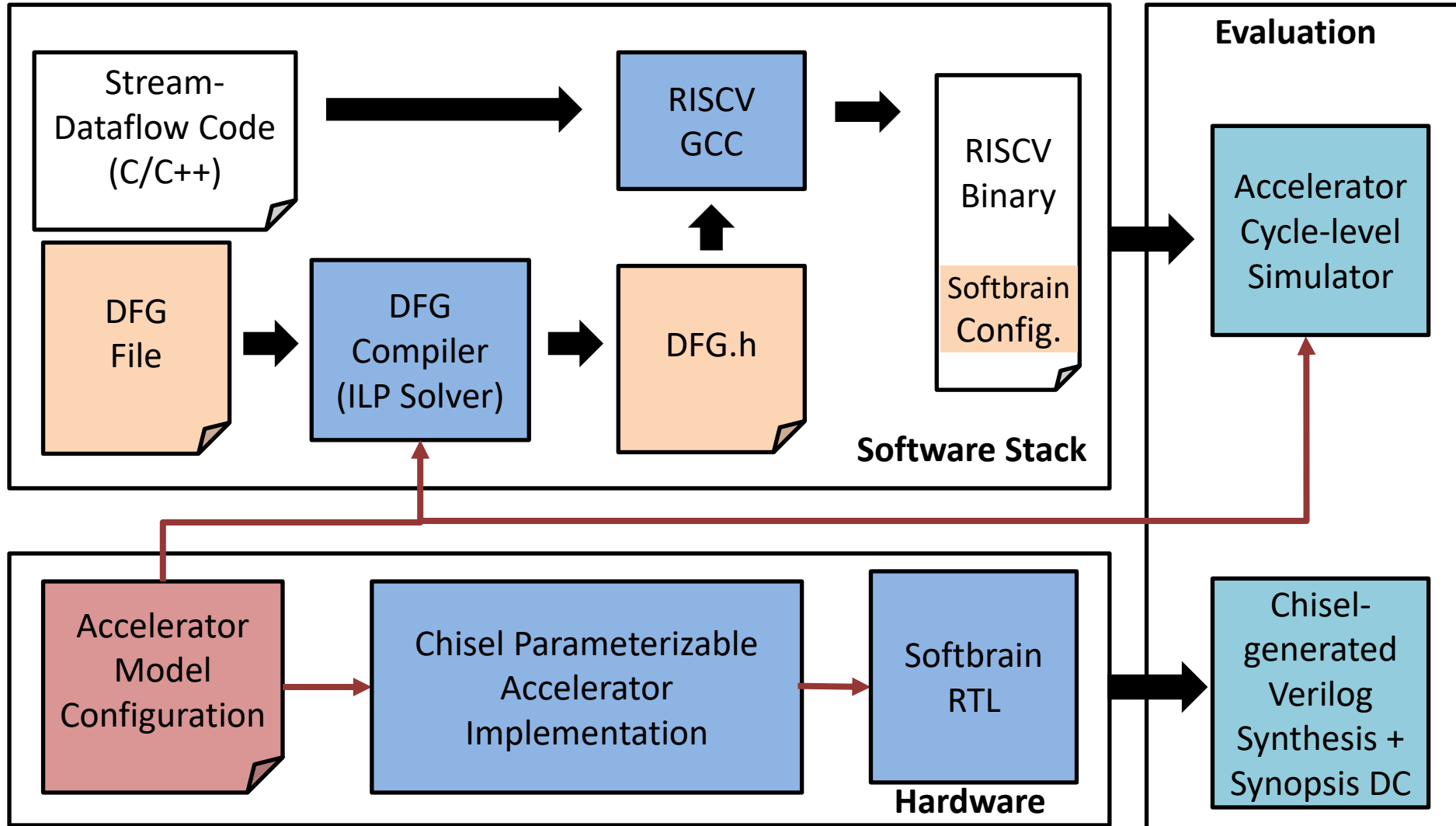- **Evaluation and Results**

# Outline

- Motivation and Overview

- Stream-Dataflow Execution Model

- Hardware-Software Interface and Example program

- Stream-Dataflow Accelerator Architecture

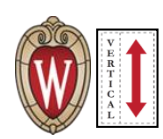- **Evaluation and Results**

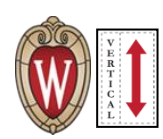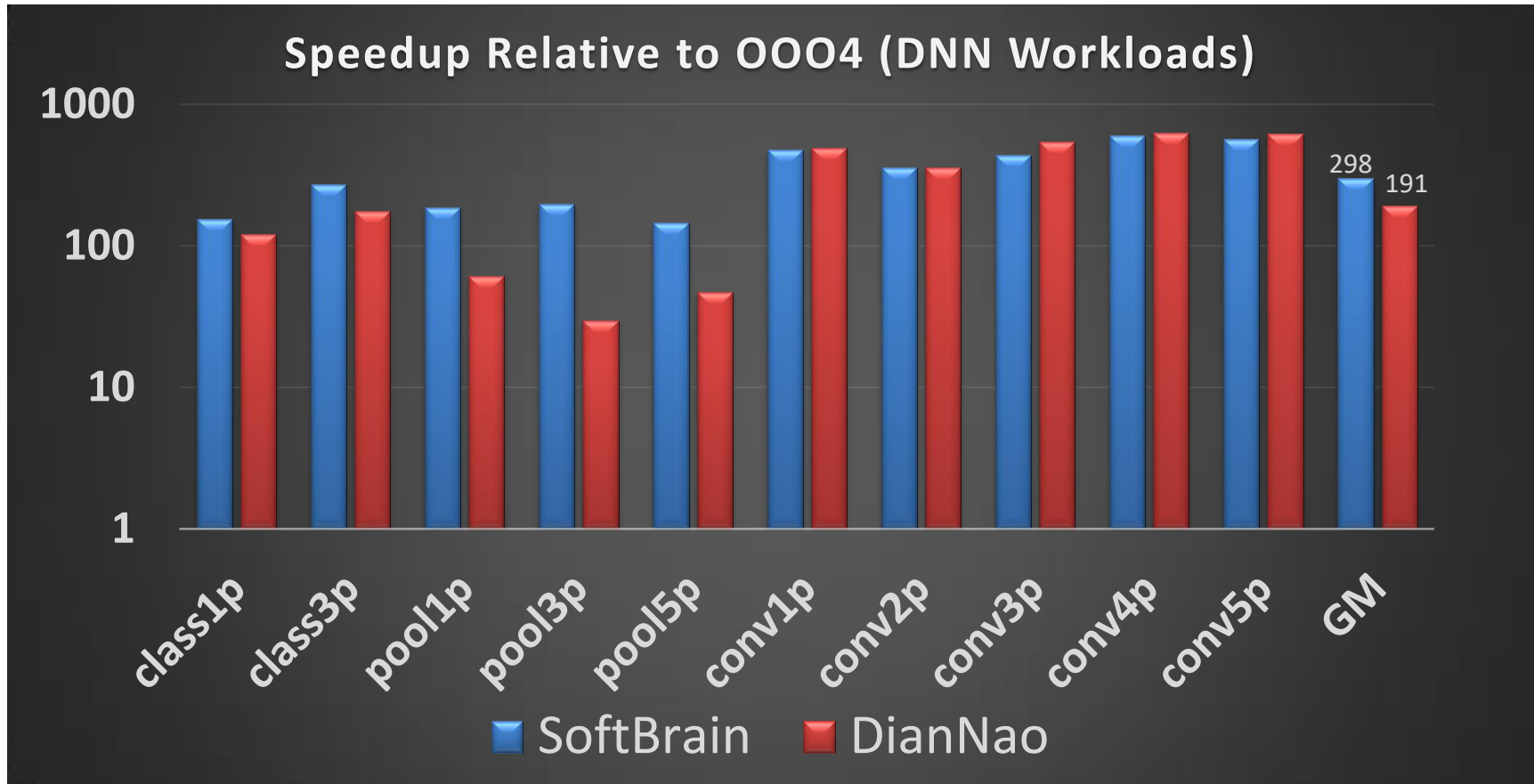# Stream-Dataflow Implementation: *Softbrain*
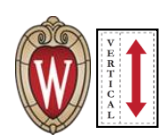
# Evaluation Methodology

- ## Workloads

  - ❑ Deep Neural Networks (DNN) – For domain provisioned comparison

  - ❑ Machsuite Accelerator Workloads – For comparison with application specific accelerators

- ## Comparison

  - ❑ Domain Provisioned Softbrain vs. DianNao DSA

  - ❑ Broadly provisioned Softbrain vs. ASIC design points – *Aladdin* generated performance, power and area
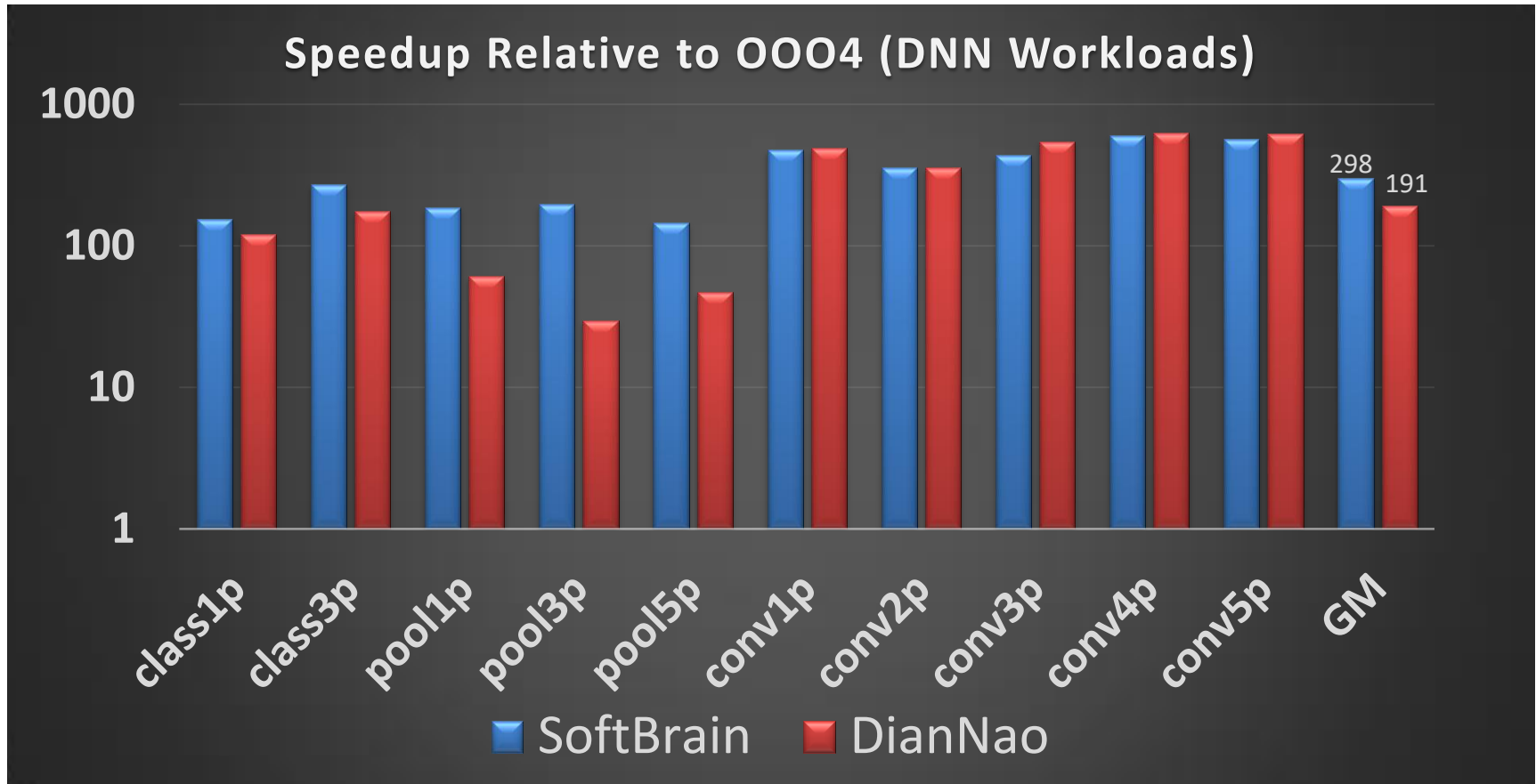
# Domain-Specific Accelerator Comparison (Softbrain vs DianNao)

Speedup Relative to OOO4 (DNN Workloads)

# Domain-Specific Accelerator Comparison (Softbrain vs DianNao)



**Speedup Relative to OOO4 (DNN Workloads)**

Categories: class1p, class3p, pool1p, pool3p, pool5p, conv1p, conv2p, conv3p, conv4p, conv5p, GM

GM: SoftBrain 298, DianNao 191

Legend: SoftBrain (blue), DianNao (red)

DianNao Area: **2.16 mm²**          Softbrain Area: **3.76 mm²**
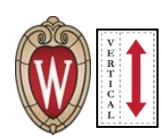
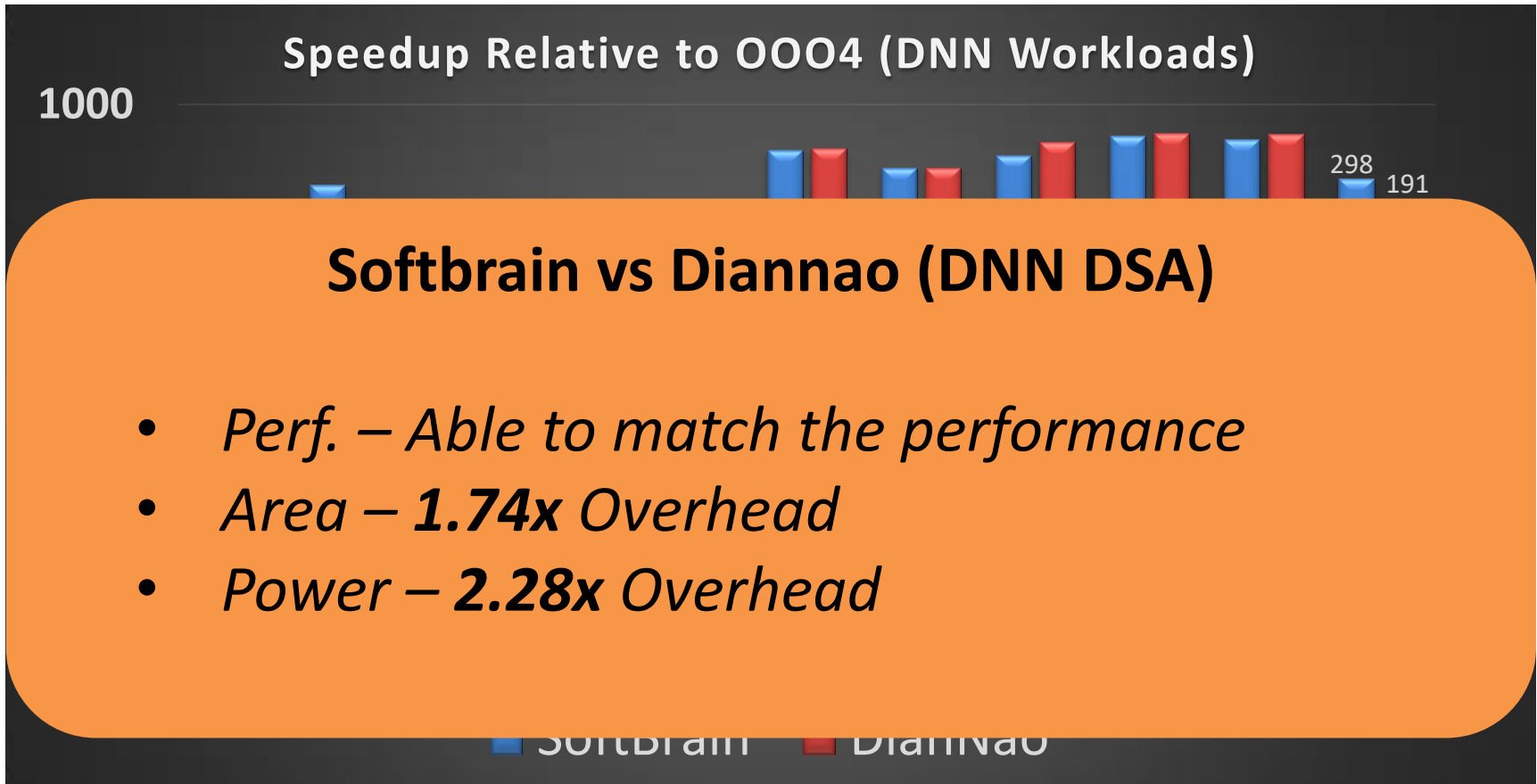DianNao Power: **420 mW**          Softbrain Power: **950 mW**

# Domain-Specific Accelerator Comparison (Softbrain vs DianNao)



**Speedup Relative to OOO4 (DNN Workloads)**

1000

298    191

## Softbrain vs Diannao (DNN DSA)

- *Perf. – Able to match the performance*
- *Area – **1.74x** Overhead*
- *Power – **2.28x** Overhead*
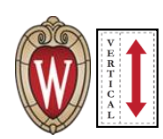
SoftBrain    DianNao

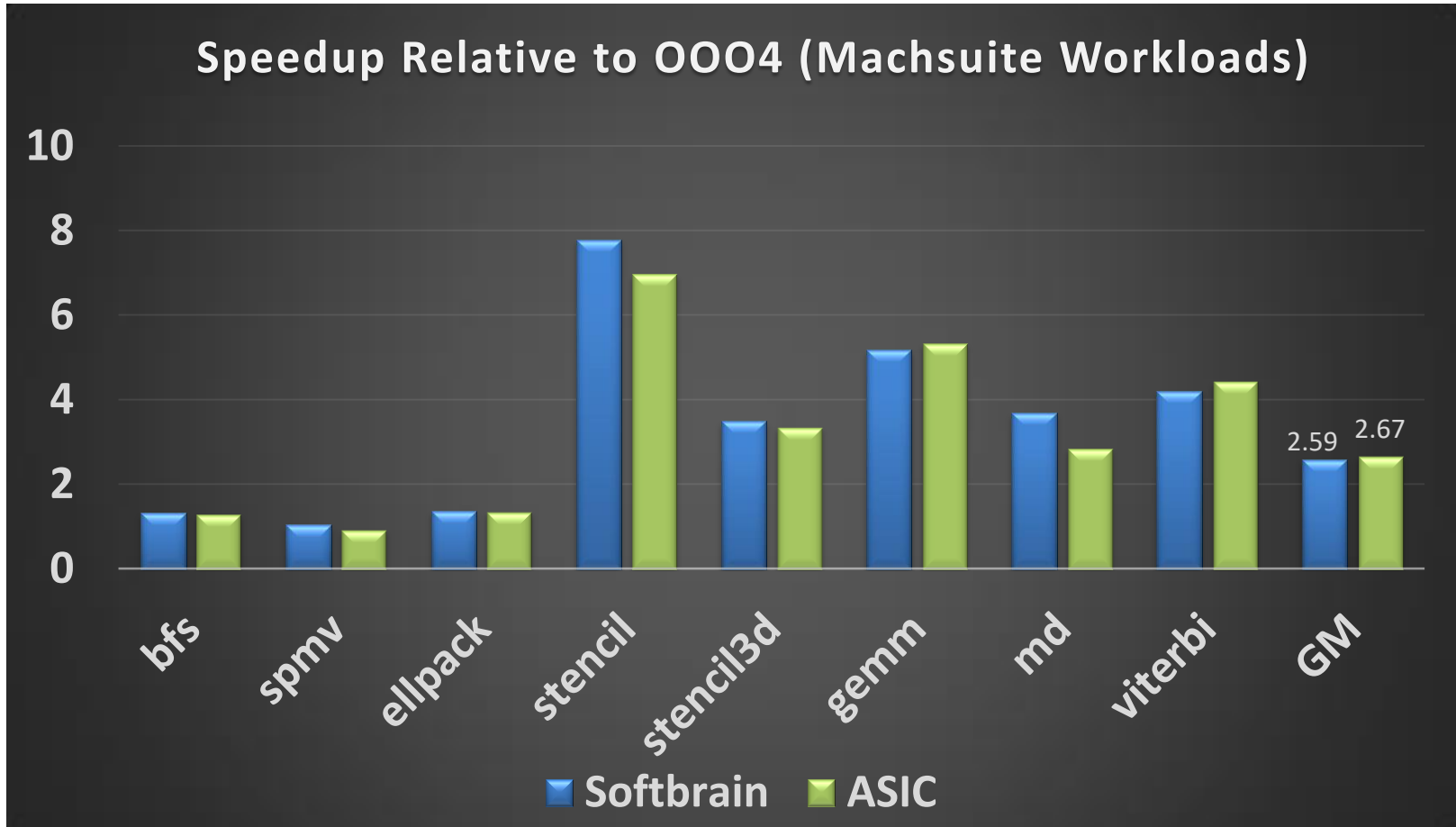DianNao Area: **2.16 mm²**    Softbrain Area: **3.76 mm²**

DianNao Power: **420 mW**    Softbrain Power: **950 mW**

# Softbrain vs ASIC Designs Comparison
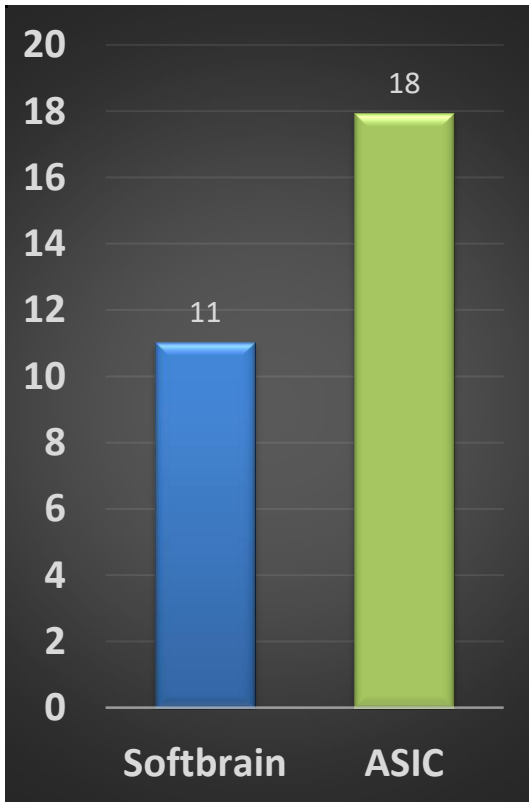


**Speedup Relative to OOO4 (Machsuite Workloads)**

Legend: ■ Softbrain ■ ASIC

GM: Softbrain 2.59, ASIC 2.67

Aladdin* generated ASIC design points – Resources constrained to be in ~15% of Softbrain Perf. to do iso-performance analysis

*Aladdin: A Pre-RTL, Power-Performance Accelerator Simulator Enabling Large Design Space Exploration of Customized Architectures. Sophia Shao , .et. al
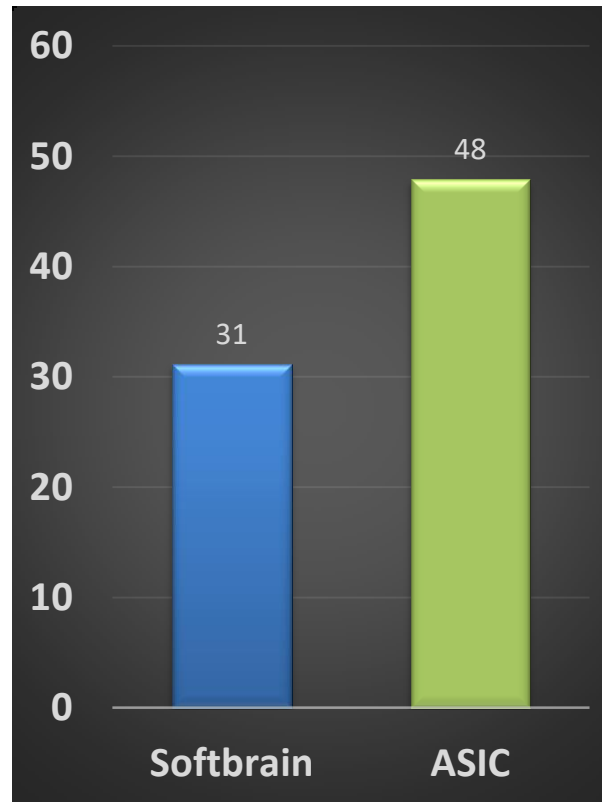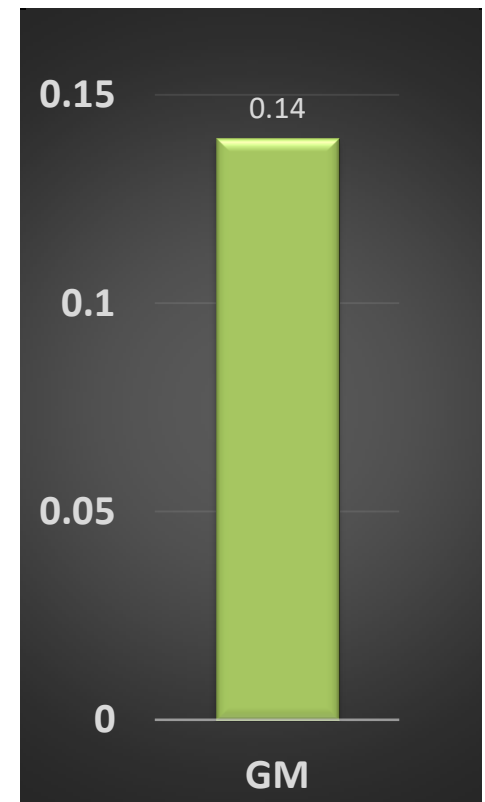
# Softbrain vs ASIC Comparison

# Softbrain vs ASIC Comparison

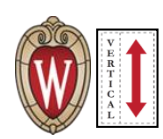| Power Efficiency Relative to OOO4 (GM) | Energy Efficiency Relative to OOO4 (GM) | ASIC Area Relative to Softbrain (GM) |
|---|---|---|

## Softbrain vs ASIC designs

- *Perf. – Able to match the performance*
- *Power – ~**1.6x** overhead*
- *Energy Efficiency – ~**1.5x** overhead*
- *Area – ~**8x** overhead\**

*\*All 8 ASICs combined → 2.15x more area than Softbrain*
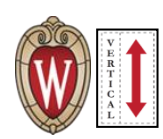
# Conclusion

# Conclusion

- Stream-Dataflow Acceleration
  - Stream-Dataflow **Execution Model** – Abstracts typical accelerator computation phases using a dataflow graph
  - Stream-Dataflow **ISA Encoding** and **Hardware-Software Interface** – Exposes parallelism available in these phases
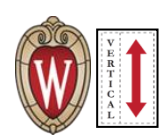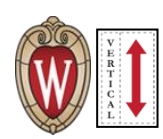
# Conclusion

- Stream-Dataflow Acceleration
  - Stream-Dataflow **Execution Model** – Abstracts typical accelerator computation phases using a dataflow graph
  - Stream-Dataflow **ISA Encoding** and **Hardware-Software Interface** – Exposes parallelism available in these phases

- Stream-Dataflow Accelerator Architecture
  - CGRA and vector ports for pipelined vector-dataflow computation
  - Highly parallel stream-engines for low-power stream communication

# Conclusion

- Stream-Dataflow Acceleration
  - Stream-Dataflow **Execution Model** – Abstracts typical accelerator computation phases using a dataflow graph
  - Stream-Dataflow **ISA Encoding** and **Hardware-Software Interface** – Exposes parallelism available in these phases

- Stream-Dataflow Accelerator Architecture
  - CGRA and vector ports for pipelined vector-dataflow computation
  - Highly parallel stream-engines for low-power stream communication

- Stream-Dataflow Prototype & Implementation – Softbrain
  - Matches performance of domain provisioned accelerator (DianNao DSA) with **~2x** overheads in area and power
  - Compared to application specific designs (ASICs), Softbrain has **~2x** overheads in power and **~8x** in area
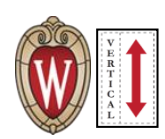
# Conclusion

- Stream-Dataflow Acceleration
  - Stream-Dataflow **Execution Model** – Abstracts typical accelerator computation phases using a dataflow graph
  - Stream-Dataflow **ISA Encoding** and **Hardware-Software Interface** –
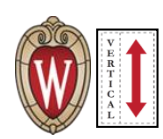
- 

- 
  - Compared to application specific designs (ASICs), Softbrain has **~2x** overheads in power and **~8x** in area

**Getting There !!**

*A good enabler for exploring general purpose programmable hardware acceleration ….*
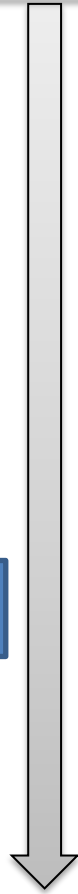
# Backup

# Traditional Arch.

# Accelerator (DSA)

Programs

General Language

Compiler

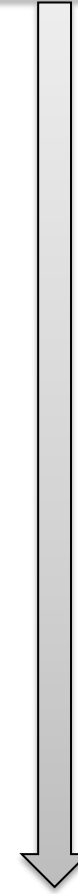General ISA

General Purpose Hardware

Domain-Specific Programs

**Tiny H/W-S/W Interface**

Application/Domain Specific Hardware

UCLA

**Traditional Arch.**

**Accelerator (DSA)**

Programs

Domain-Specific Programs

General Language

**Tiny H/W-S/W Interface**

Compiler

General ISA

General Purpose Hardware

Application/Domain Specific Hardware

10-100x Performance/Power or Performance/Area (completely lose generality/programmability)

**Traditional Arch.**

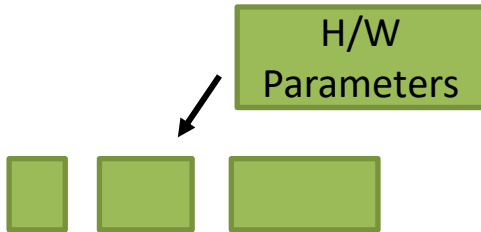**Progammable Hardware Accelerator**

**Accelerator (DSA)**
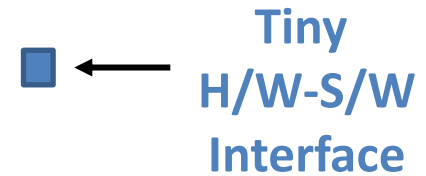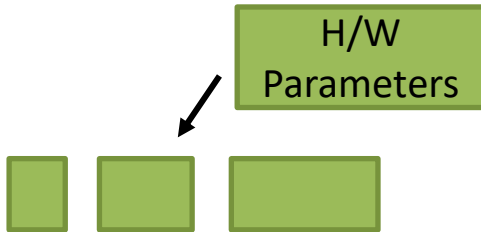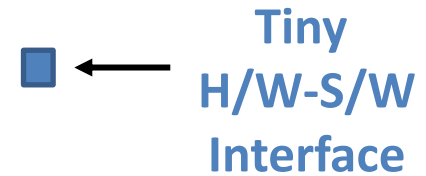
Programs

General Language

Compiler

General ISA

General Purpose Hardware

Programs ("Specialized")

H/W Parameters

Re-Configurable Hardware

Domain-Specific Programs

Tiny H/W-S/W Interface

Application/Domain Specific Hardware

Can the specialized programs be adapted in a domain-agnostic way with this interface?
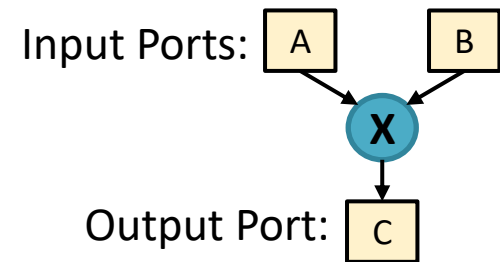
# Stream-Dataflow Execution Model Detailed Example

# Stream-Dataflow Execution Model Detailed Example
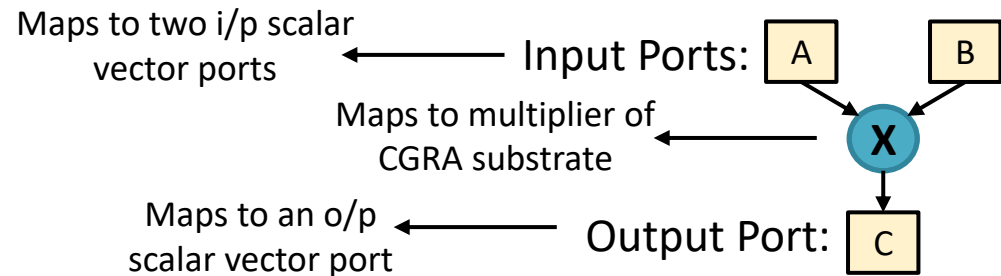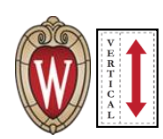
C[i] = A[i] * B[i]

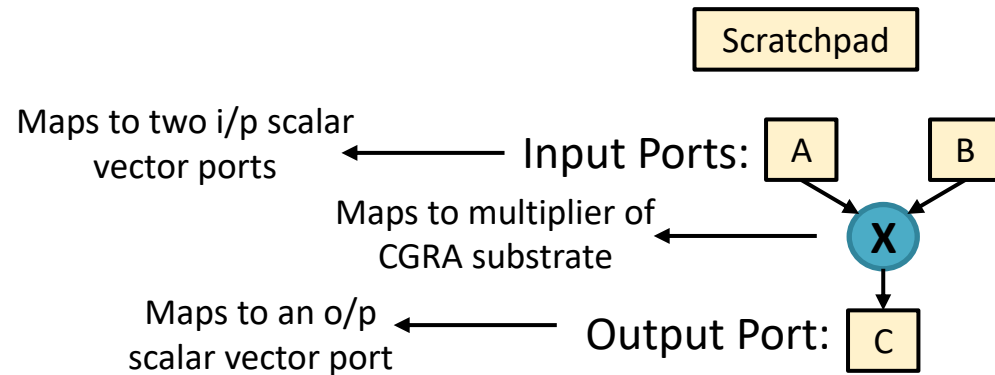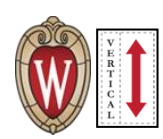# Stream-Dataflow Execution Model Detailed Example

$$C[i] = A[i] * B[i]$$

Maps to two i/p scalar vector ports ← Input Ports: [A] [B]

Maps to multiplier of CGRA substrate ← (X)

Maps to an o/p scalar vector port ← Output Port: [C]

# Stream-Dataflow Execution Model Detailed Example

**UCLA**

$$C[i] = A[i] * B[i]$$

Scratchpad

Maps to two i/p scalar vector ports ← Input Ports: A  B

Maps to multiplier of CGRA substrate ← X

Maps to an o/p scalar vector port ← Output Port: C

# Stream-Dataflow Execution Model Detailed Example
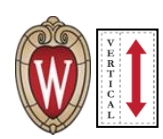
$$C[i] = A[i] * B[i]$$

**Stream Commands**

Time →

Program Order ↓

# Stream-Dataflow Execution Model Detailed Example

**Legend:**

| | | | |
|---|---|---|---|
| Enqueued | □ | Barrier | ⊙ |
| Dispatched | ○ | Dependency | → |
| Resource idle | ······· | Iter. boundary | ╱ |
| Resource in use | —— | | |
| All data at dest. | ● | | |

**C[i] = A[i] * B[i]**

**Stream Commands**

Time →

Program Order ↓

Scratchpad

A    B

X

C

**CGRA fabric state**  · · · · · · · · · · · ·

**Low-power core state**  ——

# Stream-Dataflow Execution Model Detailed Example

**Legend:**

| | | | |
|---|---|---|---|
| Enqueued | □ | Barrier | ⊙ |
| Dispatched | ○ | Dependency | → |
| Resource idle | ······· | Iter. boundary | ╱ |
| Resource in use | —— | | |
| All data at dest. | ● | | |

$C[i] = A[i] * B[i]$

**Stream Commands**

Time →

C1) Mem → Scratch

Program Order ↓

Scratchpad

A    B

X

C

**CGRA fabric state**

**Low-power core state**

Command generation

# Stream-Dataflow Execution Model Detailed Example

**UCLA**

## Legend:

| | | | |
|---|---|---|---|
| Enqueued | ☐ | Barrier | ⊙ |
| Dispatched | ○ | Dependency | → |
| Resource idle | ····· | Iter. boundary | ╱ |
| Resource in use | —— | | |
| All data at dest. | ● | | |

$C[i] = A[i] * B[i]$

Time →

**Stream Commands**

C1) Mem → Scratch

C2) Scratch Wr Barrier

C3) Scratch → Port A

**Program Order**

Scratchpad

A    B

X

C

**CGRA fabric state**

**Low-power core state**

Command generation

# Stream-Dataflow Execution Model Detailed Example

**Legend:**

| | | | |
|---|---|---|---|
| Enqueued | □ | Barrier | ⊙ |
| Dispatched | ○ | Dependency | → |
| Resource idle | ······ | Iter. boundary | ╱ |
| Resource in use | —— | | |
| All data at dest. | ● | | |

$C[i] = A[i] * B[i]$

**Stream Commands**  →  Time

C1) Mem → Scratch

C2) Scratch Wr Barrier

C3) Scratch → Port A

C4) Mem → Port B

Program Order

Scratchpad

A    B

X

C

**CGRA fabric state**

**Low-power core state**

Command generation

# Stream-Dataflow Execution Model Detailed Example

**Legend:**

| | | | |
|---|---|---|---|
| Enqueued | □ | Barrier | ☉ |
| Dispatched | ○ | Dependency | → |
| Resource idle | ······· | Iter. boundary | ╱ |
| Resource in use | —— | | |
| All data at dest. | ● | | |

$$C[i] = A[i] * B[i]$$

**Stream Commands**

Time →

C1) Mem → Scratch

C2) Scratch Wr Barrier

C3) Scratch → Port A

C4) Mem → Port B

C5) Port C → Mem

Program Order ↓

**CGRA fabric state**

**Low-power core state**

Command generation

# Stream-Dataflow Execution Model Detailed Example

**UCLA**

**Legend:**

Enqueued □    Barrier ⊙
Dispatched ○    Dependency →
Resource idle ·······    Iter. boundary ╱
Resource in use ──────
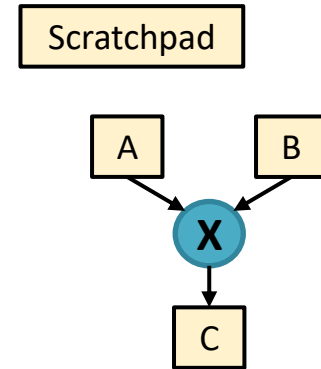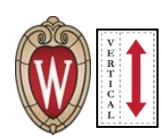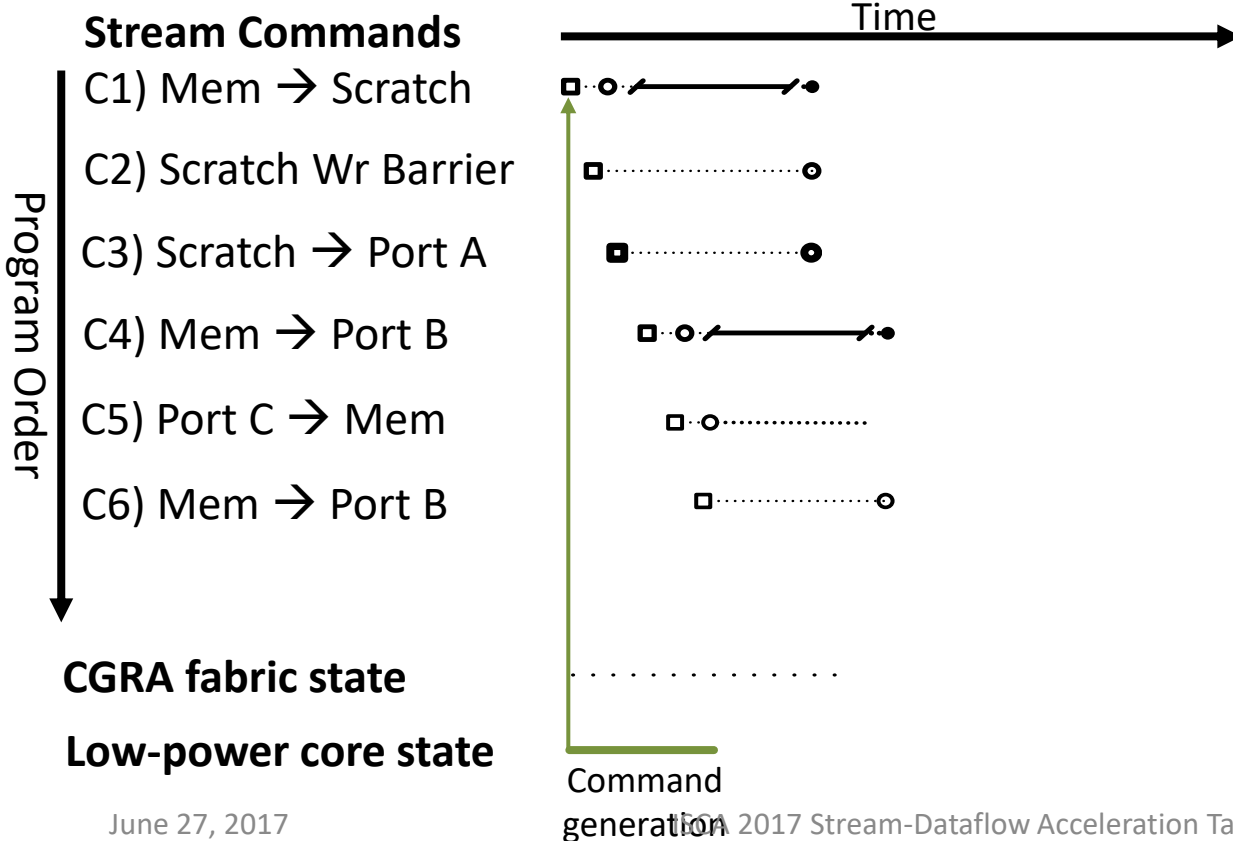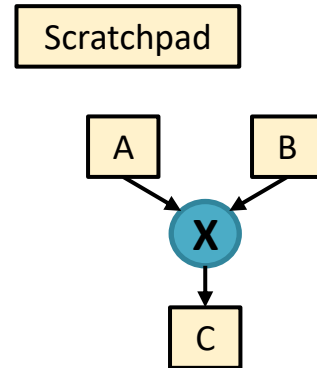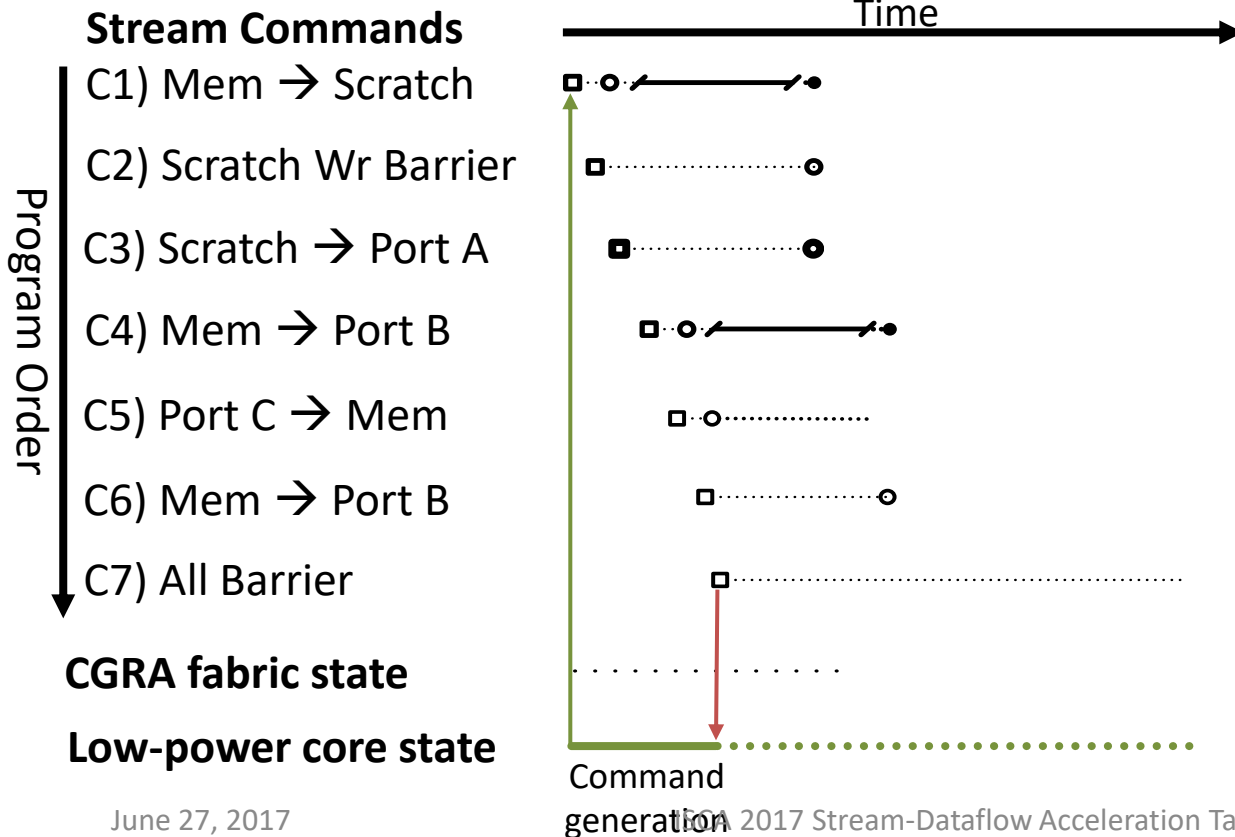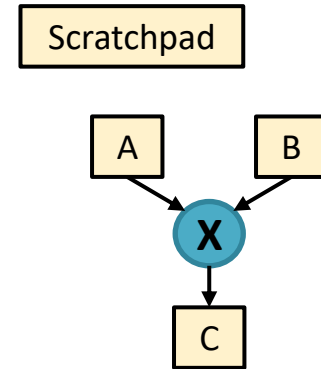All data at dest. ●

$C[i] = A[i] * B[i]$

**Stream Commands**    Time →

C1) Mem → Scratch

C2) Scratch Wr Barrier

C3) Scratch → Port A

C4) Mem → Port B

C5) Port C → Mem

C6) Mem → Port B

Program Order

Scratchpad

A    B

X

C

**CGRA fabric state**

**Low-power core state**

Command generation

# Stream-Dataflow Execution Model Detailed Example

**Legend:**

| | | | |
|---|---|---|---|
| Enqueued | □ | Barrier | ☉ |
| Dispatched | ○ | Dependency | → |
| Resource idle | ······· | Iter. boundary | ╱ |
| Resource in use | ——— | | |
| All data at dest. | ● | | |

$$C[i] = A[i] * B[i]$$

**Stream Commands**

Time →

**Program Order**

C1) Mem → Scratch

C2) Scratch Wr Barrier

C3) Scratch → Port A

C4) Mem → Port B

C5) Port C → Mem

C6) Mem → Port B

C7) All Barrier

**CGRA fabric state**

**Low-power core state**

Command generation

Scratchpad

A        B

X

C

# Stream-Dataflow Execution Model Detailed Example

**Legend:**

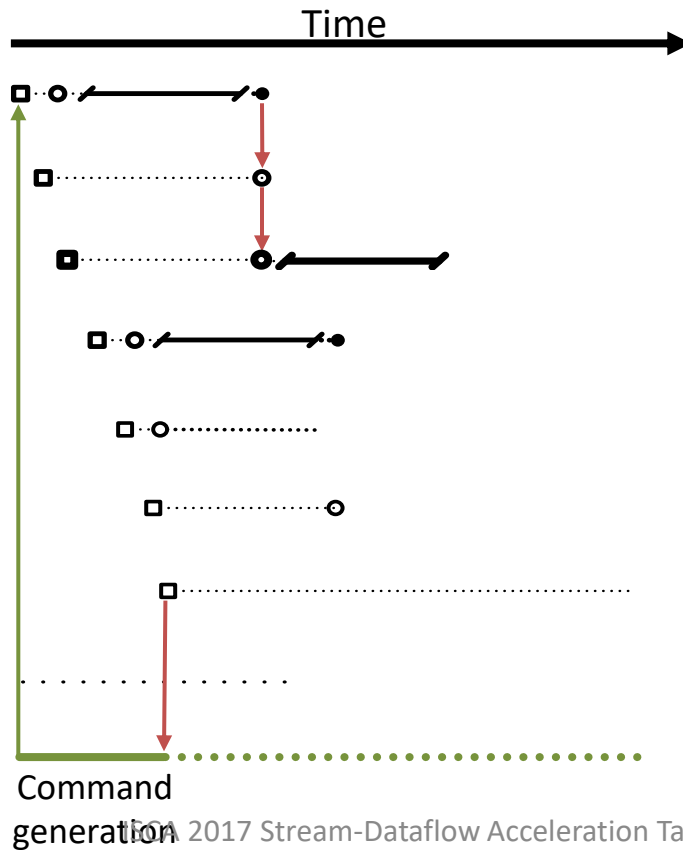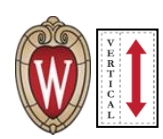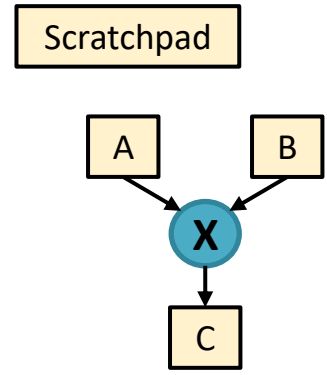| | | | |
|---|---|---|---|
| Enqueued | ☐ | Barrier | ⊙ |
| Dispatched | ○ | Dependency | → |
| Resource idle | ⋯⋯ | Iter. boundary | ╱ |
| Resource in use | ── | | |
| All data at dest. | ● | | |

$C[i] = A[i] * B[i]$

**Time** →

**Stream Commands**

- C1) Mem → Scratch
- C2) Scratch Wr Barrier
- C3) Scratch → Port A
- C4) Mem → Port B
- C5) Port C → Mem
- C6) Mem → Port B
- C7) All Barrier

**Program Order** ↓

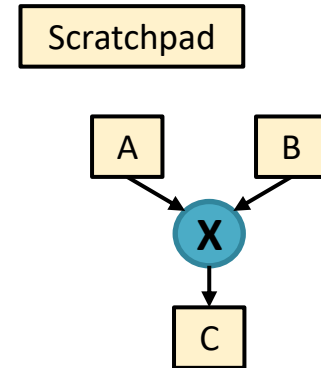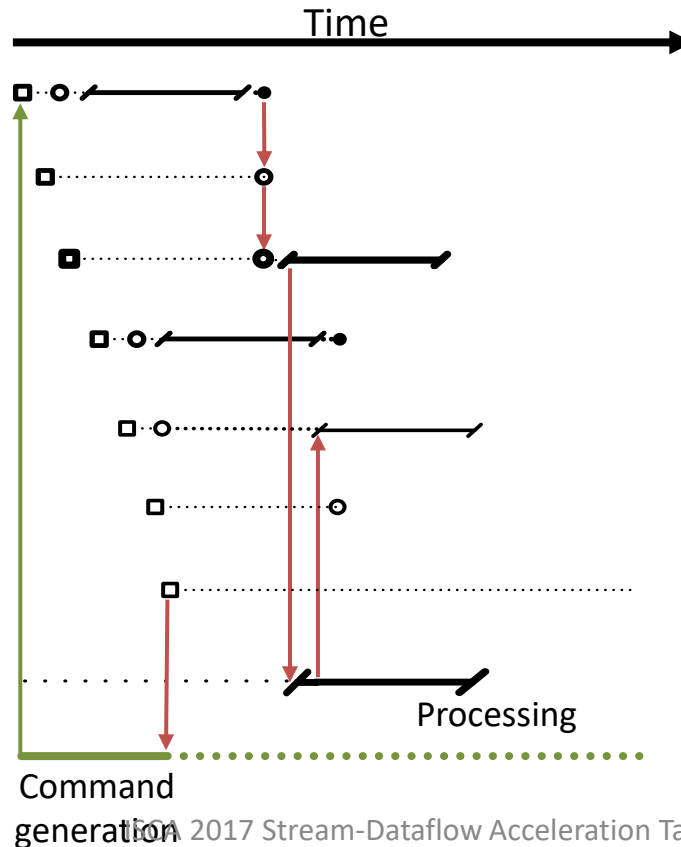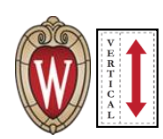**CGRA fabric state**

**Low-power core state**

Command generation

Scratchpad

A    B

X

C

**UCLA**

**Legend:**

| | | | |
|---|---|---|---|
| Enqueued | □ | Barrier | ☉ |
| Dispatched | ○ | Dependency | → |
| Resource idle | ······· | Iter. boundary | ╱ |
| Resource in use | ——— | | |
| All data at dest. | ● | | |

$C[i] = A[i] * B[i]$

**Stream Commands**

Time →

**Program Order**

C1) Mem → Scratch

C2) Scratch Wr Barrier

C3) Scratch → Port A

C4) Mem → Port B

C5) Port C → Mem

C6) Mem → Port B

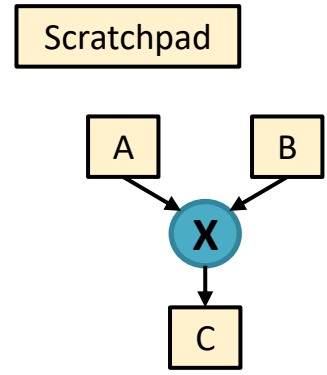C7) All Barrier

**CGRA fabric state**

**Low-power core state**

Command generation

Scratchpad

A          B

X

C

# Stream-Dataflow Execution Model Detailed Example

**Legend:**

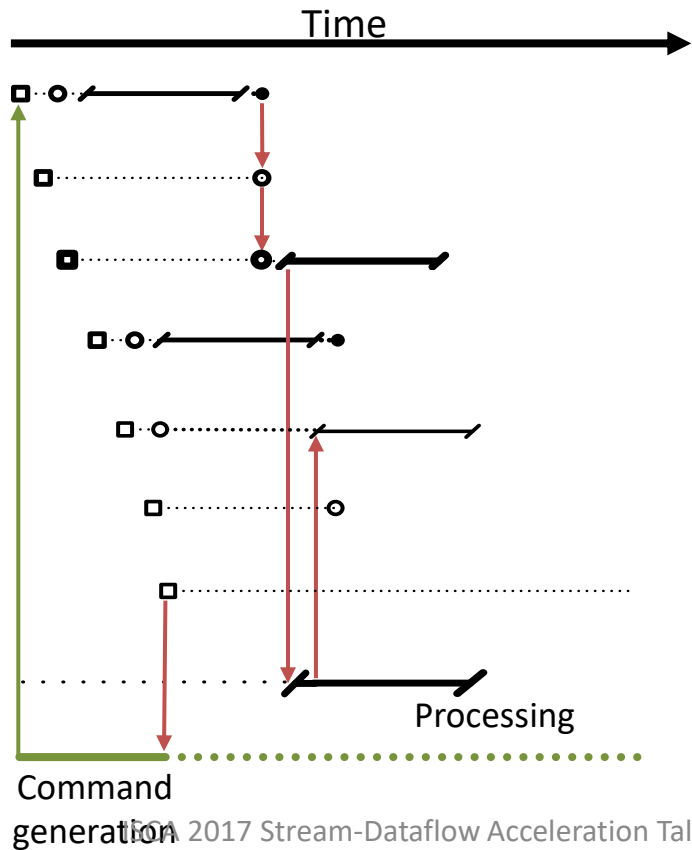| | | | |
|---|---|---|---|
| Enqueued | □ | Barrier | ☉ |
| Dispatched | ○ | Dependency | → |
| Resource idle | ······ | Iter. boundary | ╱ |
| Resource in use | —— | | |
| All data at dest. | ● | | |

$$C[i] = A[i] * B[i]$$

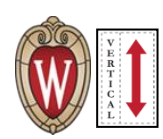**Stream Commands**

Time →

C1) Mem → Scratch

C2) Scratch Wr Barrier

C3) Scratch → Port A

C4) Mem → Port B

C5) Port C → Mem

C6) Mem → Port B

C7) All Barrier

Program Order

**CGRA fabric state**

**Low-power core state**

Processing

Command generation

Scratchpad

A       B

X

C

■

**Legend:**

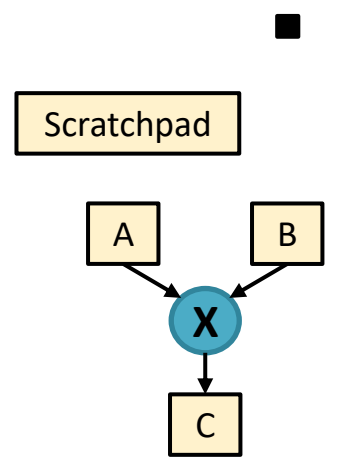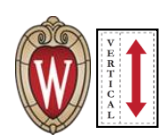| | | | |
|---|---|---|---|
| Enqueued | □ | Barrier | ☉ |
| Dispatched | ○ | Dependency | → |
| Resource idle | ⋯⋯ | Iter. boundary | ╱ |
| Resource in use | — | | |
| All data at dest. | ● | | |

$C[i] = A[i] * B[i]$

**Stream Commands**

Time

C1) Mem → Scratch

C2) Scratch Wr Barrier

C3) Scratch → Port A

C4) Mem → Port B

C5) Port C → Mem

C6) Mem → Port B

C7) All Barrier

**CGRA fabric state**

**Low-power core state**

Program Order

Processing

Command generation

Scratchpad

A    B

X

C

■

# Stream-Dataflow Execution Model Detailed Example

**Legend:**

| | | | |
|---|---|---|---|
| Enqueued | □ | Barrier | ⊙ |
| Dispatched | ○ | Dependency | → |
| Resource idle | ······ | Iter. boundary | ╱ |
| Resource in use | —— | | |
| All data at dest. | ● | | |

$C[i] = A[i] * B[i]$

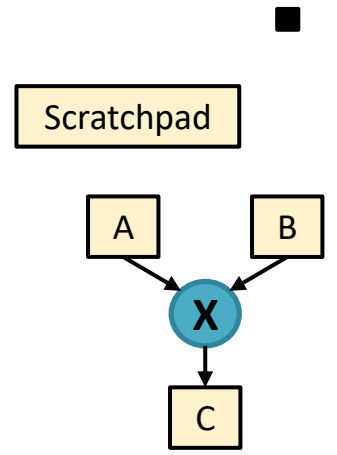**Stream Commands**

C1) Mem → Scratch

C2) Scratch Wr Barrier

C3) Scratch → Port A

C4) Mem → Port B

C5) Port C → Mem

C6) Mem → Port B

C7) All Barrier

**Program Order**

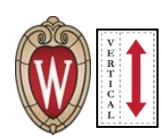**CGRA fabric state**

**Low-power core state**

Time

Processing

Command generation

# Stream-Dataflow Execution Model Detailed Example



**Legend:**

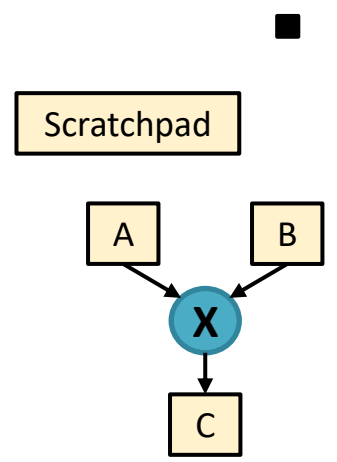| | |
|---|---|
| Enqueued ☐ | Barrier ⊙ |
| Dispatched ○ | Dependency → |
| Resource idle ⋯⋯ | Iter. boundary ╱ |
| Resource in use ── | |
| All data at dest. ● | |

$C[i] = A[i] * B[i]$

**Stream Commands**

Program Order

- C1) Mem → Scratch
- C2) Scratch Wr Barrier
- C3) Scratch → Port A
- C4) Mem → Port B
- C5) Port C → Mem
- C6) Mem → Port B
- C7) All Barrier

**CGRA fabric state**

**Low-power core state**

Time

Processing

Command generation

Scratchpad

A     B

X

C

# Stream-Dataflow Execution Model Detailed Example

**Legend:**

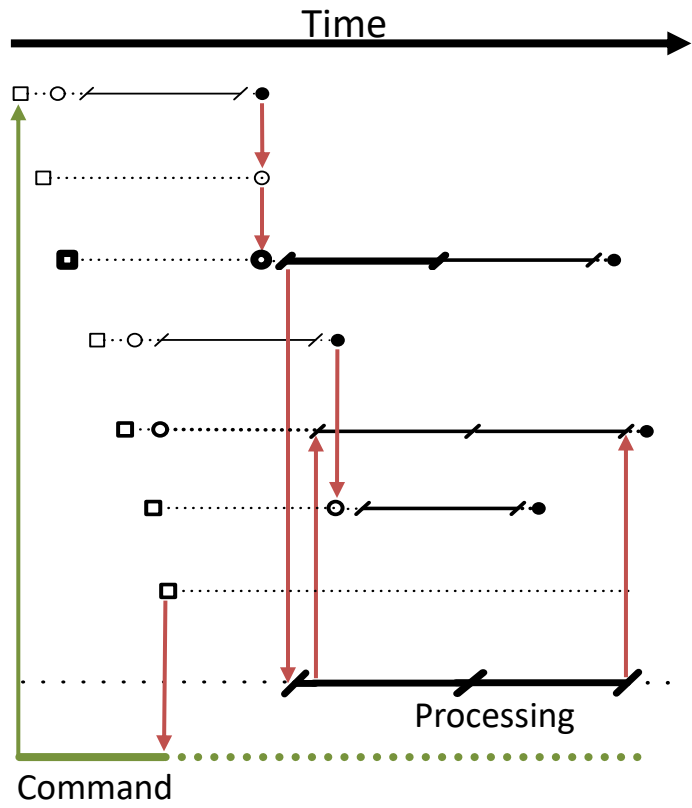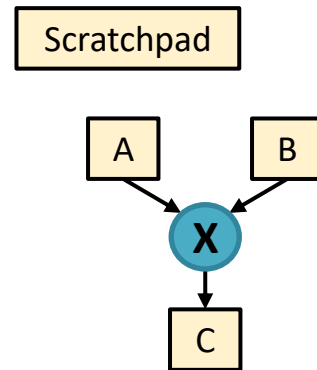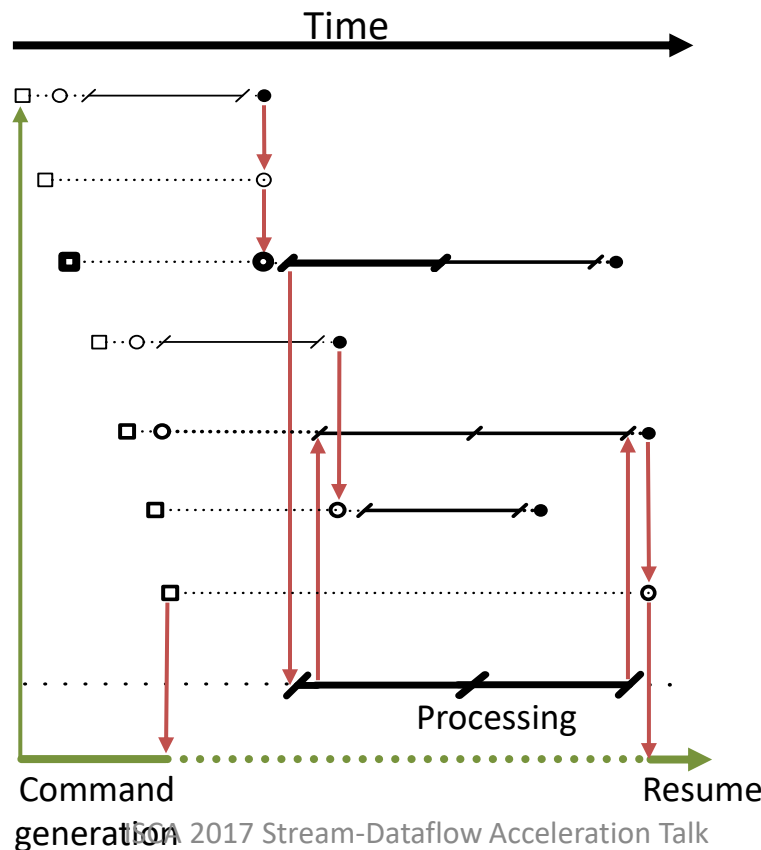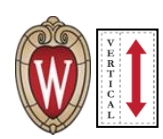| | | |
|---|---|---|
| Enqueued | ☐ | Barrier ⊙ |
| Dispatched | ○ | Dependency → |
| Resource idle | ······· | Iter. boundary ╱ |
| Resource in use | ——— | |
| All data at dest. | ● | |

$C[i] = A[i] * B[i]$

**Stream Commands**

Time →

C1) Mem → Scratch

C2) Scratch Wr Barrier

C3) Scratch → Port A

C4) Mem → Port B

C5) Port C → Mem

C6) Mem → Port B

C7) All Barrier

Program Order

Scratchpad

A    B

X

C

CGRA fabric state

Processing

Low-power core state

Command generation

# Stream-Dataflow Execution Model Detailed Example

**Legend:**

| | | |
|---|---|---|
| Enqueued | □ | Barrier ⊙ |
| Dispatched | ○ | Dependency → |
| Resource idle | ⋯⋯ | Iter. boundary ╱ |
| Resource in use | —— | |
| All data at dest. | ● | |

$$C[i] = A[i] * B[i]$$

**Stream Commands**

Program Order

C1) Mem → Scratch

C2) Scratch Wr Barrier
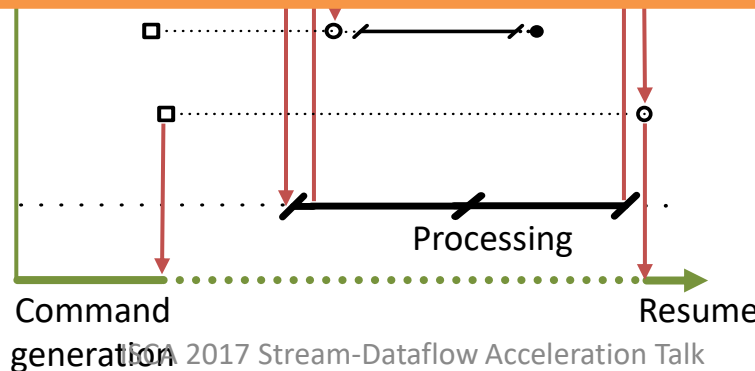
C3) Scratch → Port A

C4) Mem → Port B

C5) Port C → Mem

C6) Mem → Port B

C7) All Barrier

**CGRA fabric state**

**Low-power core state**

Time

Processing

Command generation

# Stream-Dataflow Execution Model Detailed Example



**Legend:**

| | |
|---|---|
| Enqueued □ | Barrier ⊙ |
| Dispatched ○ | Dependency → |
| Resource idle ······· | Iter. boundary ╱ |
| Resource in use ── | |
| All data at dest. ● | |

C[i] = A[i] * B[i]

**Stream Commands**

- C1) Mem → Scratch
- C2) Scratch Wr Barrier
- C3) Scratch → Port A
- C4) Mem → Port B
- C5) Port C → Mem
- C6) Mem → Port B
- C7) All Barrier

**CGRA fabric state**

**Low-power core state**

Time

Program Order

Processing

Command generation

Resume

Scratchpad

A          B

X

C

**Legend:**

| | | | |
|---|---|---|---|
| Enqueued | □ | Barrier | ⊙ |
| Dispatched | ○ | Dependency | → |
| Resource idle | ······· | Iter. boundary | ╱ |

## Stream-Dataflow Accelerator Potential

*1. Dataflow based pipelined concurrent execution*

*2. High Computation Activity Ratio:*
Number of Computations/Stream Commands

Program Order

C6) Mem → Port B
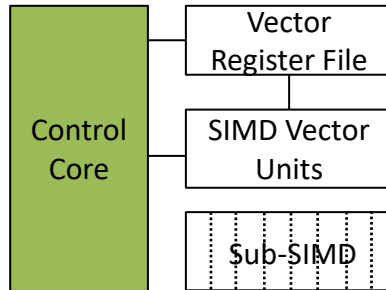
C7) All Barrier

**CGRA fabric state**

**Low-power core state**

Processing

Command generation

Resume

# Inefficiencies in Data-Parallel Architectures



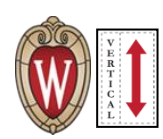|  | **SIMD & Short Vector SIMD** | **SIMT** | **Vector Thread** | **Spatial Dataflow** |
|---|---|---|---|---|
| *Addressing & Communication* | • Unaligned addressing<br>• Complex scatter-gather<br>• Mask & merge instructions | • Redundant address generation<br>• Address coalescing across threads<br>• Non-decoupled access-execute phases | • Redundant address generation | • Redundant address generation<br>• Inefficient memory b/w for local accesses |
| *Resource Utilization & Latency hiding* | • Core-issue width<br>• Fixed vector width<br>• Core to reorder instructions | • Thread scheduling<br>• Multi-ported large register file & cache pressure | • Redundant dispatchers<br>• Core issue width and re-ordering | • Redundant dispatch |
| *Irregular execution support* | • Inefficient general pipeline | • Warp divergence hardware support | • Re-convergence for diverged vector threads | - |

# Stream-Dataflow Accelerator Architecture Opportunities

- Reduce address generation & duplication overheads

- Distributed control to boost pipelined concurrent execution

- High utilization of execution resources w/o massive multi-threading, reducing cache pressure or using multi-ported scratchpad

- Decouple access and execute phases of programs

- Able to be easily customizable/configurable for new application domain
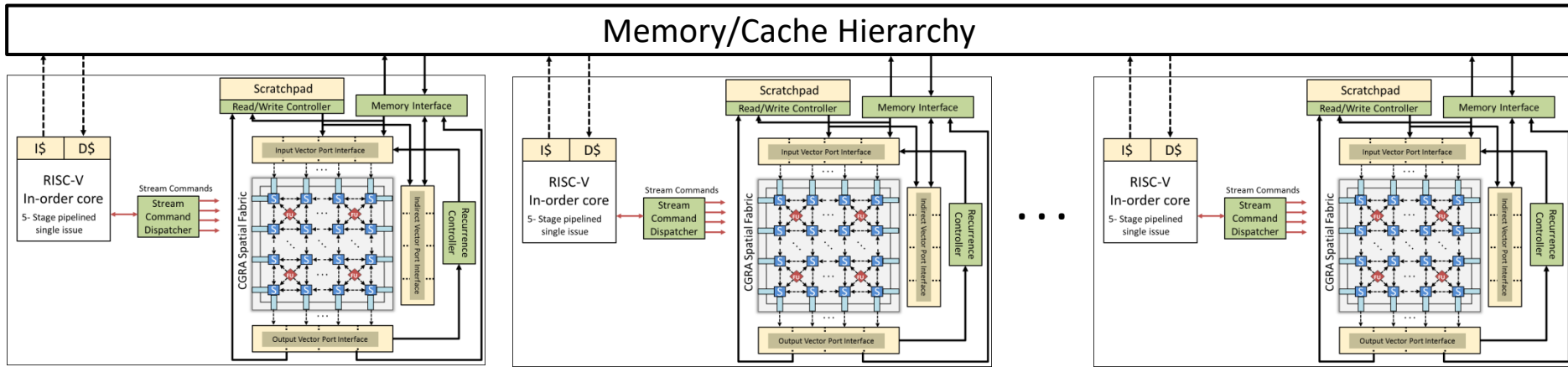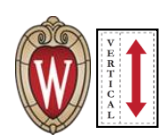
**Stream Dataflow**

| Command Core | | |
|---|---|---|
| | | Scratchpad |
| | Vector Interface | |
| | Coarse-Grained Reconfigurable Arch. | |
| | Vector Interface | |
| | Memory Interface | |

# Stream-Dataflow Accelerator Architecture

**512b** **- - - -** **64b** ——— Stream Command
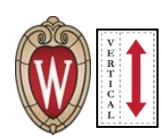
## Multi-Tile Stream-Dataflow Accelerator



- Each tile is connected to higher-L2 cache interface

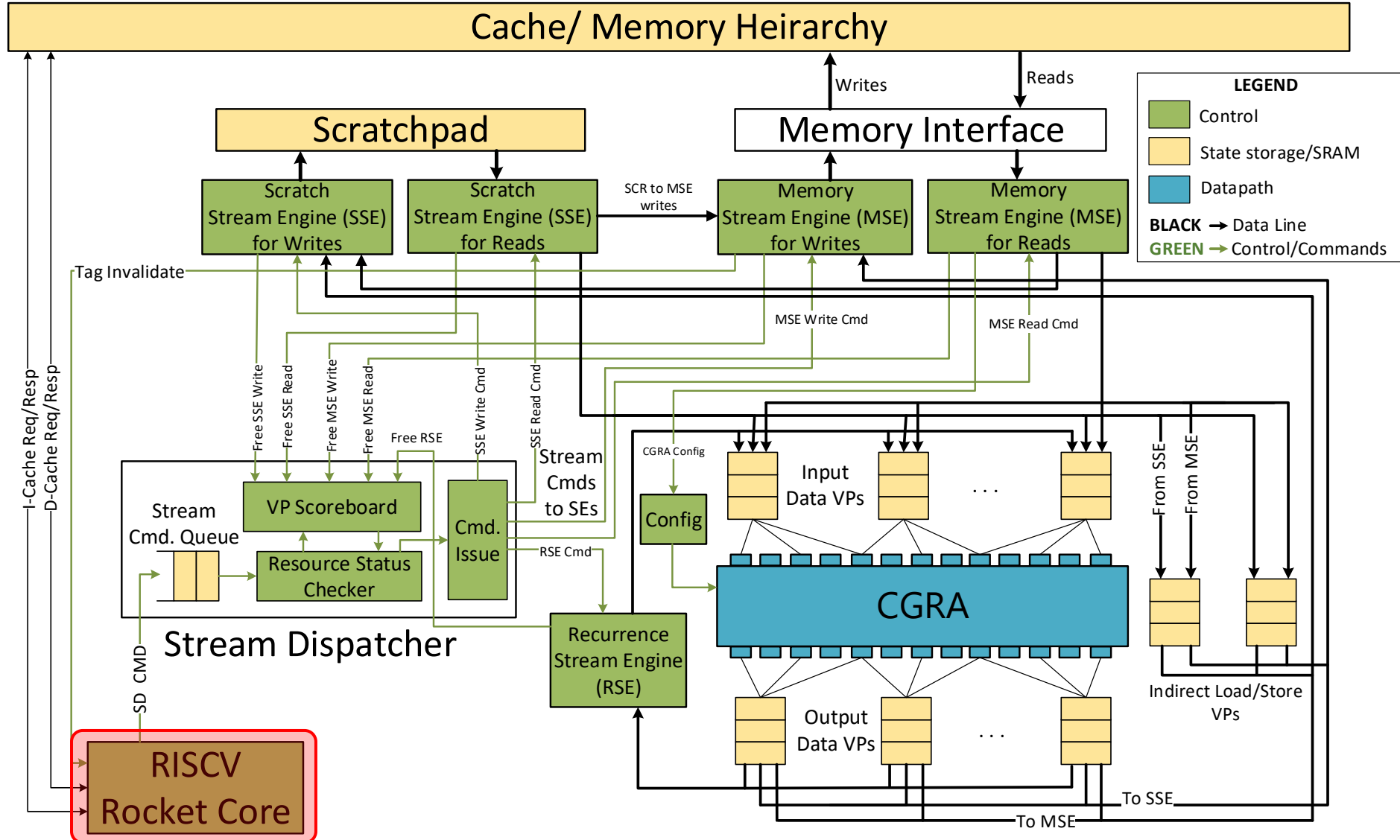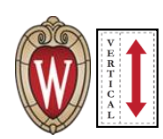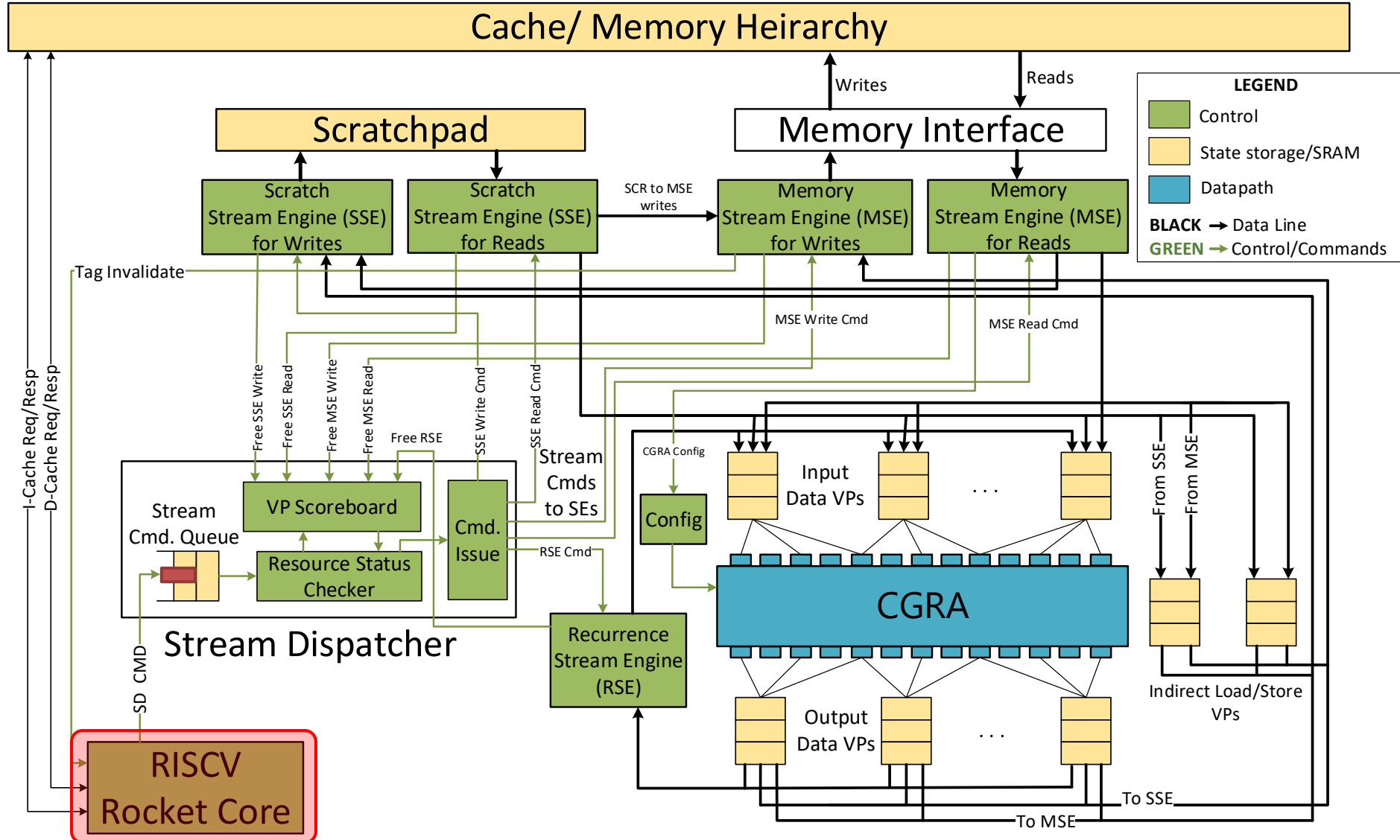- Need a simple scheduler logic to schedule the offloaded stream-dataflow kernels to each tile

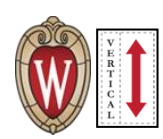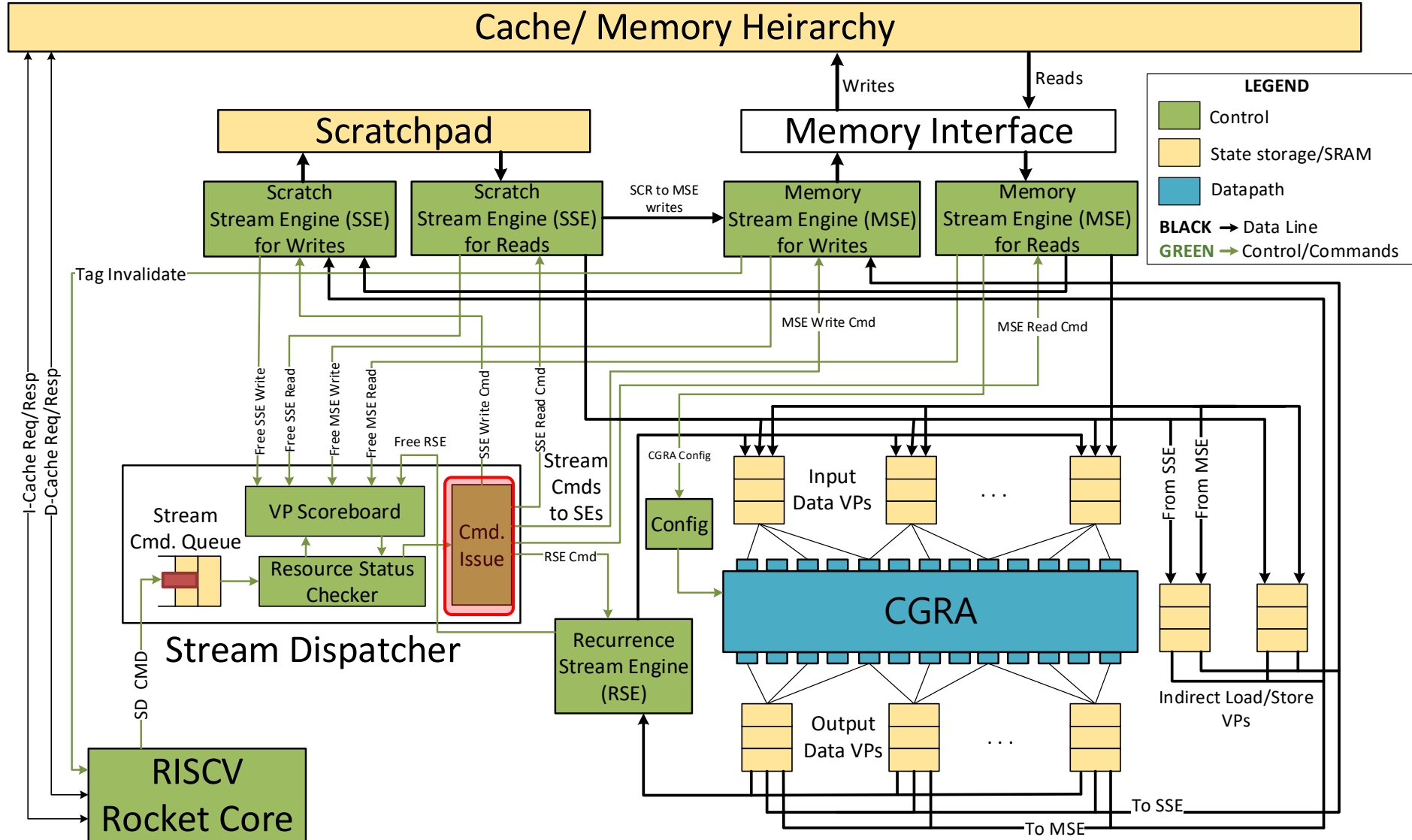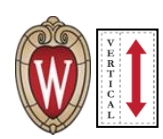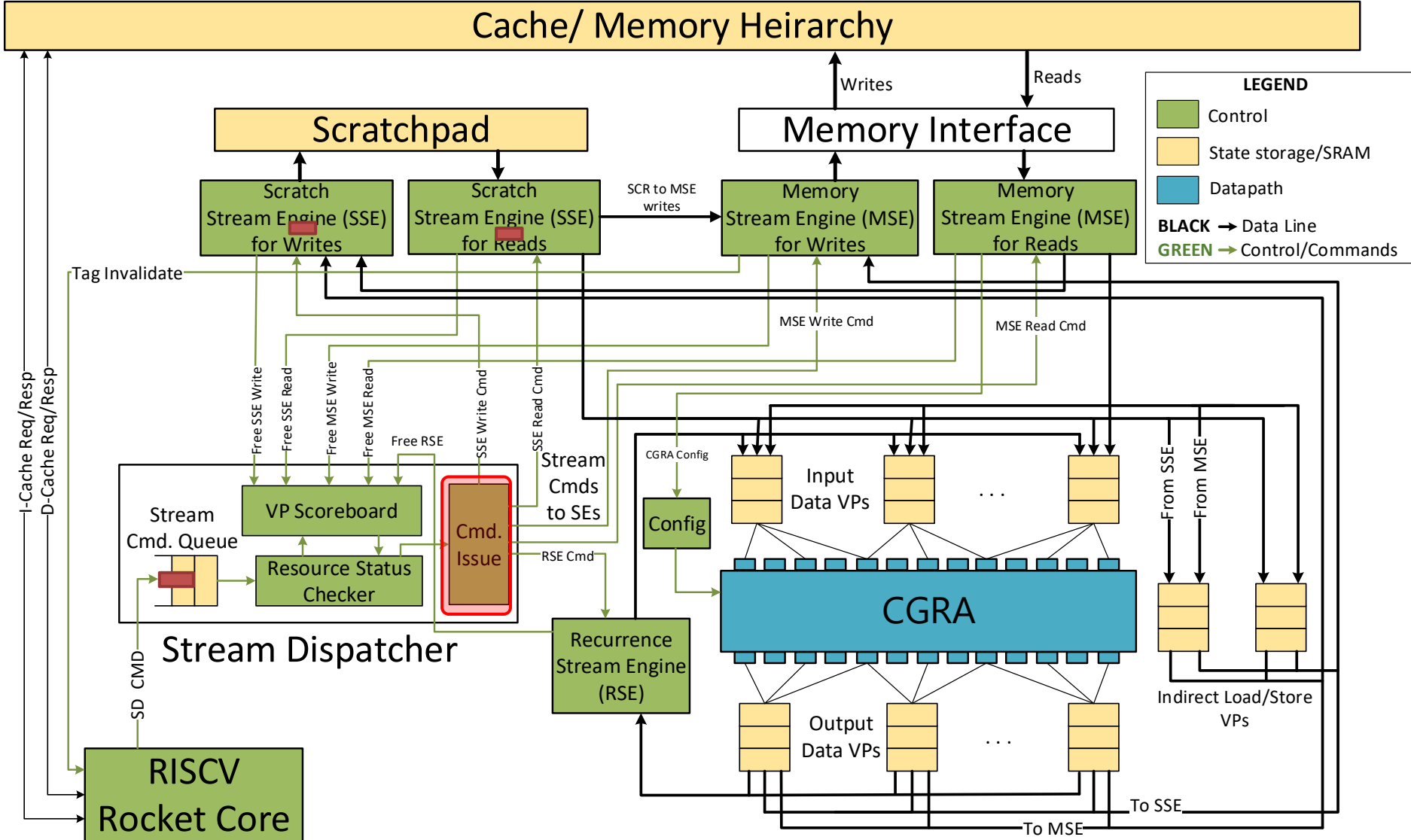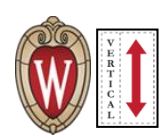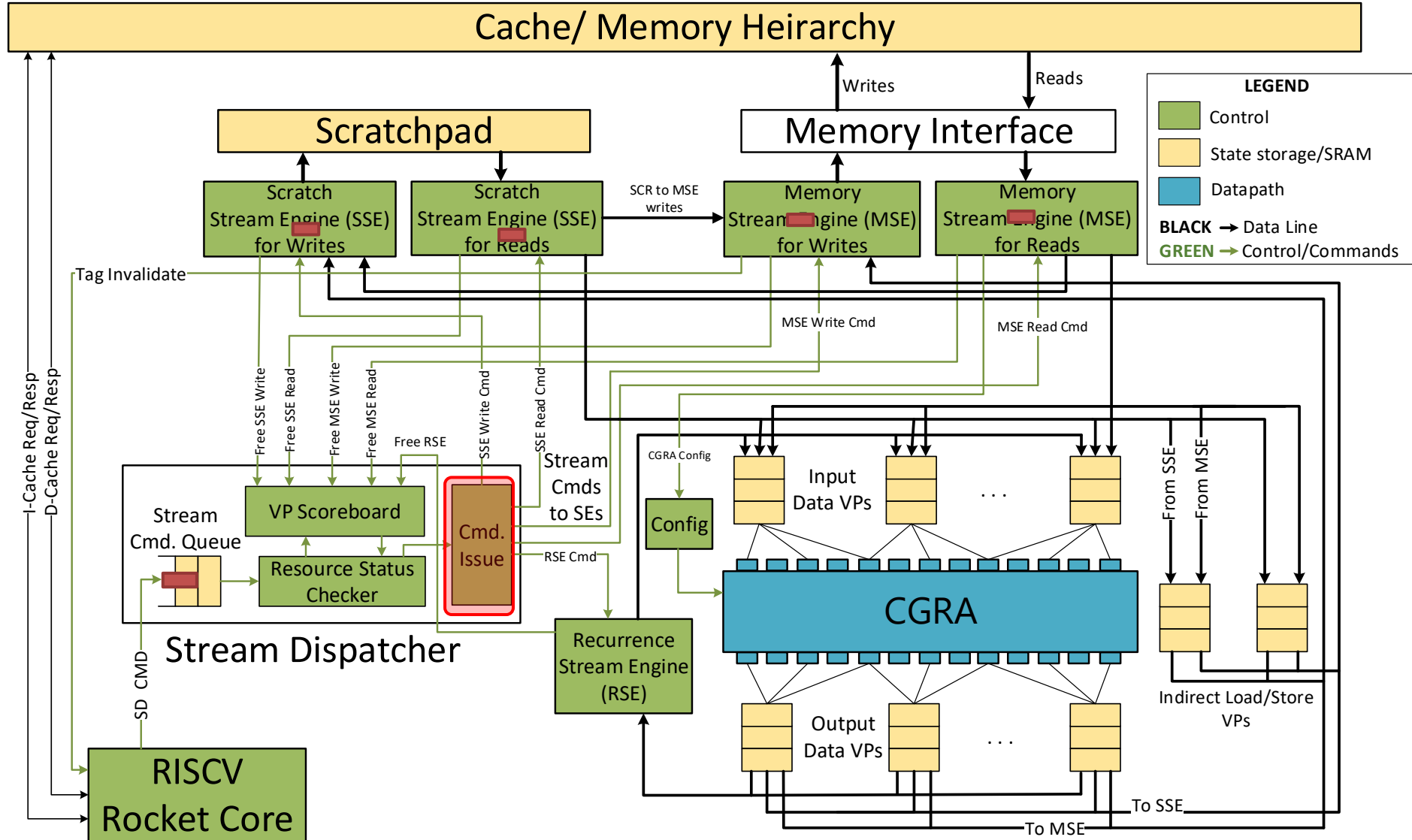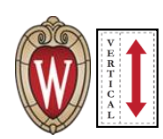# Micro-Architecture of Stream-Dataflow Accelerator (*Softbrain*)

**UCLA**

**Cache/ Memory Heirarchy**

**LEGEND**
- Control (green)
- State storage/SRAM (yellow)
- Datapath (blue)

**BLACK** → Data Line
**GREEN** → Control/Commands

**Scratchpad**

**Memory Interface**

Writes
Reads

**Scratch Stream Engine (SSE) for Writes**

**Scratch Stream Engine (SSE) for Reads**

SCR to MSE writes

**Memory Stream Engine (MSE) for Writes**

**Memory Stream Engine (MSE) for Reads**

Tag Invalidate

MSE Write Cmd
MSE Read Cmd

I-Cache Req/Resp
D-Cache Req/Resp

Free SSE Write
Free SSE Read
Free MSE Write
Free MSE Read
Free RSE
SSE Write Cmd
SSE Read Cmd

Stream Cmds to SEs

CGRA Config

Input Data VPs

...

From SSE
From MSE

**VP Scoreboard**

**Stream Cmd. Queue**

**Cmd. Issue**

**Config**

**Resource Status Checker**

RSE Cmd

**CGRA**

**Stream Dispatcher**

**Recurrence Stream Engine (RSE)**

Indirect Load/Store VPs

SD  CMD

**RISCV Rocket Core**
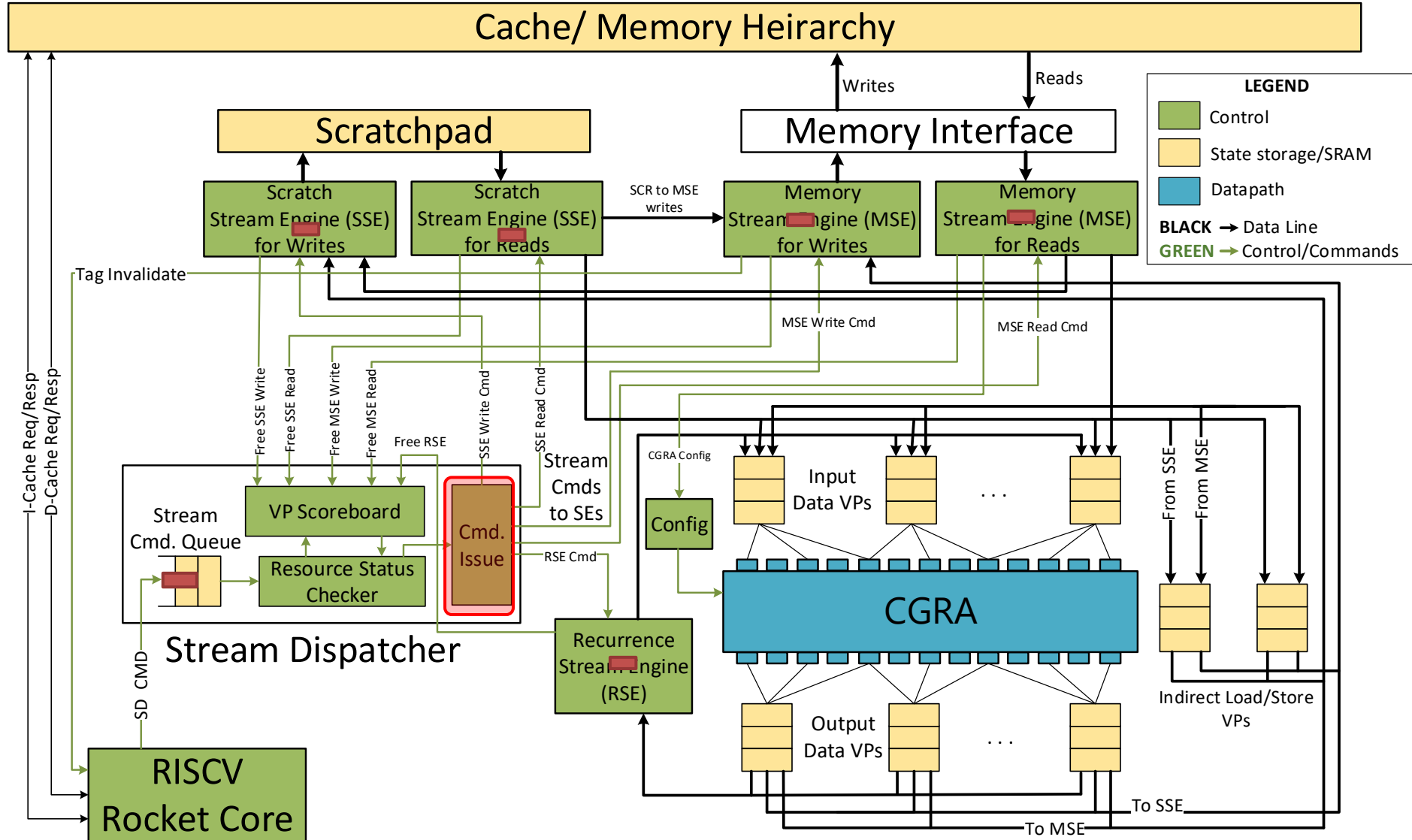
Output Data VPs

...

To SSE
To MSE

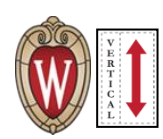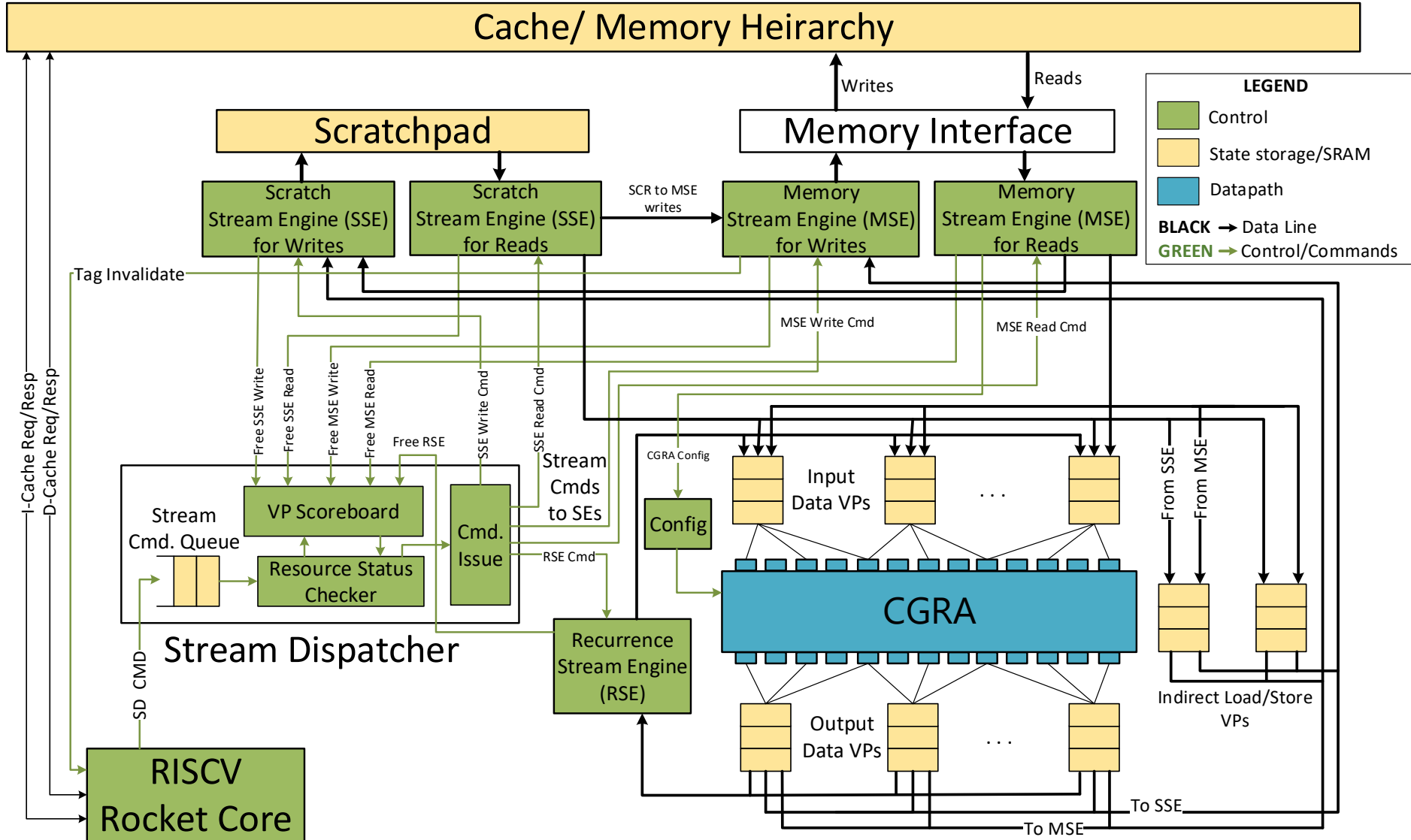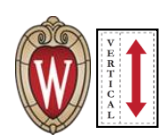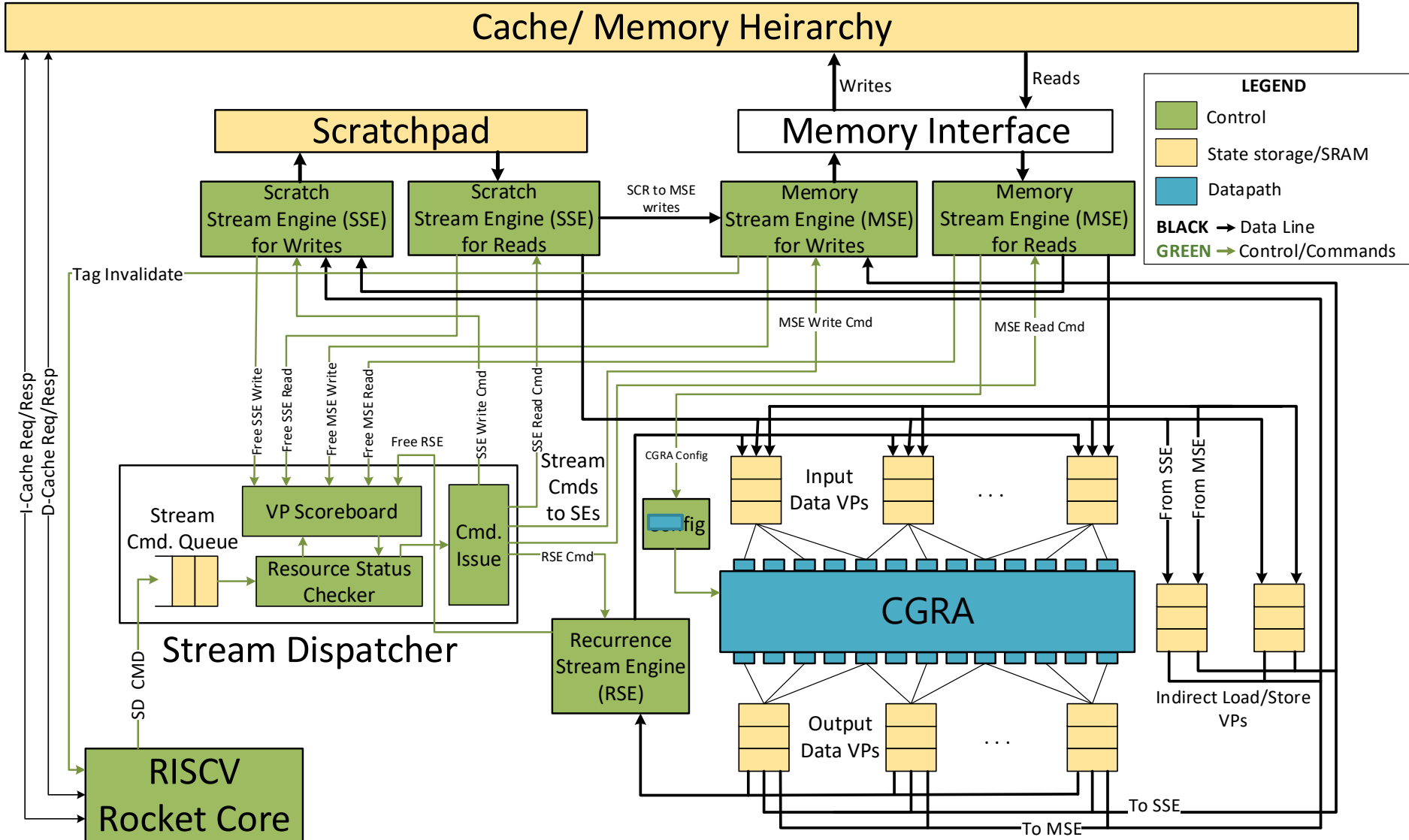# Micro-Architecture of Stream-Dataflow Accelerator (*Softbrain*)

# Micro-Architecture of Stream-Dataflow Accelerator (*Softbrain*)

# Micro-Architecture of Stream-Dataflow Accelerator (*Softbrain*)

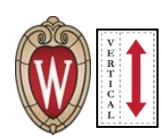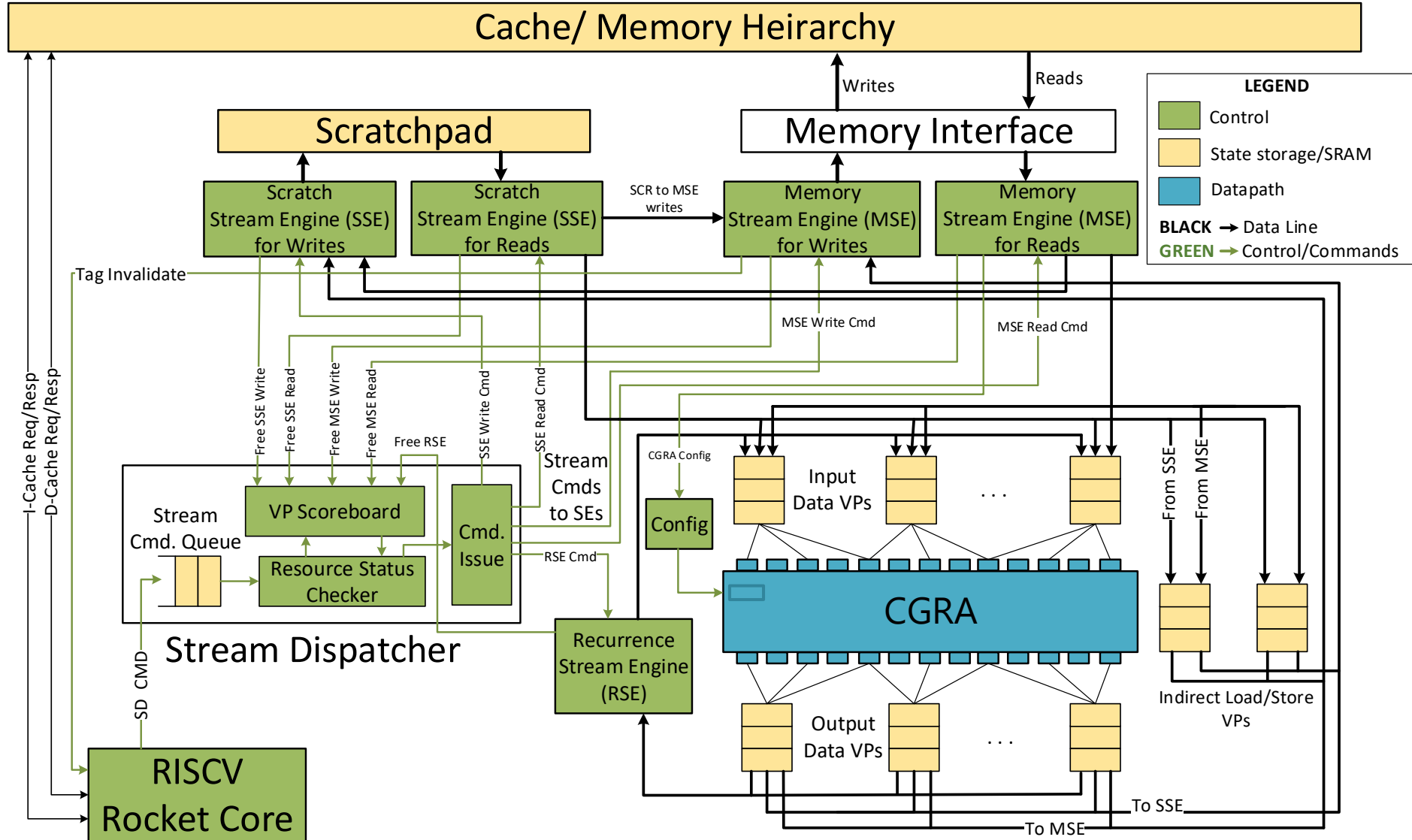# Micro-Architecture of Stream-Dataflow Accelerator (*Softbrain*)

# Micro-Architecture of Stream-Dataflow Accelerator (*Softbrain*)

# Micro-Architecture of Stream-Dataflow Accelerator (*Softbrain*)

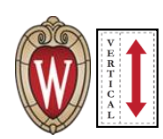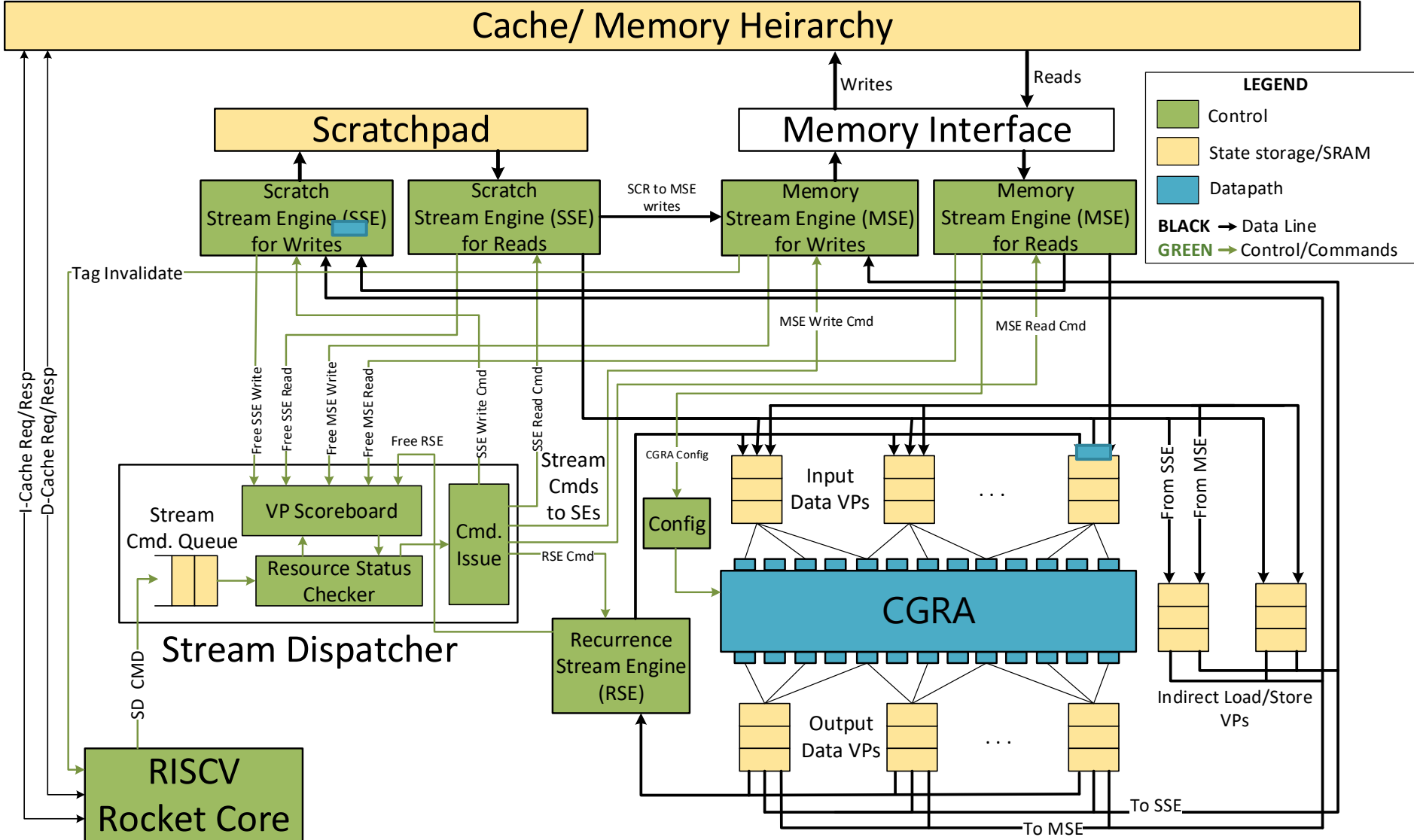# Micro-Architecture of Stream-Dataflow Accelerator (*Softbrain*)

**UCLA**

# Micro-Architecture of Stream-Dataflow Accelerator (*Softbrain*)

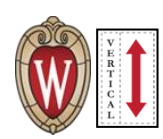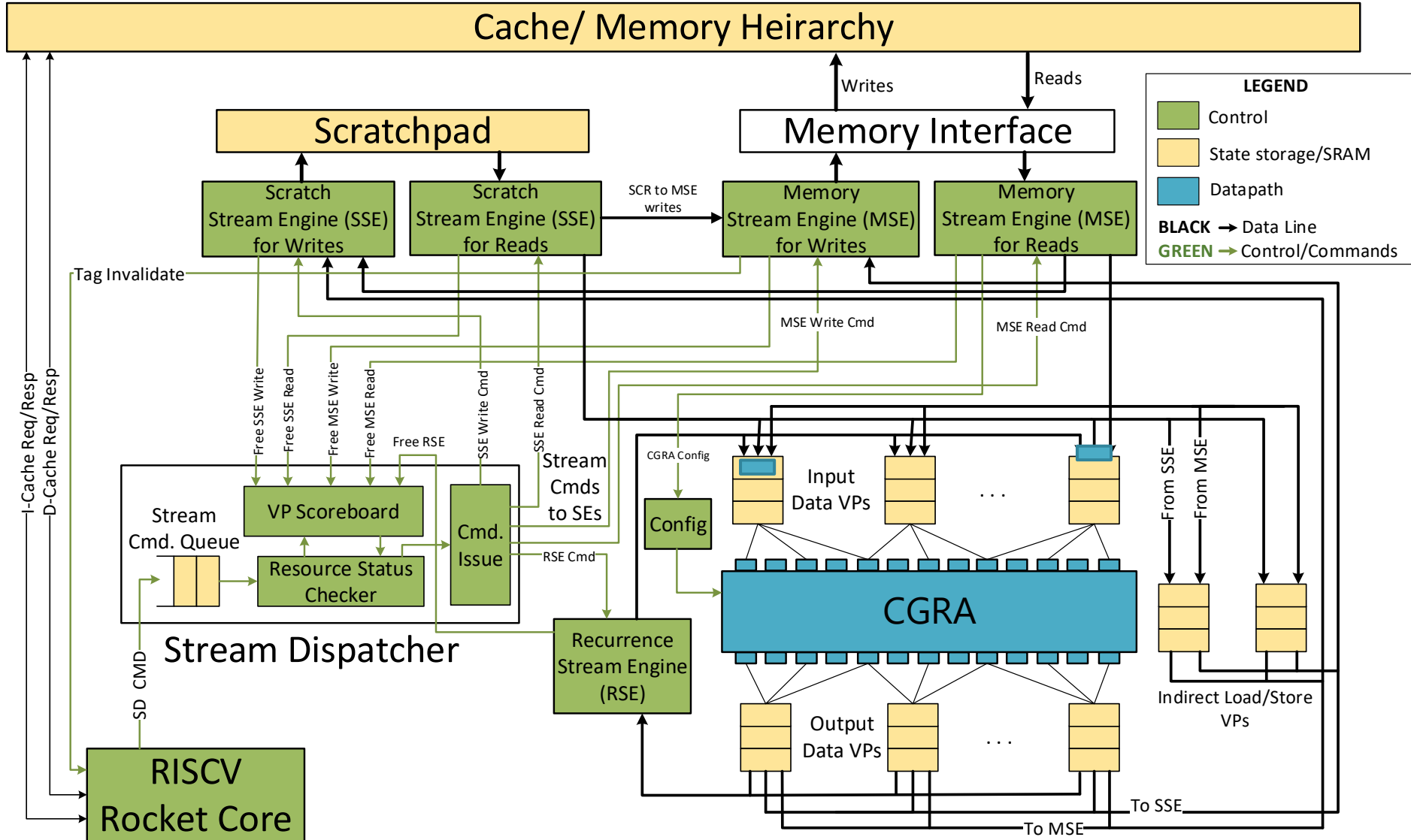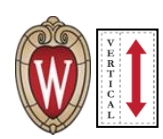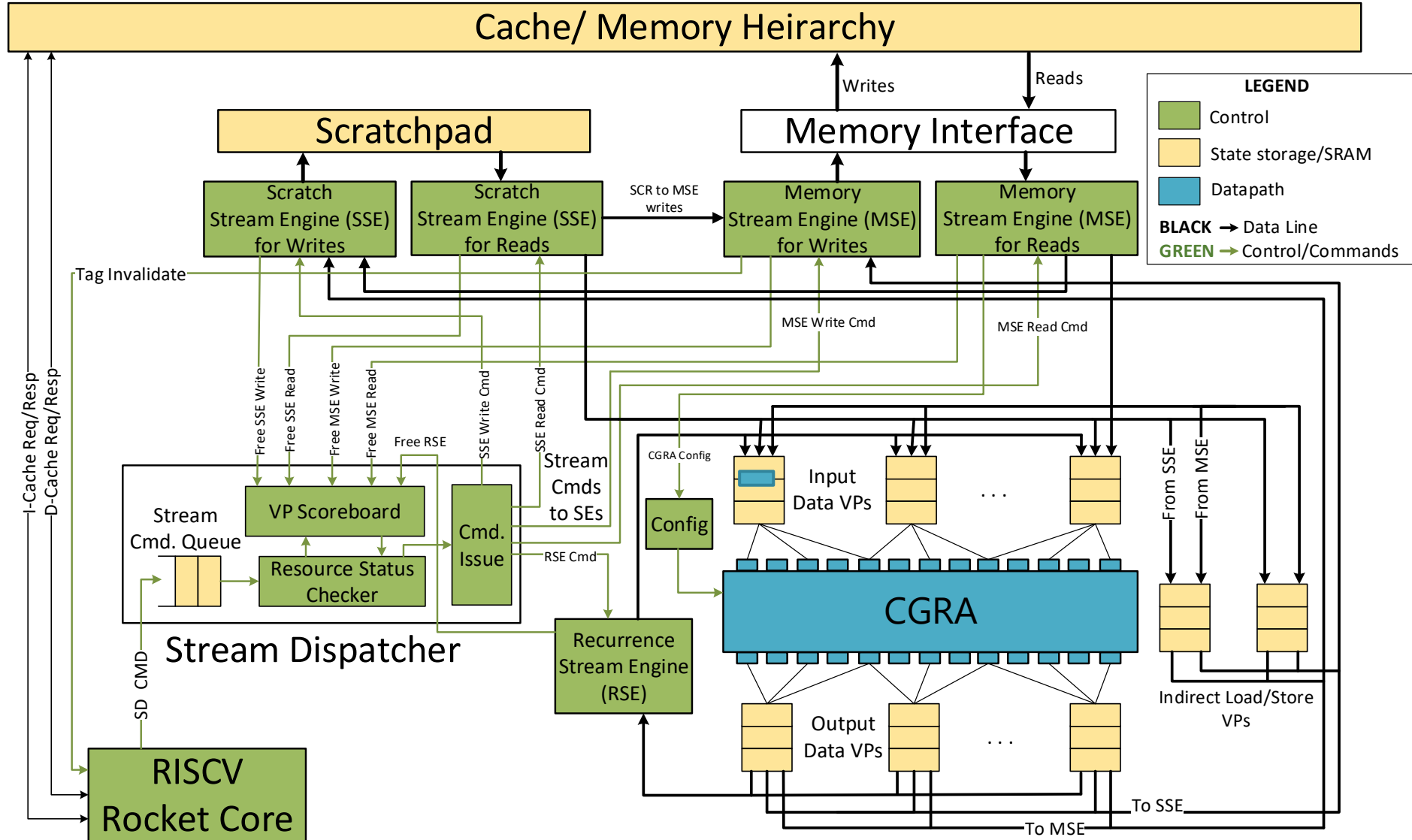# Micro-Architecture of Stream-Dataflow Accelerator (*Softbrain*)

UCLA



**LEGEND**

- Control
- State storage/SRAM
- Datapath

**BLACK** ➡ Data Line
**GREEN** ➡ Control/Commands

Cache/ Memory Heirarchy

Scratchpad

Memory Interface

Writes    Reads

Scratch Stream Engine (SSE) for Writes

Scratch Stream Engine (SSE) for Reads

SCR to MSE writes

Memory Stream Engine (MSE) for Writes

Memory Stream Engine (MSE) for Reads

Tag Invalidate

MSE Write Cmd    MSE Read Cmd

Free SSE Write
Free SSE Read
Free MSE Write
Free MSE Read
Free RSE
SSE Write Cmd
SSE Read Cmd

Stream Cmds to SEs

CGRA Config

Input Data VPs    . . .

From SSE
From MSE

Stream Cmd. Queue

VP Scoreboard

Cmd. Issue

Config

Resource Status Checker

RSE Cmd

CGRA

Stream Dispatcher

SD CMD

Recurrence Stream Engine (RSE)

Indirect Load/Store VPs

I-Cache Req/Resp
D-Cache Req/Resp

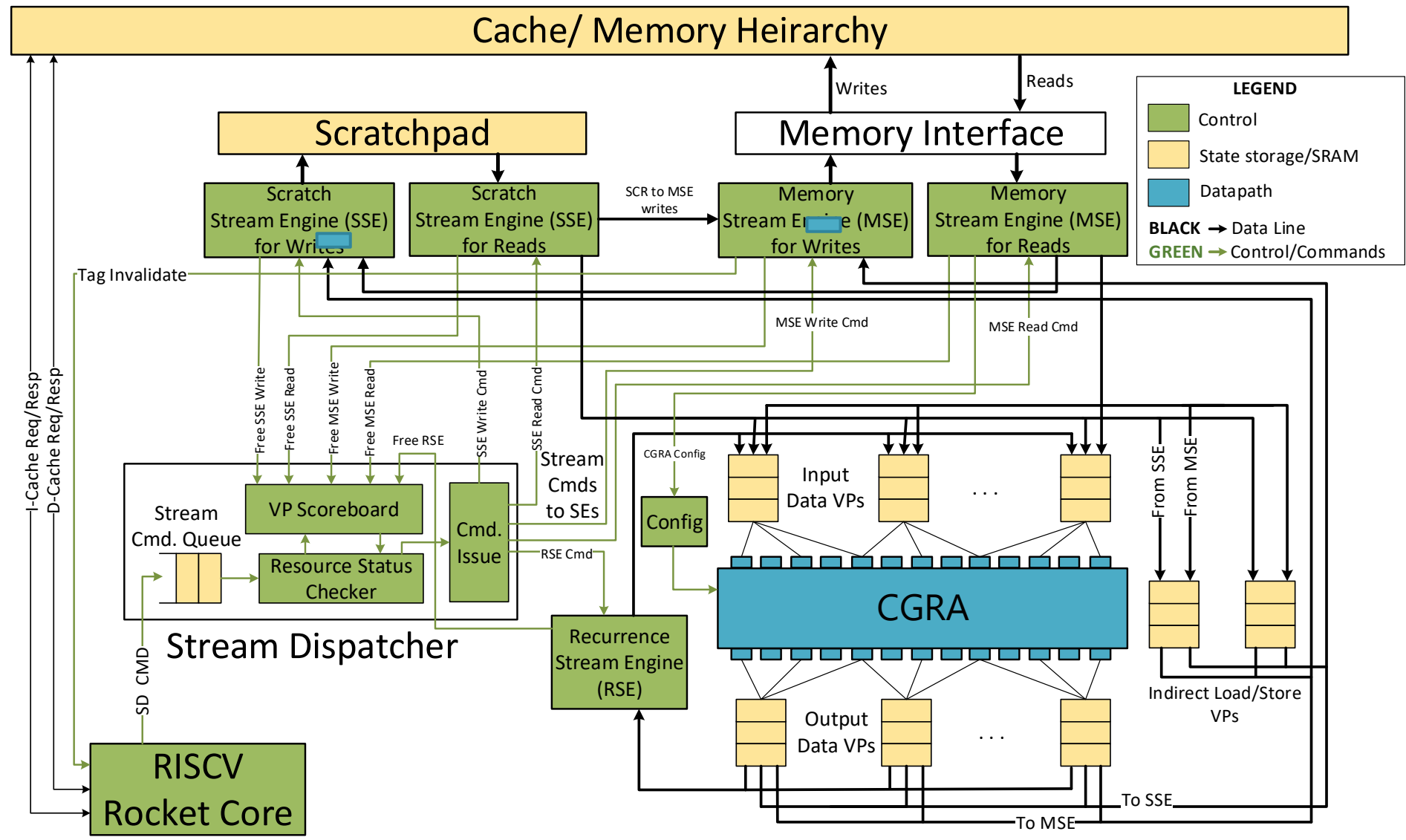RISCV Rocket Core

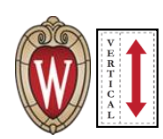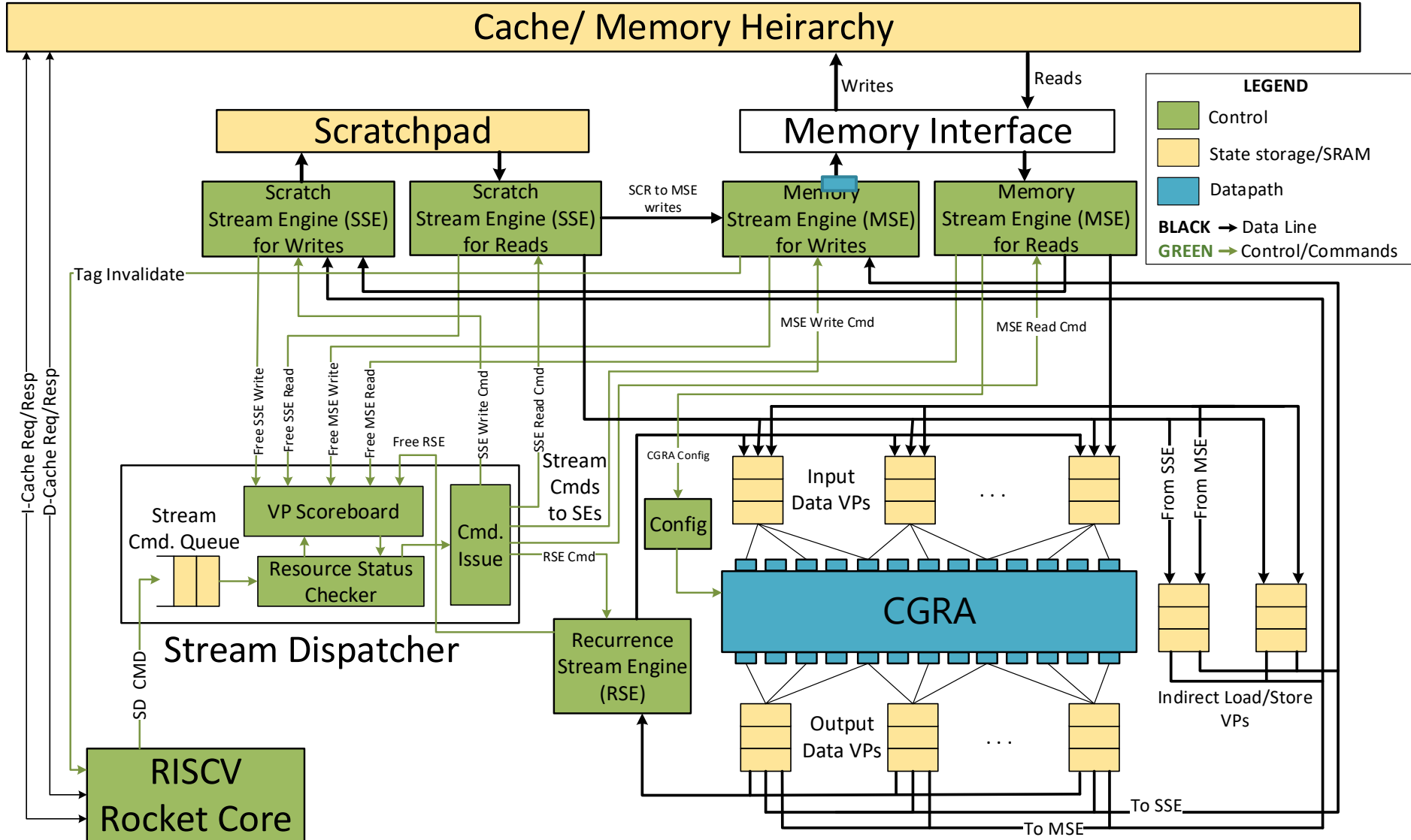Output Data VPs    . . .

To SSE
To MSE

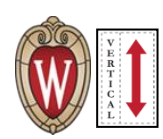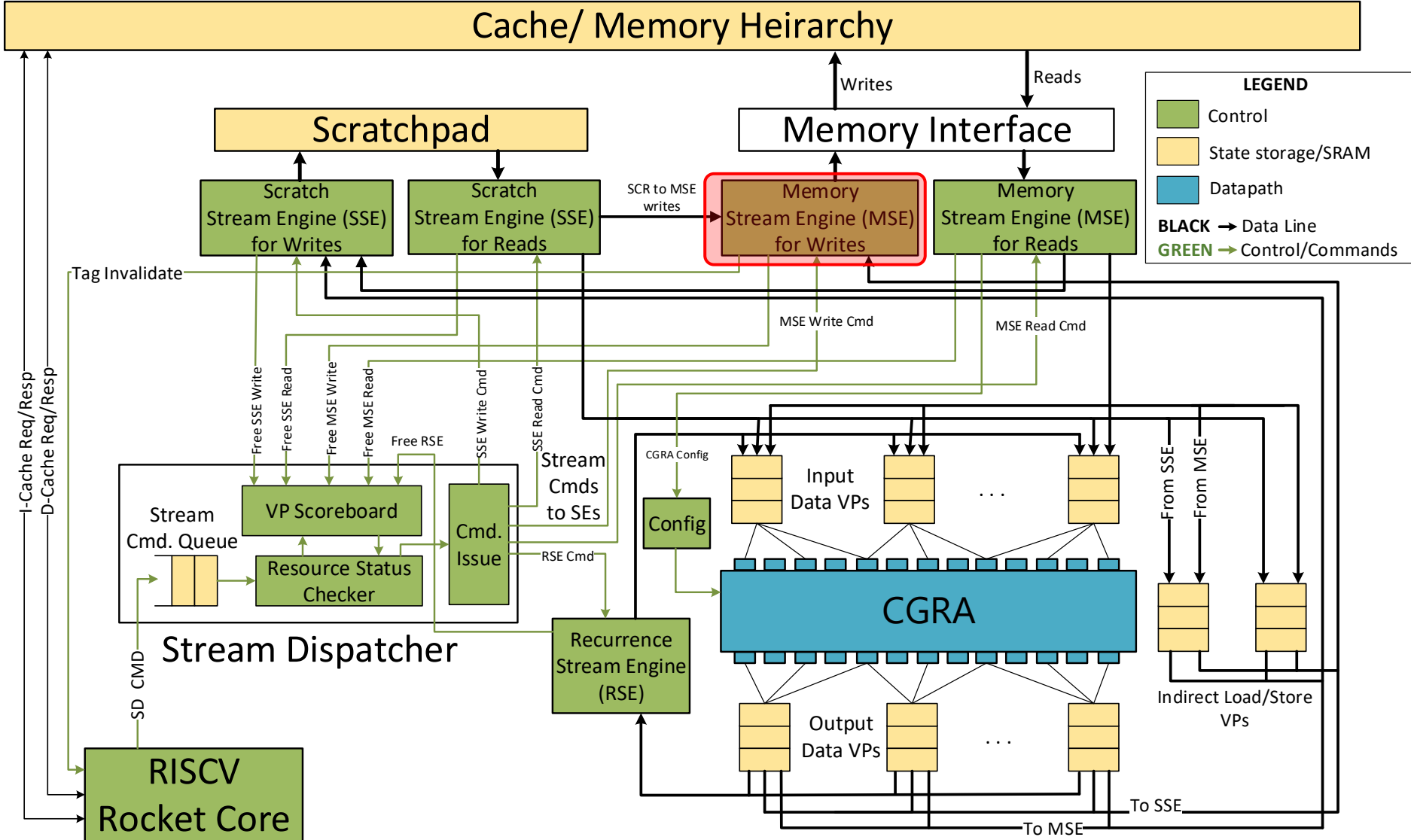# Micro-Architecture of Stream-Dataflow Accelerator (*Softbrain*)

# Micro-Architecture of Stream-Dataflow Accelerator (*Softbrain*)

UCLA



**LEGEND**

- Control
- State storage/SRAM
- Datapath

**BLACK** → Data Line
**GREEN** → Control/Commands
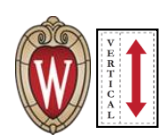
# Micro-Architecture of Stream-Dataflow Accelerator (*Softbrain*)

# Micro-Architecture of Stream-Dataflow Accelerator (*Softbrain*)

**UCLA**

Cache/ Memory Heirarchy

Scratchpad

Memory Interface

Writes

Reads

**LEGEND**
- Control
- State storage/SRAM
- Datapath

**BLACK** → Data Line
**GREEN** → Control/Commands

Scratch Stream Engine (SSE) for Writes

Scratch Stream Engine (SSE) for Reads

SCR to MSE writes

Memory Stream Engine (MSE) for Writes

Memory Stream Engine (MSE) for Reads

Tag Invalidate

MSE Write Cmd

MSE Read Cmd

Free SSE Write
Free SSE Read
Free MSE Write
Free MSE Read
Free RSE
SSE Write Cmd
SSE Read Cmd

Stream Cmds to SEs

CGRA Config

Input Data VPs

From SSE
From MSE

Stream Cmd. Queue

VP Scoreboard

Cmd. Issue

Config

Resource Status Checker

RSE Cmd

**CGRA**

Stream Dispatcher

Recurrence Stream Engine (RSE)

Indirect Load/Store VPs

SD CMD

I-Cache Req/Resp
D-Cache Req/Resp

**RISCV Rocket Core**

Output Data VPs

To SSE
To MSE

# Micro-Architecture of Stream-Dataflow Accelerator (*Softbrain*)

# Micro-Architecture of Stream-Dataflow Accelerator (*Softbrain*)
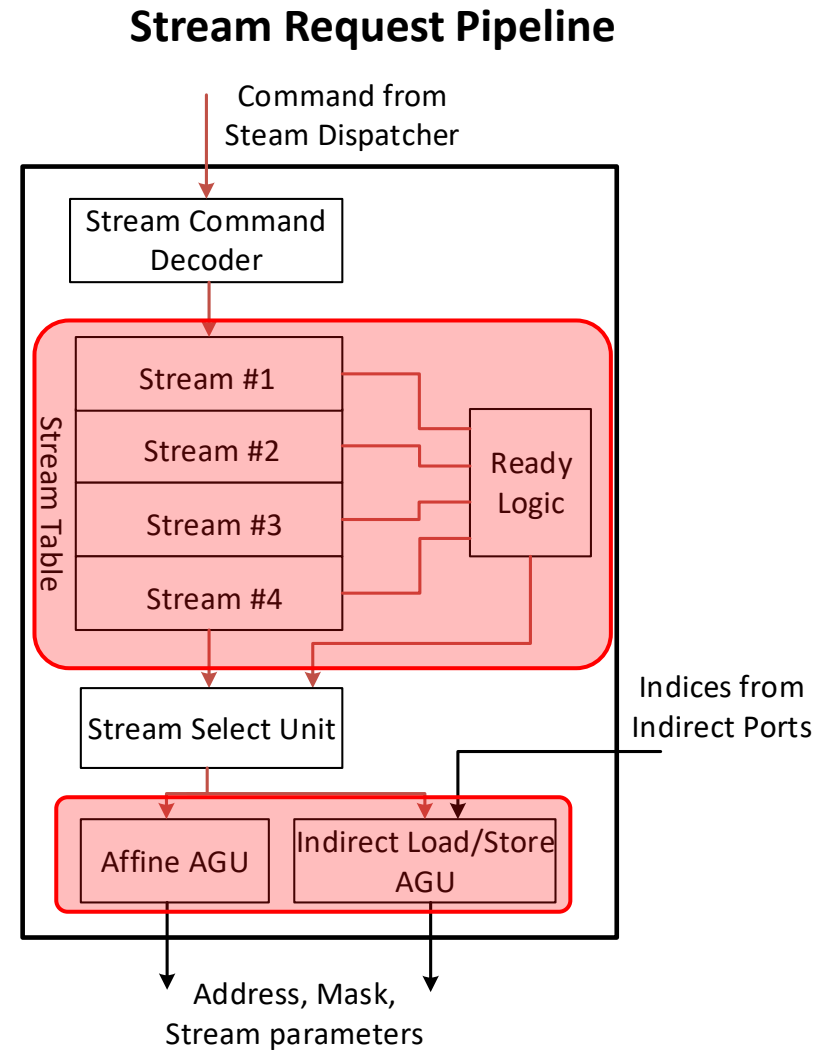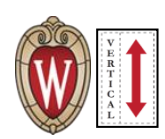
ISCA 2017 Stream-Dataflow Acceleration Talk

# Softbrain Stream Engine Request Pipeline

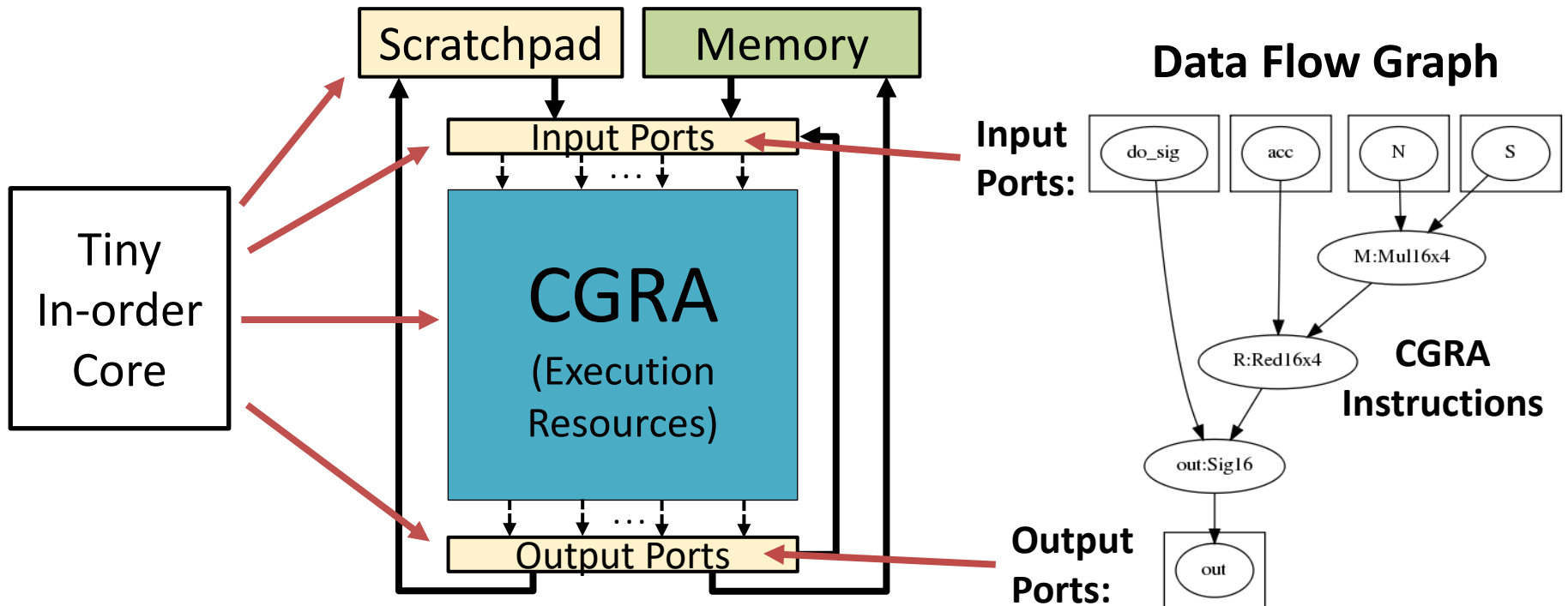- Responsible for address generation for both *affine* and *non-affine* data-streams

- Priority based selection among multiple queued data-steams

- Affine streams – Affine Address Generation Unit (AGU) generates memory addresses

- Non-affine AGU gets addresses and offsets from indirect vector ports

- Similar stream request pipeline is used for scratchpad stream-engines with minimal changes

**Stream Request Pipeline**

# Programming Stream-Dataflow Accelerator

1. Specify Datapath for the CGRA
   - Simple Dataflow Language for DFG

2. Orchestrate the parallel execution of hardware components
   - Coarse-grained stream commands using the stream-interface

# Example Code: Dot Product

**Original Program**

```
for(int i = 0 to N) {
    dot_prod += a[i] * b[i]
}
```

**Computation Graph:**



### Scalar

```
for(i = 0 to N) {
    Send a[i] → P1
    Send b[i] → P2
}
Get P3 -> result
```

~2N Instructions

### Vector

```
for(i = 0 to N, i+=vec_len) {
    Send a[i:i+vec_len] → P1
    Send b[i:i+vec_len] → P2
}
Get P3 -> result
```

~2N/vec_len Instructions

### Stream-Dataflow

```
Send a[i:i+N] → P1
Send b[i:i+N] → P2
Get P3 -> result
```

~3 Instructions

# Existing Architectures for Data Parallel

## Vector Processor

(eg. ARM Neon, X86 SSE)



## Spatial Processor

(eg. Tilera, TRIPS, Wavescalar)



- Amortized Instruction Issue
- Efficient Vector-Memory

- Efficient Dataflow b/t Units
- Flexible Computation Patterns

# Existing Architectures for Data Parallel

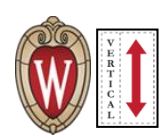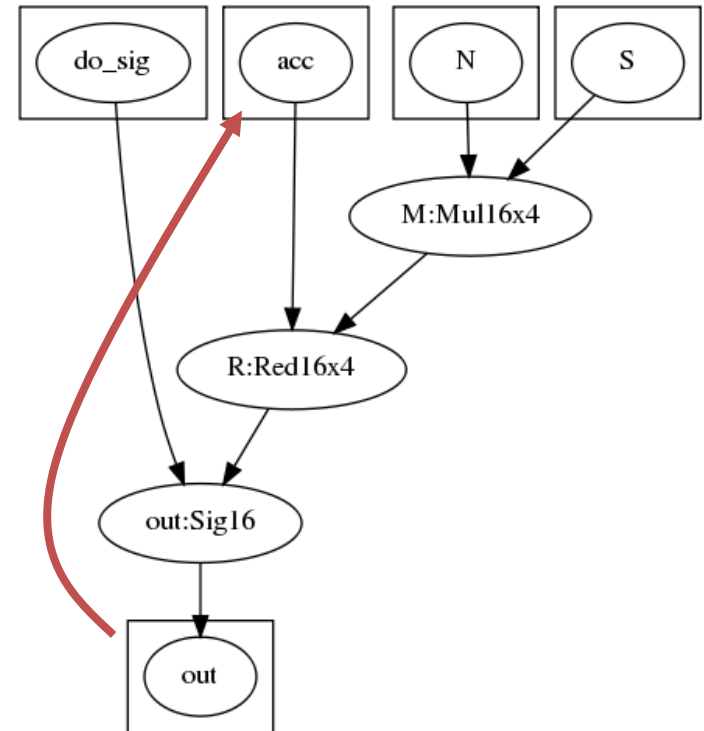## Vector Processor

(eg. ARM Neon, X86 SSE)

Vector Lane

Vector Issue

Reg. — LSU | FUs

Reg. — LSU | FUs

Reg. — LSU | FUs

Vector Mem

## Spatial Processor

(eg. Tilera, TRIPS, Wavescalar)

Independent PEs

Issue — Reg. — LSU | FUs

Issue — Reg. — LSU | FUs

Issue — Reg. — LSU | FUs

Issue — Reg. — LSU | FUs

Scalar Mem

**Vectorized memory interface + Spatial Datapath + Amortized Issue**

# Dataflow Graph (DFG) for CGRA



```
Input: do_sig
Input: acc
Input: N
Input: S
M = Mul16x4(N, S)
R = Red16x4(M, acc)
out = Sig16(R, do_sig)
Output: out
```

# Stream Dataflow Program:

```
uint16_t synapse[Nn][Ni];
uint16_t neuron_i[Ni];
uint16_t neuron_n[Nn];

SD_CONFIG(dfg_config, dfg_size);

SD_DMA_READ(synapse,  8,   8,Ni*Nn/4,P_dfg_S);
SD_DMA_READ(neuron_i, 0,Ni*2,Nn,        P_dfg_N);

for (n = 0; n < Nn/nthreads; n++) {
  SD_CONST(P_dfg_acc,0,1);
  SD_RECURRENCE(P_dfg_out,Ni/4-1,Port_acc);
  SD_CONST(P_dfg_do_sig,0,Ni/4-1);
  SD_CONST(P_dfg_do_sig,1,1);
  SD_DMA_WRITE(P_dfg_out,2,2,1,&neuron_n[n]);
}

SD_WAIT_ALL();
```

# Performance Considerations

- Goal: Fully Pipeline the Largest Data Flow Graph!

- Primary Bottlenecks:

  ☐ Size of Data Flow Graph

    Increase through Loop Unrolling/Stripmining

  ☐ General Core (for Issuing Streams)

    Increase "length" of streams

  ☐ Memory/Cache Bandwidth

    Use Scratchpad for reused Data

  ☐ Recurrence Serialization Overhead

    Either: 1. Increase Parallel Computations (tiling)
            2. Use internal accumulation

# Optimized DFG

```
InputVec: N [0, 1, 2, 3, 4, 5, 6, 7]
InputVec: S [0, 1, 2, 3, 4, 5, 6, 7]
Input: reset

M0 =Mul16x4(N0, S0)
M1 =Mul16x4(N1, S1)
M2 =Mul16x4(N2, S2)
M3 =Mul16x4(N3, S3)
M4 =Mul16x4(N4, S4)
M5 =Mul16x4(N5, S5)
M6 =Mul16x4(N6, S6)
M7 =Mul16x4(N7, S7)

A0 =Add16x4(M0, M1)
A1 =Add16x4(M2, M3)
A2 =Add16x4(M4, M5)
A3 =Add16x4(M6, M7)

A8 =Add16x4(A0, A1)
A9 =Add16x4(A2, A3)

A10 = Add16x4(A8, A9)

Red = Red16x4(A10)

Res = Acc16x4(Red, reset)

out=Sig16(Res)

Output: out
```



**Two optimizations:**
1. **Increased the size of the DFG**
2. **Add an accumulation step and remove recurrence accumulation.**

# Optimized Classifier Layer



Input Neurons (Ni)

×

∑

Output Neurons (Nn)

Synapses (Nn x Ni)

# Optimized Classifier Layer

```
SD_CONFIG(dfg_config,dfg_size);

SD_DMA_READ(synapse,8,8,Ni*Nn/4,P_dfg_S);
SD_DMA_SCRATCH_LOAD(neuron_i, 0,Ni*2,1,0);
SD_WAIT_SCR_WR();

SD_SCR_PORT_STREAM(0, 0, Ni*2,1, P_dfg_N);
for (n = 0; n < Nn/nthreads; n++) {
  SD_CONST(P_dfg_reset,0,Ni/4-1);
  SD_CONST(P_dfg_reset,1,1);
  SD_GARBAGE(P_dfg_out,Ni/4-1);
  SD_DMA_WRITE(P_dfg_out,2,2,1,&neuron_n[n]);
}

SD_WAIT_ALL();
```

Input Neurons (Ni)

×

Σ

Output Neurons (Nn)

Synapses (Nn x Ni)

# DianNao Power/Area Comparison

| | | area(mm$^2$) | power (mw) |
|---|---|---|---|
| Control Core + 16kB I & D$ | | 0.16 | 39.1 |
| CGRA | Network | 0.12 | 31.2 |
| | FUs (4×5) | 0.04 | 24.4 |
| | Total CGRA | 0.16 | 55.6 |
| 5×Stream Engines | | 0.02 | 18.3 |
| Scratchpad (4KB) | | 0.1 | 2.6 |
| Vector Ports (Input & Output) | | 0.03 | 3.6 |
| **1 Softbrain Total** | | **0.47** | 119.3 |
| 8 Softbrain Units | | 3.76 | 954.4 |
| DianNao | | 2.16 | 418.3 |
| Softbrain / DianNao Overhead | | 1.74 | 2.28 |

**Table 3: Area and Power Breakdown / Comparison**
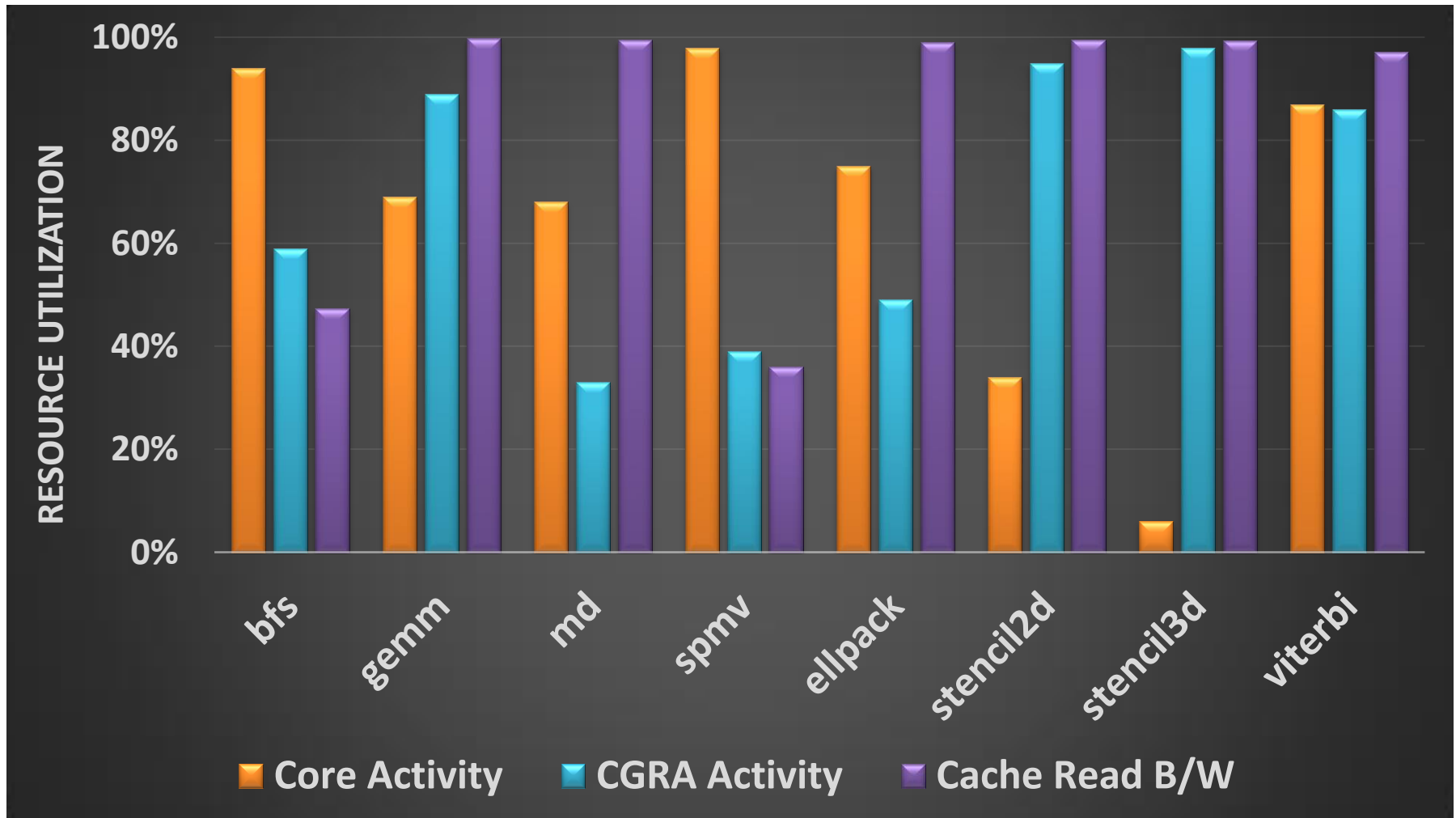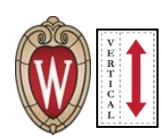**(All numbers normalized to 55nm process technology)**

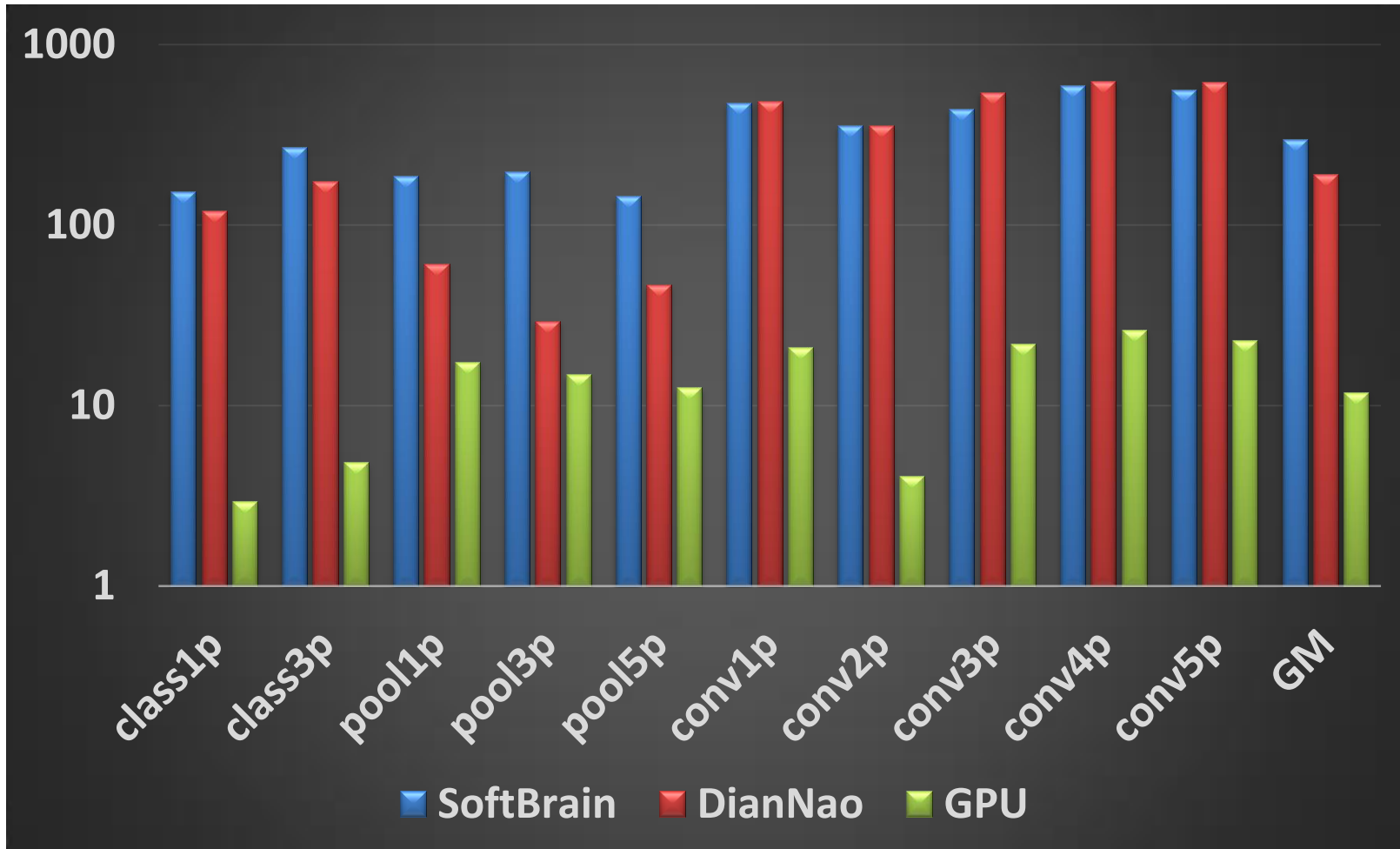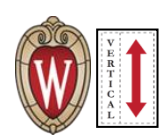# Softbrain Resource Utilization

# Softbrain Resource Utilization

# Softbrain Resource Utilization

# Softbrain Resource Utilization
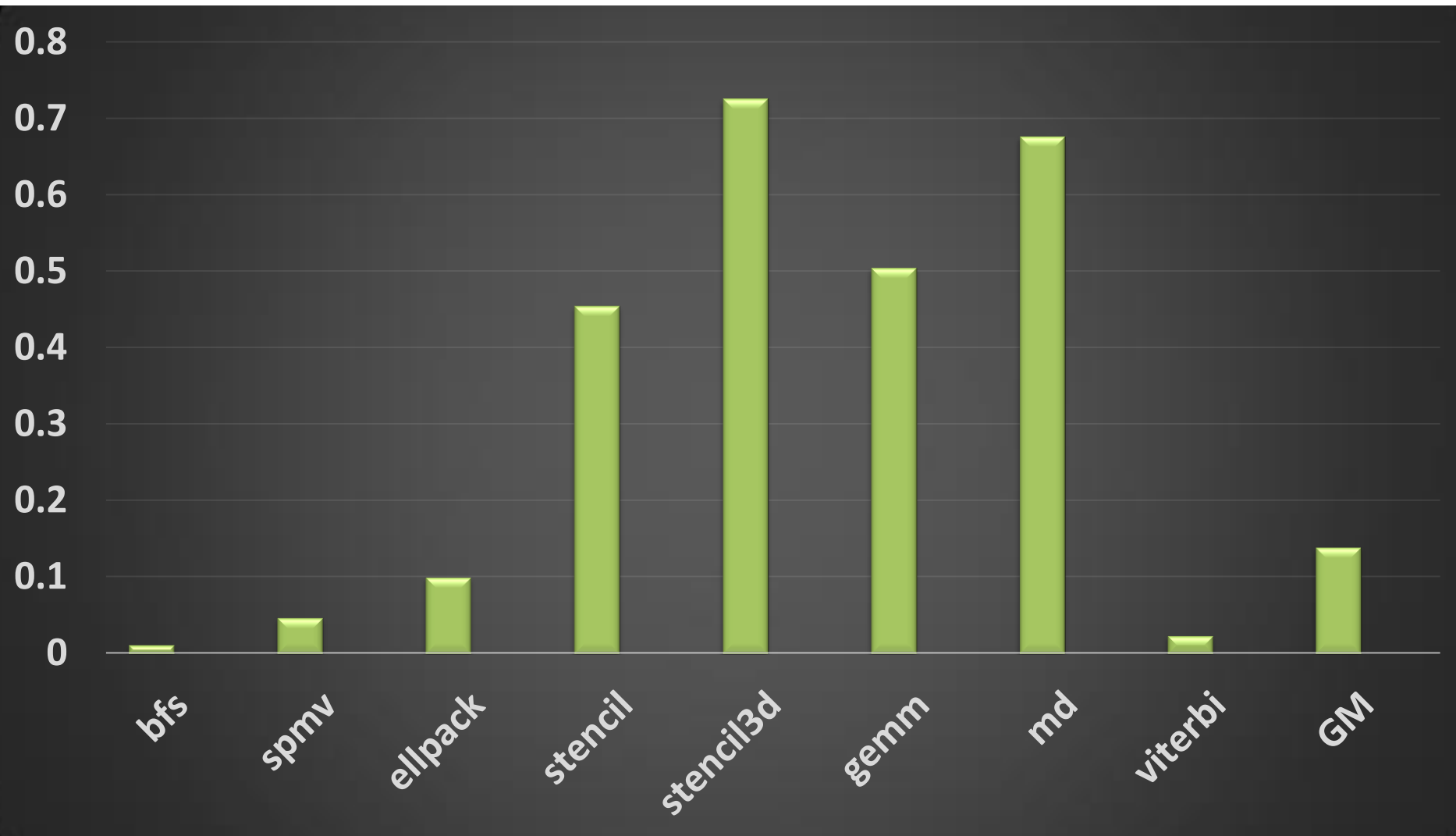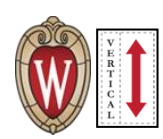
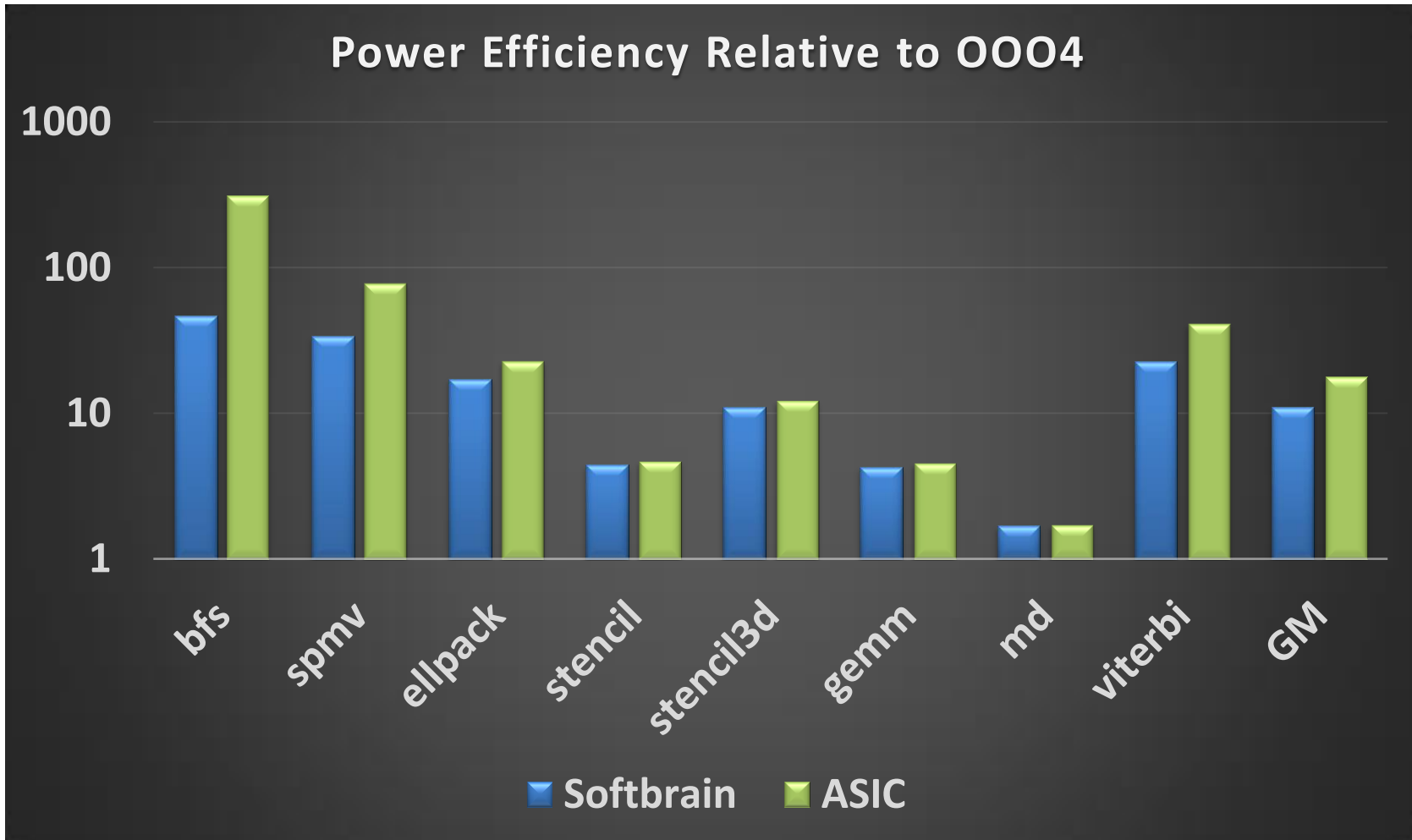# Softbrain Resource Utilization
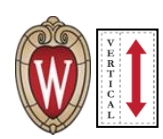
# Softbrain vs. DianNao vs. GPU

# ASIC Area Relative to Softbrain

# Softbrain vs. ASIC
# Power Efficiency Comparison

# Softbrain vs. ASIC
# Energy Efficiency Comparison