# DSAGEN: Synthesizing Programmable Spatial Accelerators

Jian Weng, Sihao Liu, Vidushi Dadu, Zhengrong Wang, Preyas Shah, Tony Nowatzki
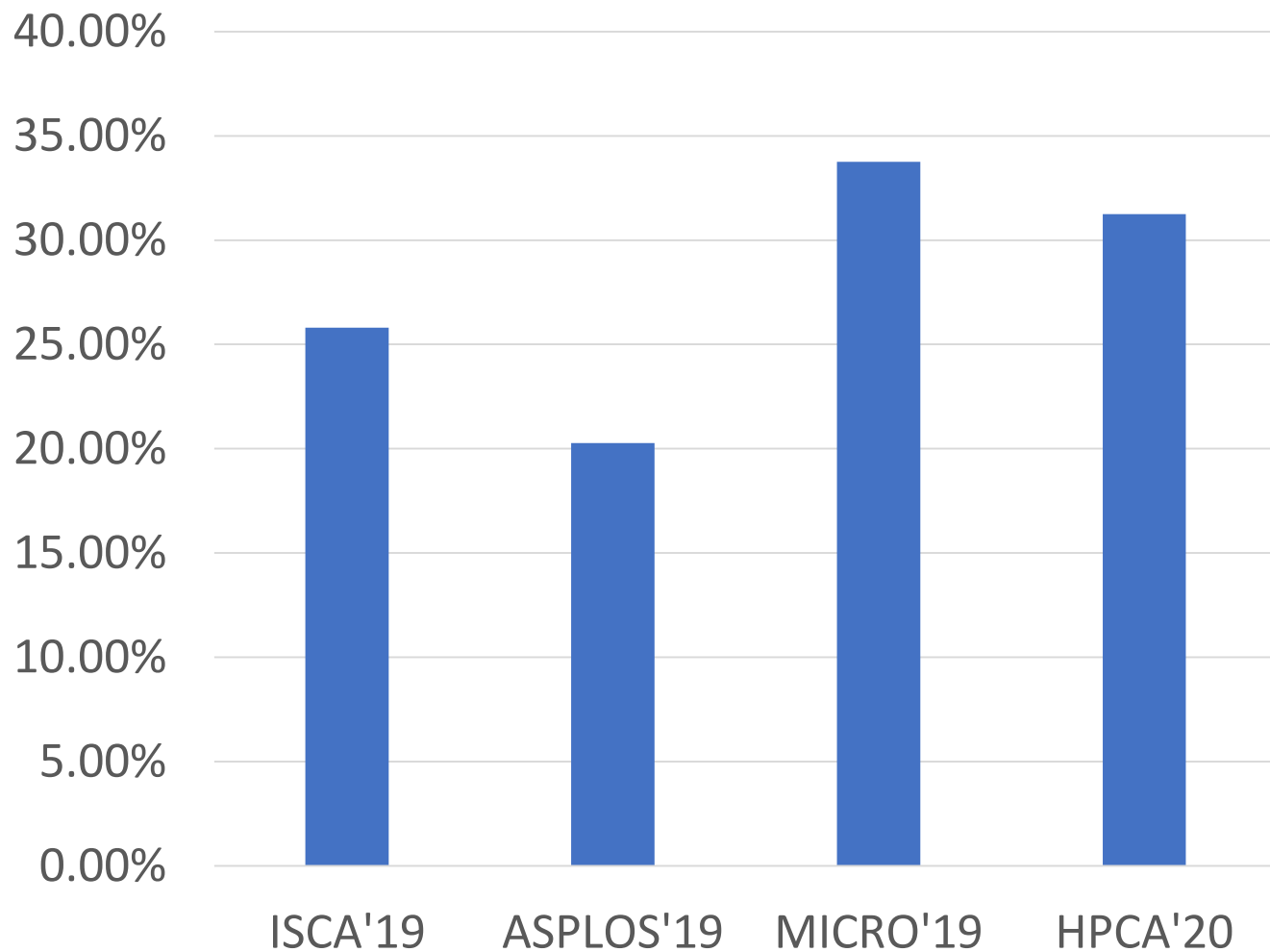
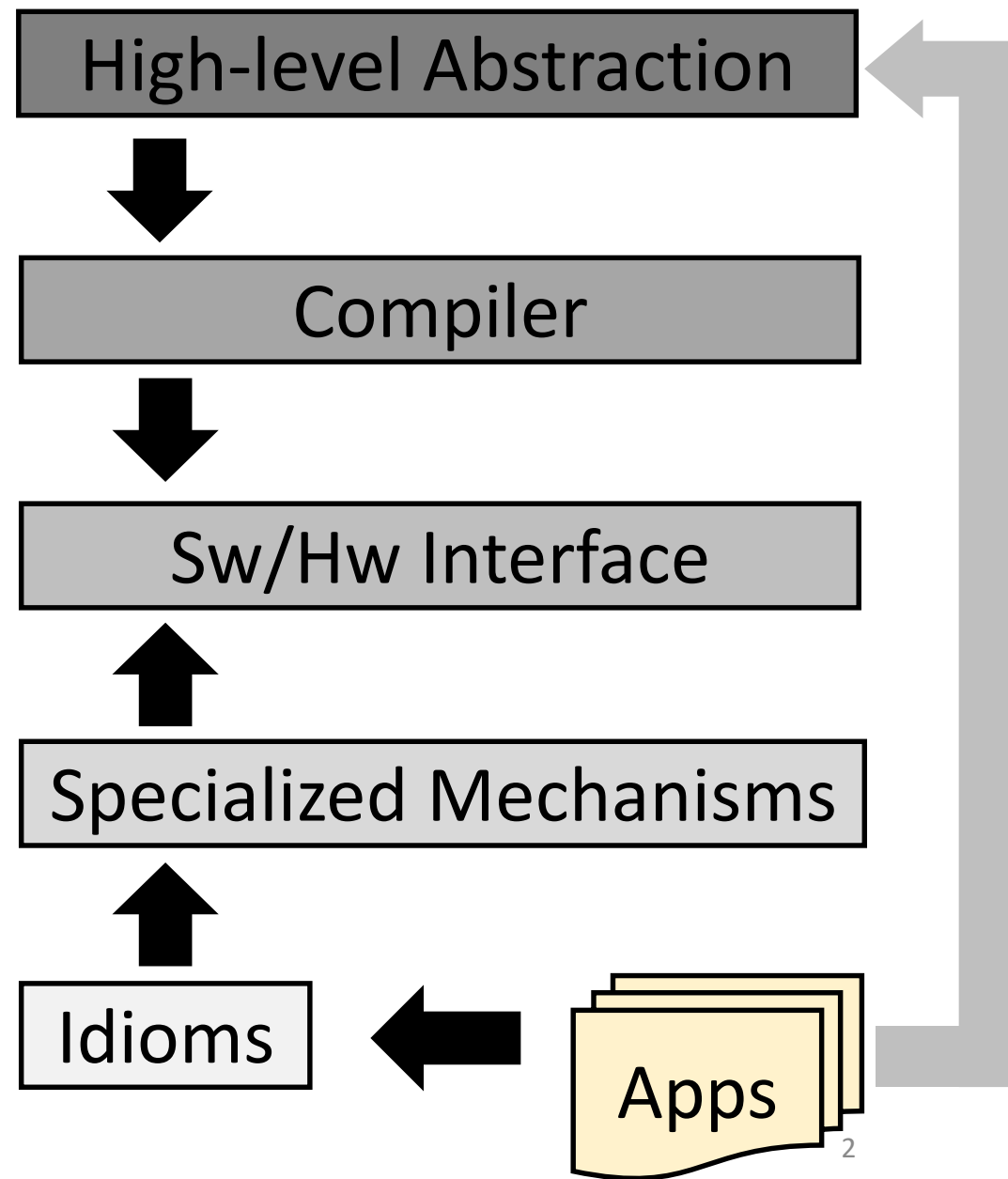University of California, Los Angeles

May 15th , 2020

# Specialized Accelerators

**Specialized architecture often occupies 1/5~1/3 of publications in top conferences.**
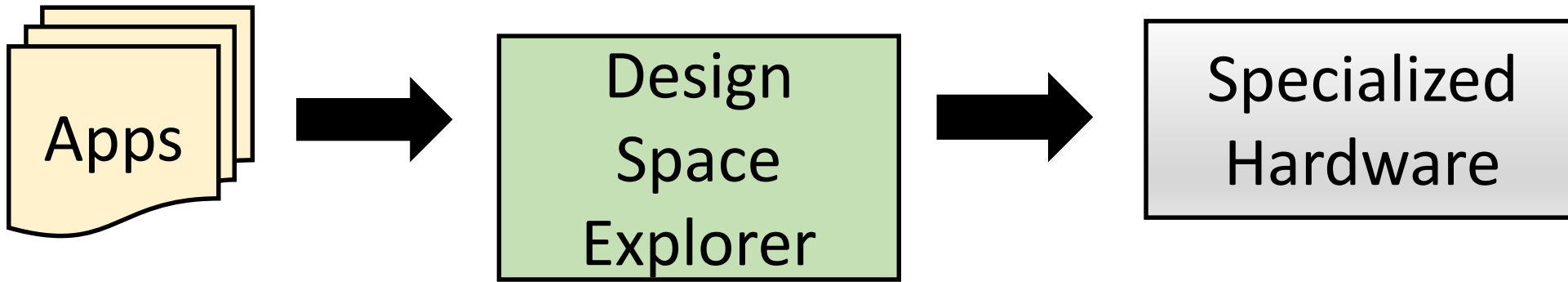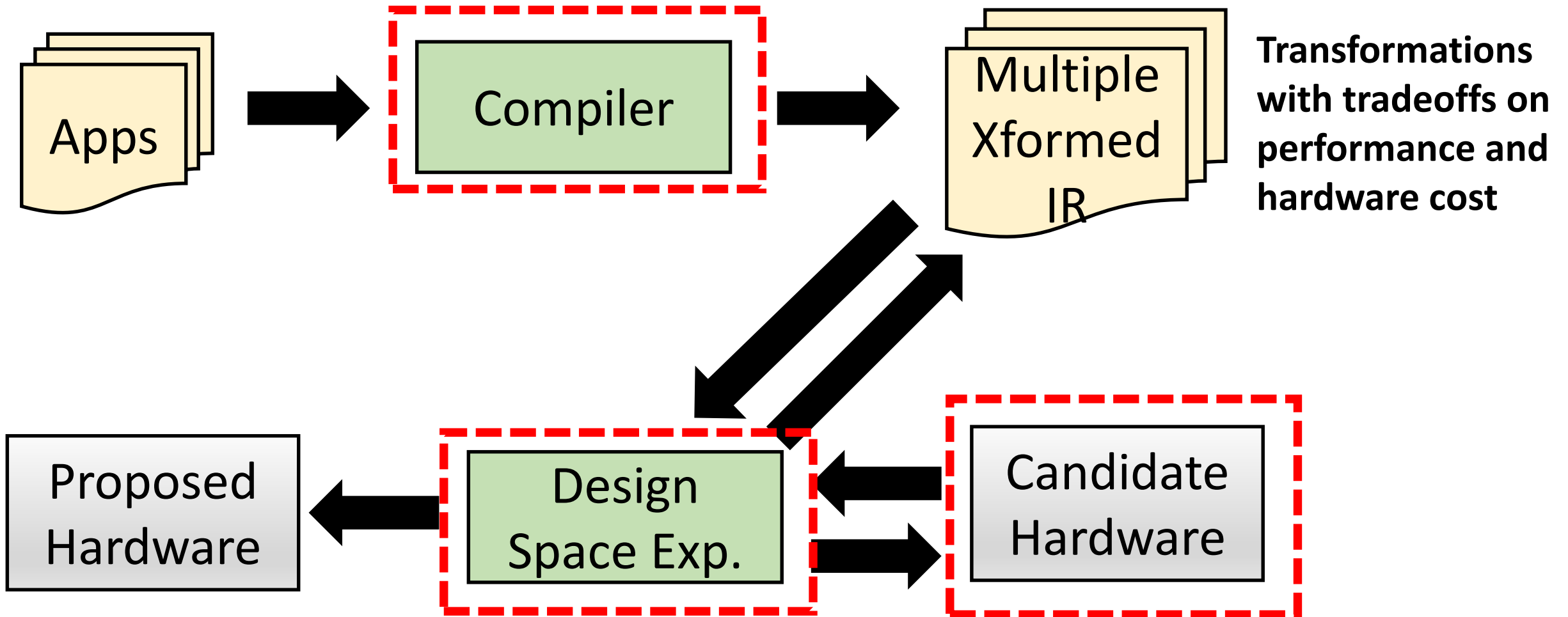


**Existing Domain-Specific Approach:**

High-level Abstraction

↓

Compiler

↓

Sw/Hw Interface

↑

Specialized Mechanisms

↑

Idioms ← Apps

# DSAGEN: Decoupled Spatial Accelerator Generator

## Domain   Specific



Apps → Design Space Explorer → Specialized Hardware

# DSAGEN: <u>D</u>ecoupled <u>S</u>patial <u>A</u>ccelerator <u>Ge</u>nerator

<span style="color:#b8cce4">Domain</span> <span style="color:#b8cce4">Specific</span>

Apps → Compiler → Multiple Xformed IR

**Transformations with tradeoffs on performance and hardware cost**

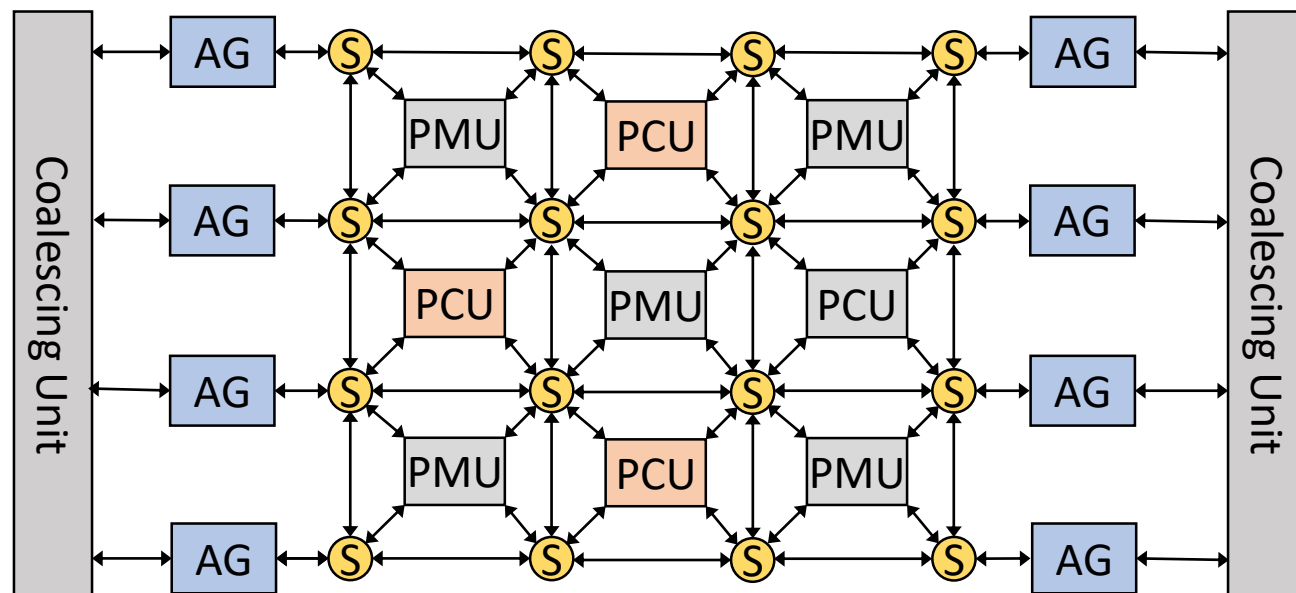Proposed Hardware ← Design Space Exp. ⟷ Candidate Hardware
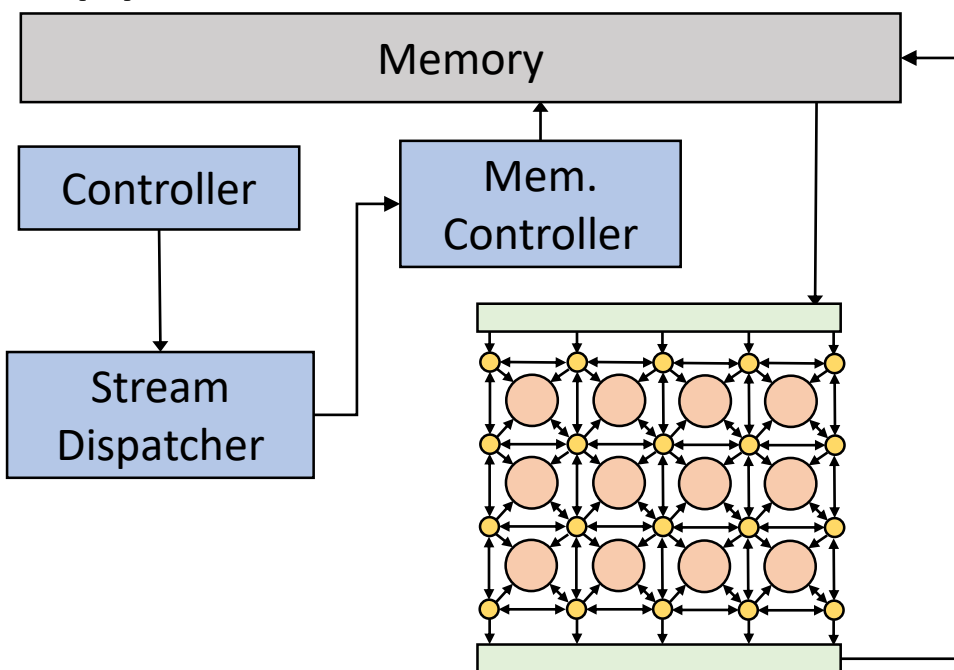
4

# Outline

- **Design Space — Decoupled-Spatial Architecture**
  - **Insight from Prior Work**
  - **The Programming Paradigm**
  - **Design Space: Hardware Primitives (& Composition)**
- Compilation
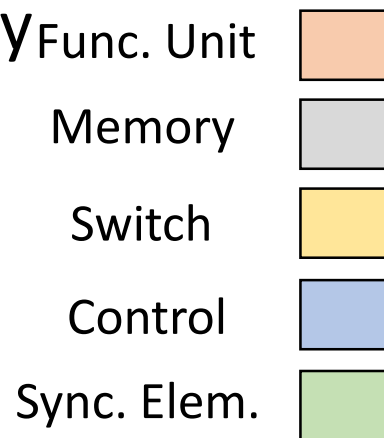- Design Space Exploration
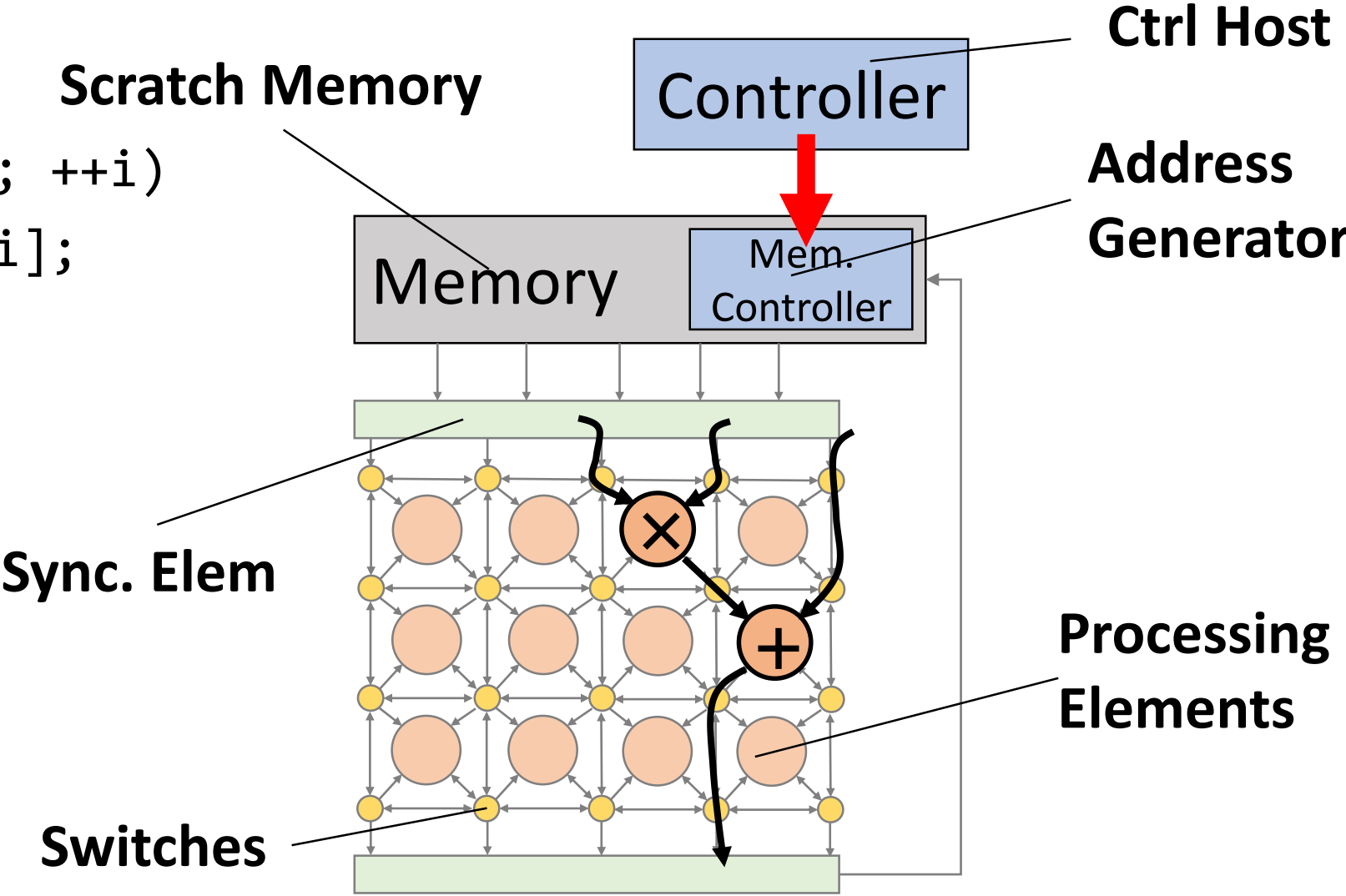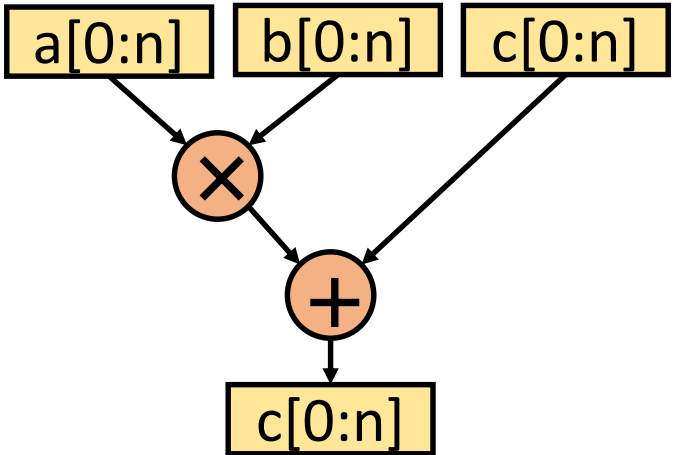- Evaluation

**(a) ASPLOS18-MAERI**



**(b) ISCA17-Plasticine**



**(c) ISCA17-Softbrain**

- Decoupled-Spatial Paradigm
  - Decoupled Compute/Memory
  - Spatially exposed resources
- Design Space
  - Composing hardware with simple primitives
  - Architecture Description Graph

Func. Unit
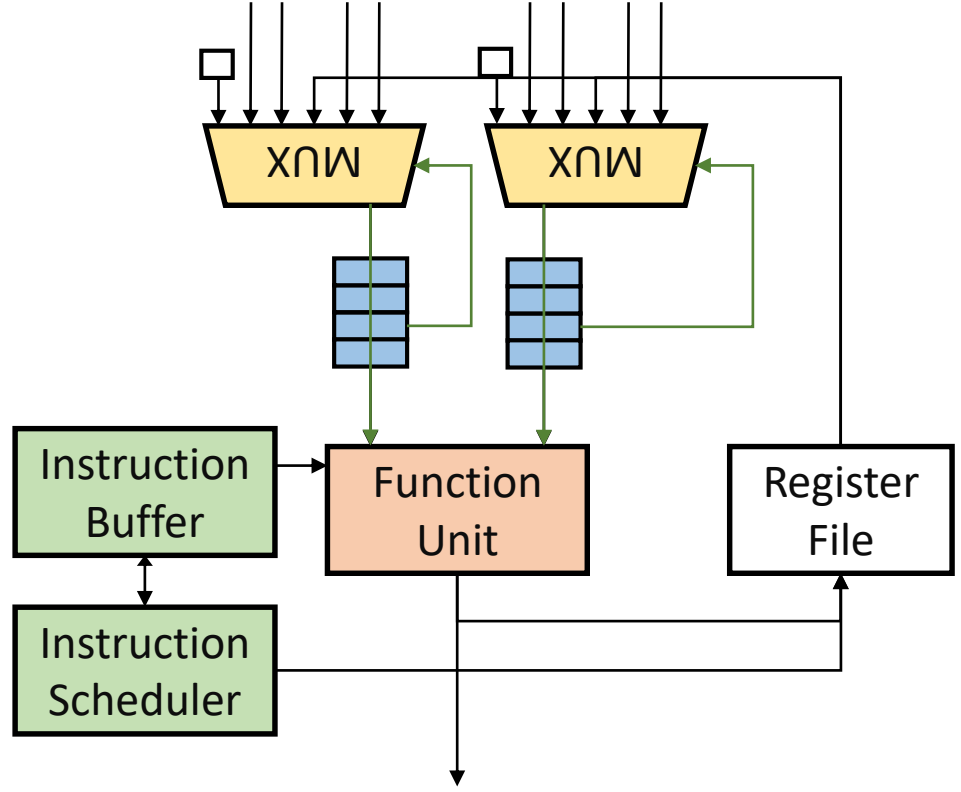
Memory

Switch

Control

Sync. Elem.

6

# Background: Decoupled-Spatial Architecture

```
for (int i = 0; i < n; ++i)
    c[i] += a[i] * b[i];
```

**Ctrl Host**

**Controller**

**Scratch Memory**

**Address Generator**

**Memory**

**Mem. Controller**

a[0:n]    b[0:n]    c[0:n]

$\times$

$+$

c[0:n]

**Sync. Elem**

$\times$

$+$

**Processing Elements**

**Switches**

# Hardware Primitives: Processing Element & Switch

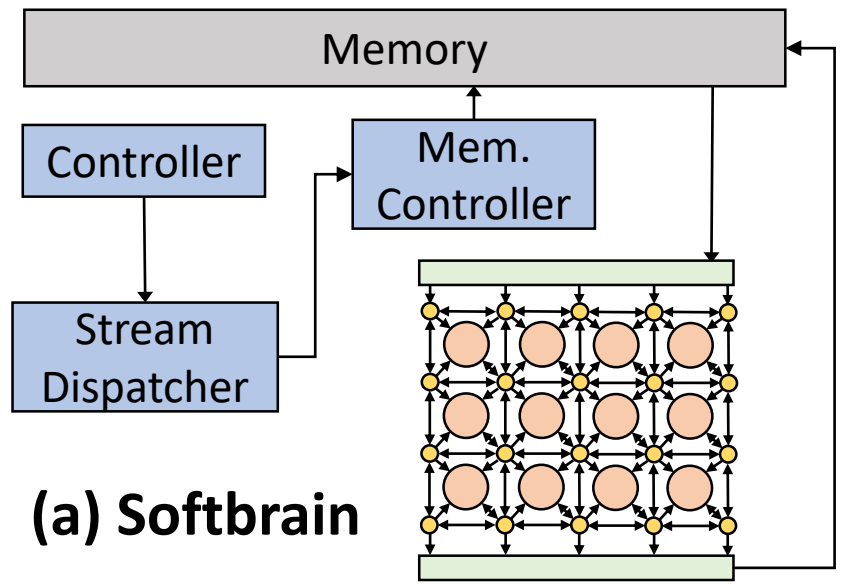| Hardware Cost: Low | Dedicated (=1) | Shared (>1) | High → |
|---|---|---|---|
| **Statically Scheduled** | **"Systolic" 1x Area**<br>+ No contention<br>- Harder to map<br>- Higher power<br>*Softbrain | **"CGRA" 2.6x Area**<br>+ Better resource utilization<br>- Harder to map<br>*Conventional CGRA | |
| **Dynamically Scheduled** | **"Ordered Dataflow" 2.1x Area**<br>+ Better flexibility<br>*SPU | **"Tagged Dataflow" 5.8x Area**<br>+ Better flexibility<br>+ Better resource utilization<br>*Triggered Instruction | |
| High ↓ | | | |

# Hardware Primitives: Memory

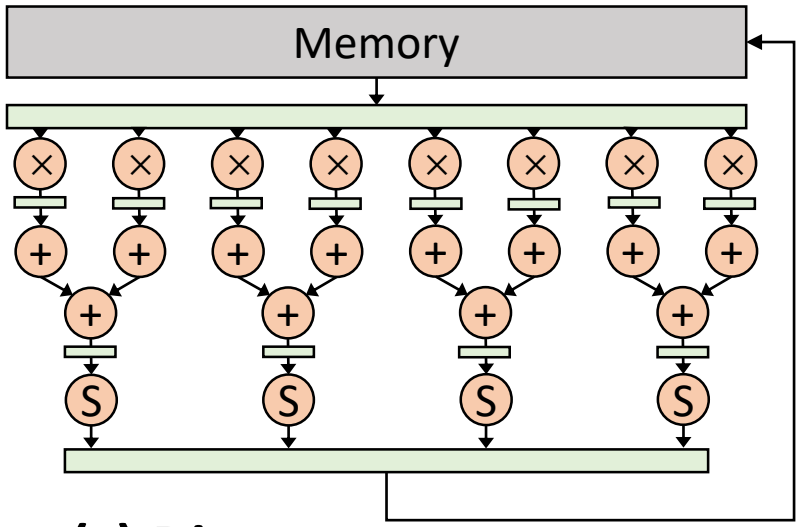- Memory
  - Size
  - Bandwidth
  - Indirect Support
    - a[b[i]]
  - Atomic Update
    - a[b[i]] += 1
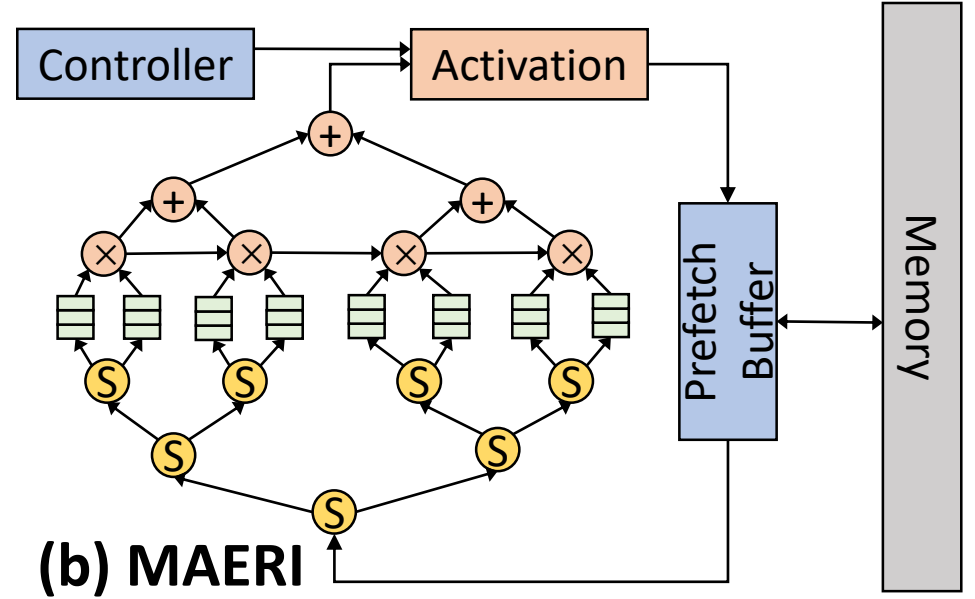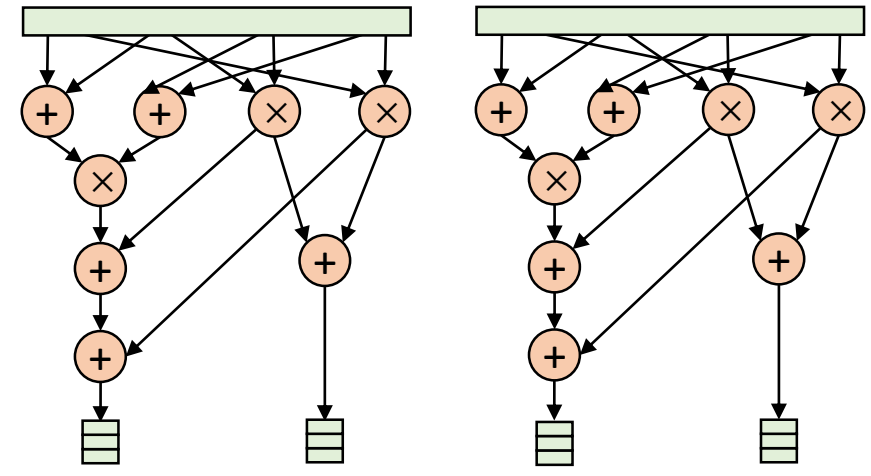
# Examples of ADG



(a) Softbrain

(b) MAERI

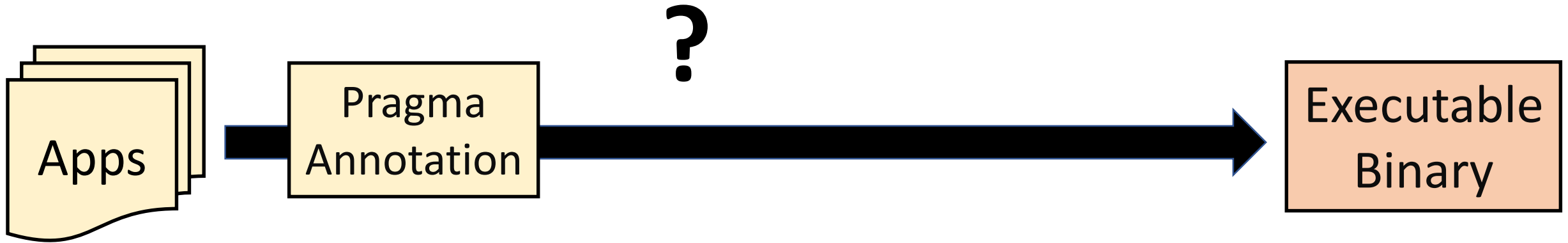(c) Diannao

(d) Data Path of Complex Mul.

# Outline

- <span style="color:lightgray">Decoupled-Spatial Architecture</span>
- **Compilation**
  - **High-Level Abstraction**
  - **Hardware-Aware Modular Compilation**
- <span style="color:lightgray">Design Space Exploration</span>
- <span style="color:lightgray">Evaluation</span>

# Compiling High-Level Lang. to Decoupled Spatial

**?**

Apps → Pragma Annotation → Executable Binary

**How to abstract diverse underlying features with a <span style="color:red">unified</span> high-level interface?**

- Programmer Hints
  - Which code regions are offloaded onto the spatial accelerator.
  - Which memory accesses can be decoupled intrinsics.
  - Which offloaded regions should be concurrent.

# An example of pragma annotation

```
#pragma config  ← The offloaded region in this compound body are concurrent
{
    #pragma stream  ← The memory accesses below will be stricted
    for (i=0; i<n; ++i)
        #pragma offload  ← The computational instructions below will be offloaded
        for (j=0; j<n; ++j)
            a[i*n+j] += b[c[j]] * d[i*n+j];
}
```
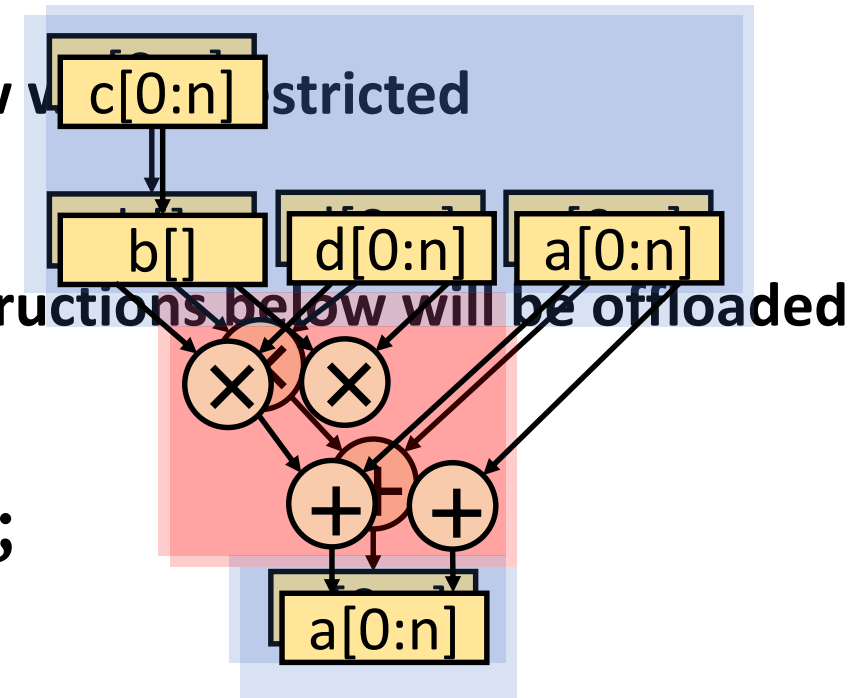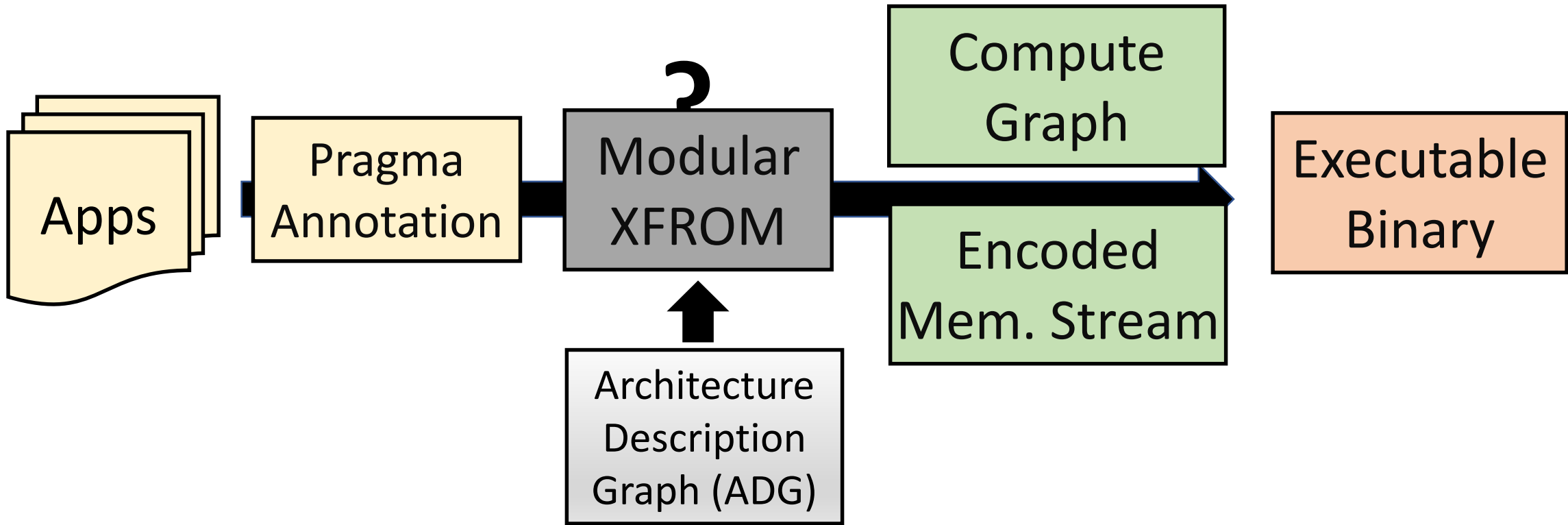
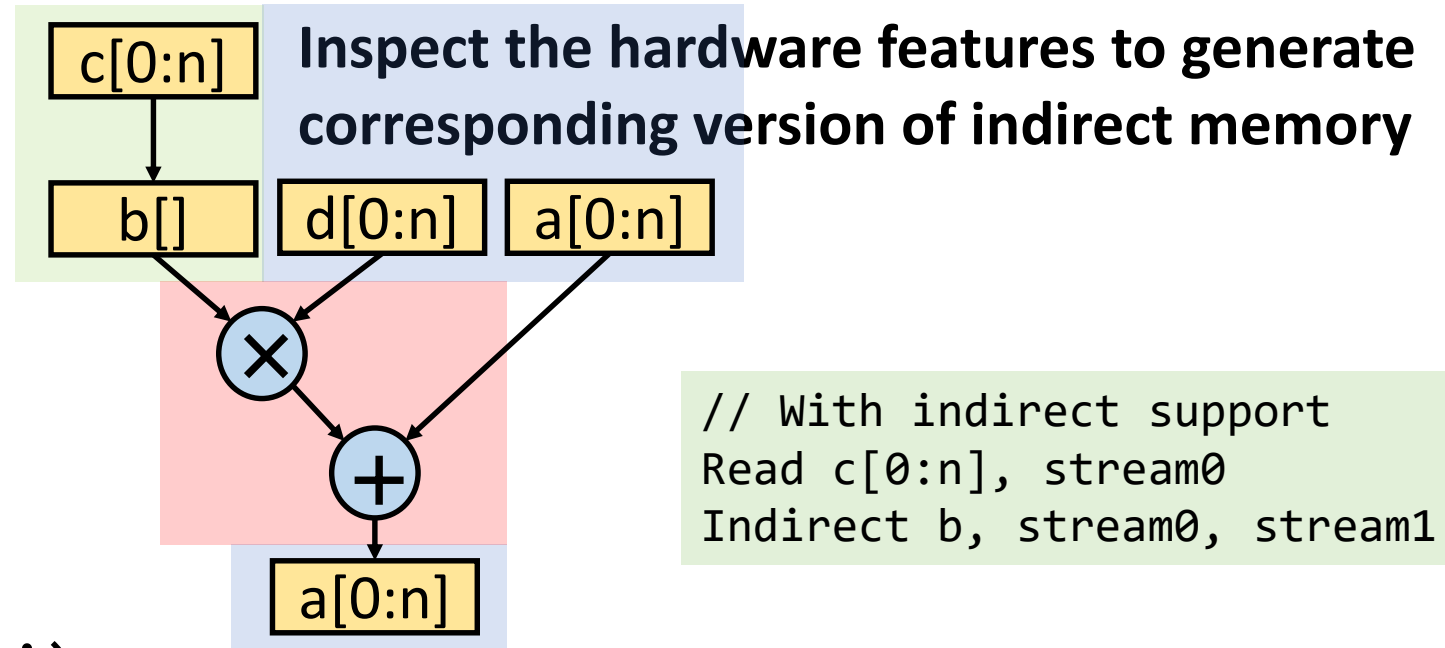# Compiling High-Level Lang. to Decoupled Spatial



**How to hide the diversity of underlying hardware?**

- Modular Transformation
  - Specialized Hardware features often dictate the code transformation
  - A fallback is required when the hardware feature is not available
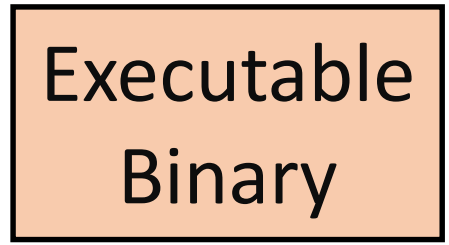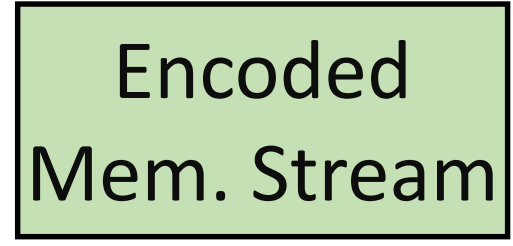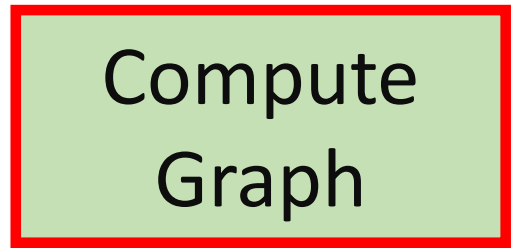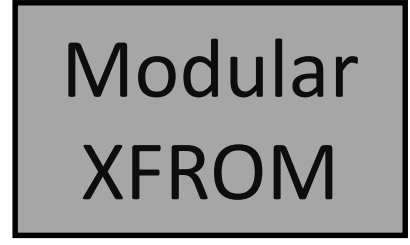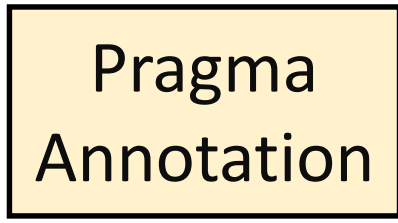
14

# Modular Transformation

```
#pragma config
{
    #pragma stream
    for (i=0; i<n; ++i)
        #pragma offload
        for (j=0; j<n; ++j)
            a[i*n+j] += b[c[j]] * d[i*n+j];
}
```



**Inspect the hardware features to generate corresponding version of indirect memory**

c[0:n]

b[]   d[0:n]   a[0:n]

⊗

⊕

a[0:n]

```
// With indirect support
Read c[0:n], stream0
Indirect b, stream0, stream1
```
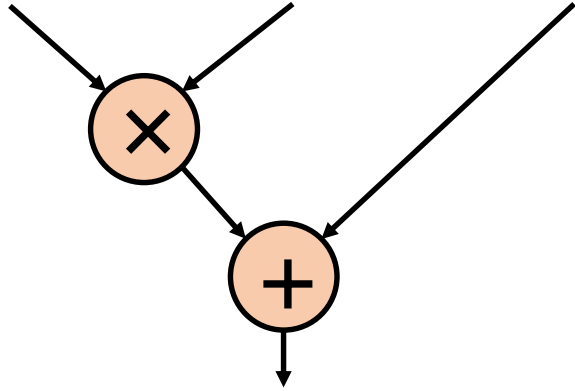
```
// Without indirect support
for (j=0; j<n; ++j)
    Scalar b[c[j]], stream0
```
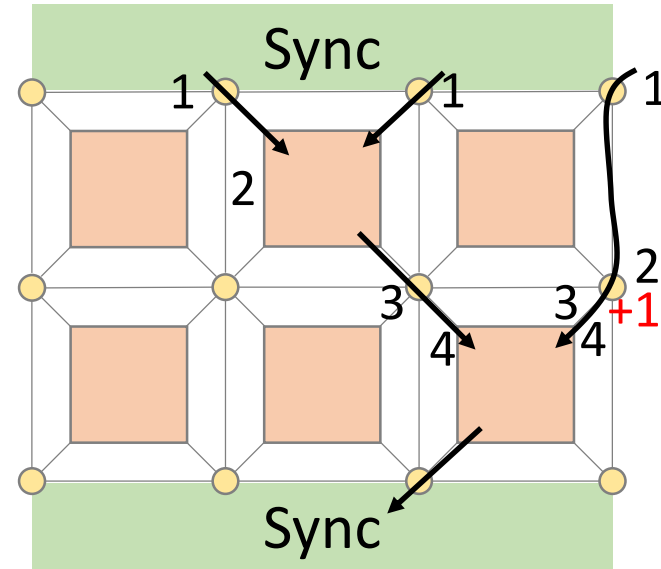
# Compiling High-Level Lang. to Decoupled Spatial

Apps

Pragma Annotation

Modular XFROM

Compute Graph

Encoded Mem. Stream

Executable Binary

**How is the dependence graph of computational instructions mapped?**

# Spatial Mapping



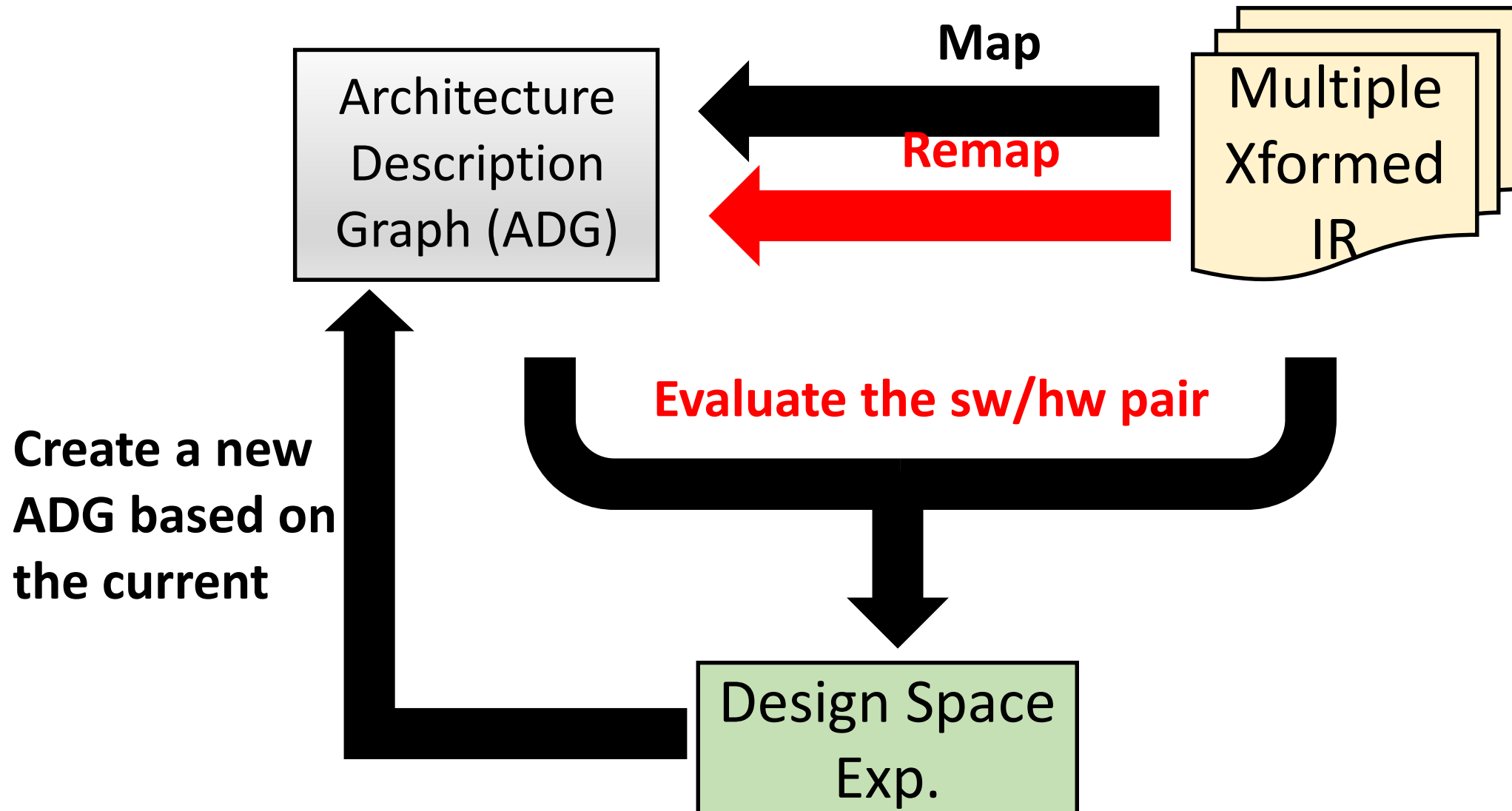**How is the dependence graph of computational instructions mapped?**

1. Placement: Map instruction to PE's with corresponding capability.

2. Routing: Routing the dependence edges thru the spatial network.

3. Timing: If necessary, balance the timing of data arrival

- If one of 1-3 is not successful, revert some nodes and repeat 123

# Outline

- Decoupled-Spatial Architecture
- Compilation
- **Design Space Exploration**
  - **Drive the Search**
  - **Evaluating Design Points**
  - **Repairing the Mapping**
- Evaluation

# Design Space Exploration

Architecture Description Graph (ADG)

**Map**

**Remap**

Multiple Xformed IR

**Evaluate the sw/hw pair**

**Create a new ADG based on the current**
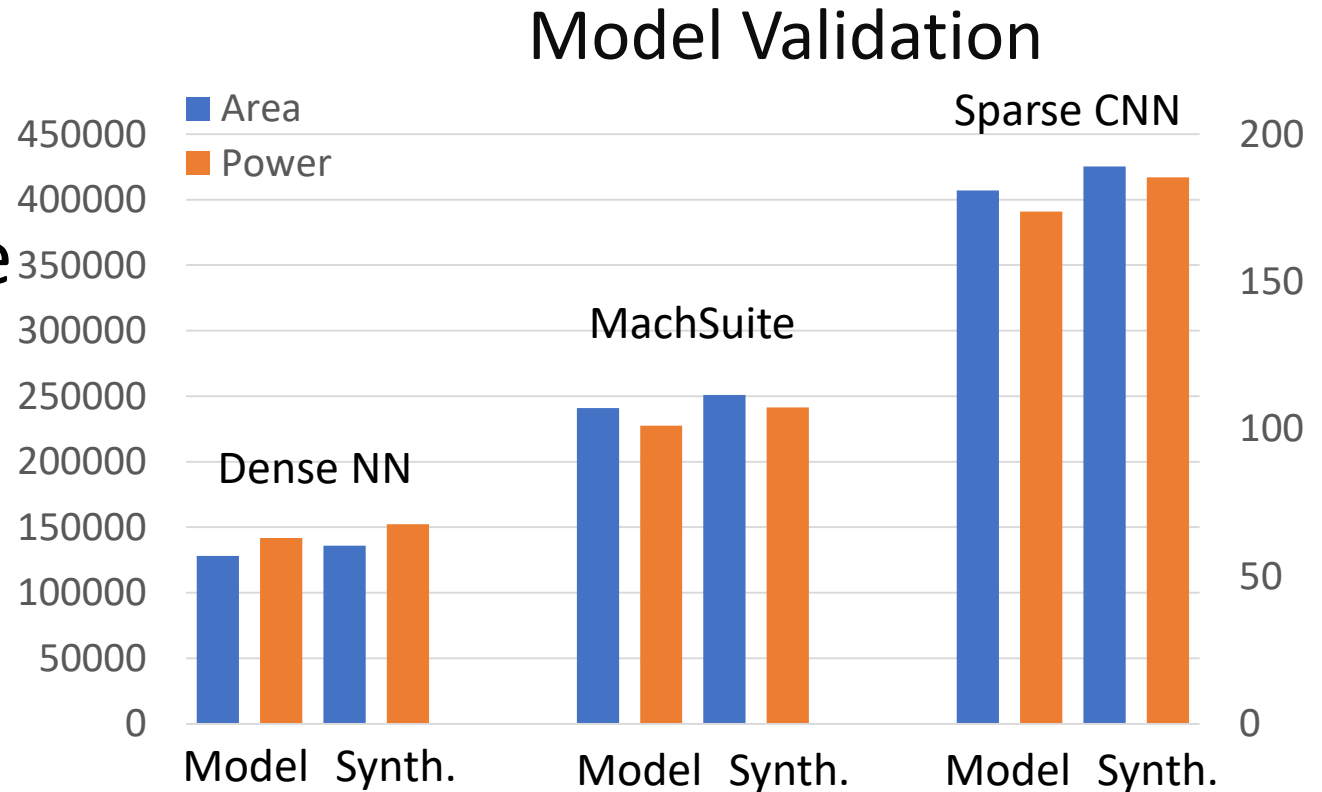
Design Space Exp.

# Estimation Model

- Performance
  - Spatial architecture essentially enables hardware specialized sw-pipelining
  - The ratio of data availability determines the performance
  - Perf=#Inst * (Activity Ratio)

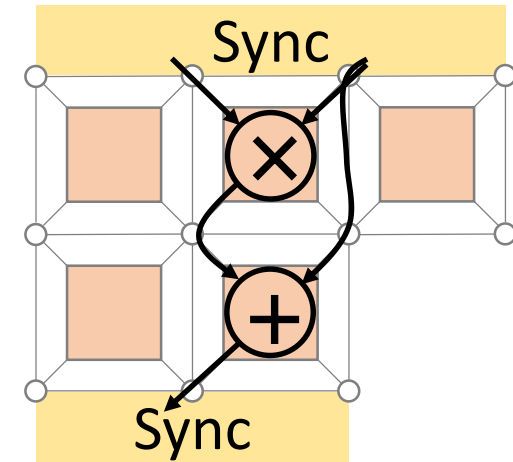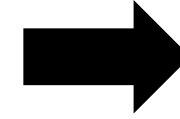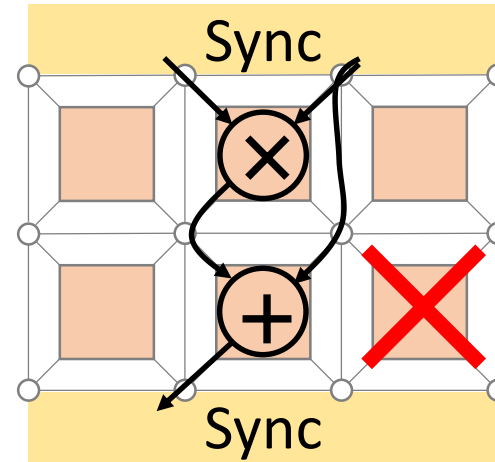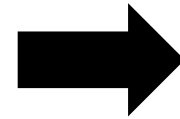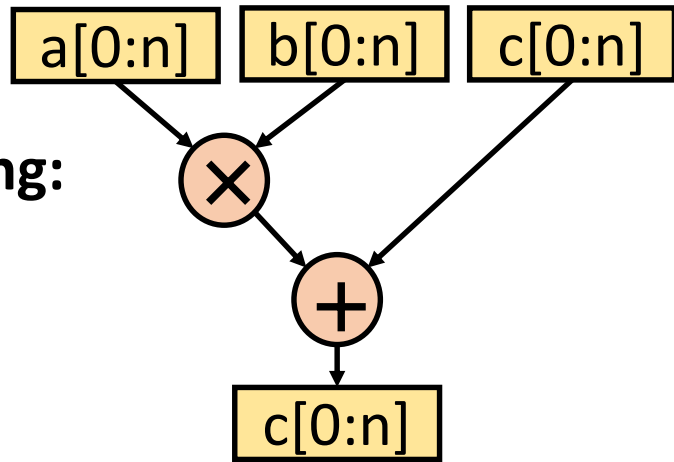The model has mean performance error of 7%, and with maximum error 30%.

- Power/Area
  - Synthesis can be time consuming
  - A regression model can predict the trend of hardware cost

Model Validation



Area
Power

Sparse CNN

MachSuite

Dense NN
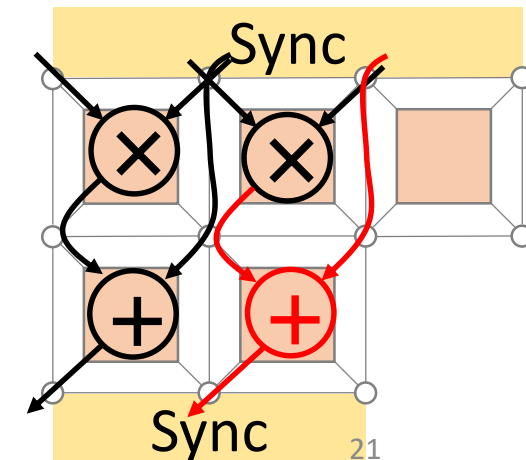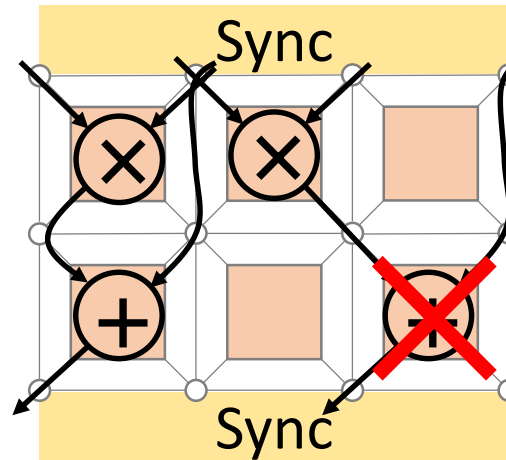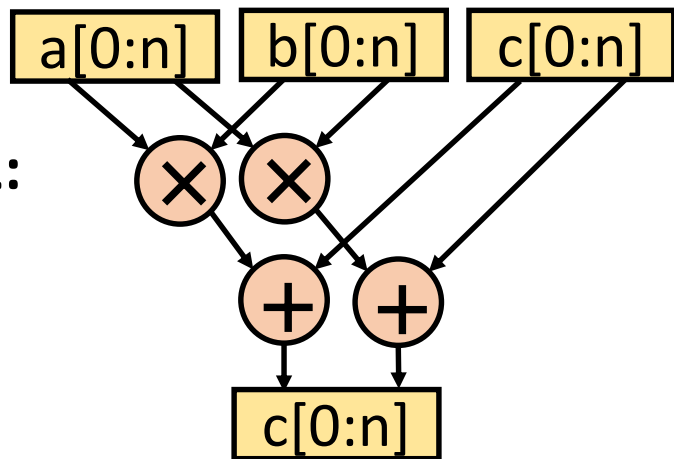
Model  Synth.    Model  Synth.    Model  Synth.

# Repairing the Spatial Mapping

```
// Original Code
for (i=0; i<n; ++i)
    c[i]+=a[i]*b[i];
```

**No Unrolling:**

**Unroll by 2:**

# Hardware/Software Interface Generation

- How to configure accelerator with arbitrary topology?
  - Reuse the data path for configuration
  - Find path(s) that cover(s) all the components
  - A heuristic based heuristic algorithm to minimize the longest path of configuration

- For a graph with m nodes covered by n paths, the longest path cannot be shorter than $\lceil \frac{m}{n} \rceil$.

- We only introduces 40% overhead over the ideal bound.

# Outline

- Decoupled-Spatial Architecture
- Compilation
- Design Space Exploration
- **Evaluation**
  - **Methodology**
  - **Compiler**
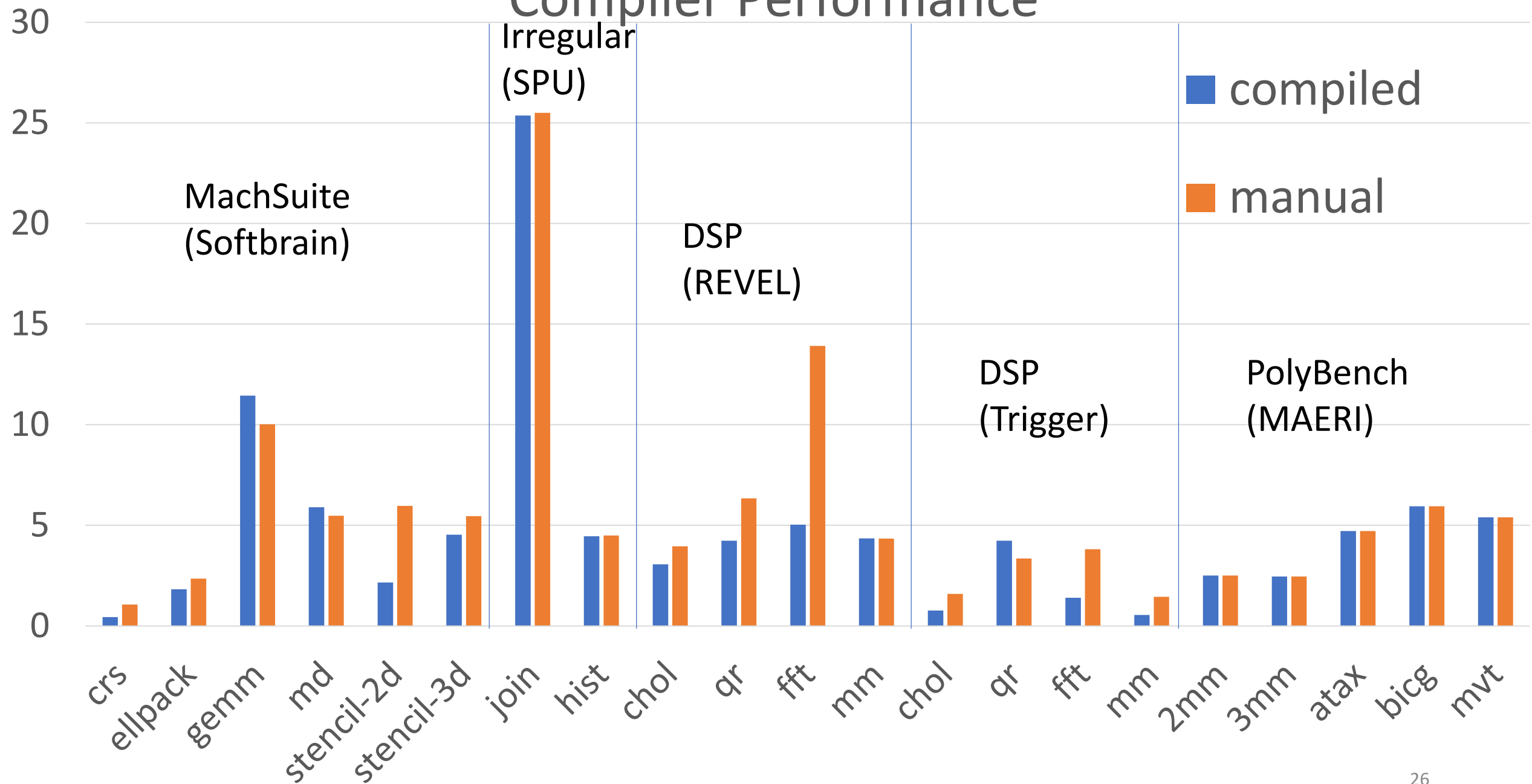  - **Design Space Exploration**

# Methodology

- Performance
  - Gem5 RISCV in-order core integrated with a cycle-accurate spatial accelerator simulator
    - The in-order core is extended with stream decoupled ISA
- Power/Area
  - All the components are implemented in Chisel RTL
  - Synthesized in Synopsys DC 28nm @1.00GHz
  - SRAM power/area are estimated by CACTI 7.0

# Compiler Performance

- Softbrain — MachSuite
  - Versatile accelerator can handle moderate irregularity
- SPU — Histogram, and Key Join
  - Accelerator specialized for irregular workloads
- REVEL and Trigger — DSP
  - Accelerator specialized for imperfect loop body
- MAERI — PolyBench
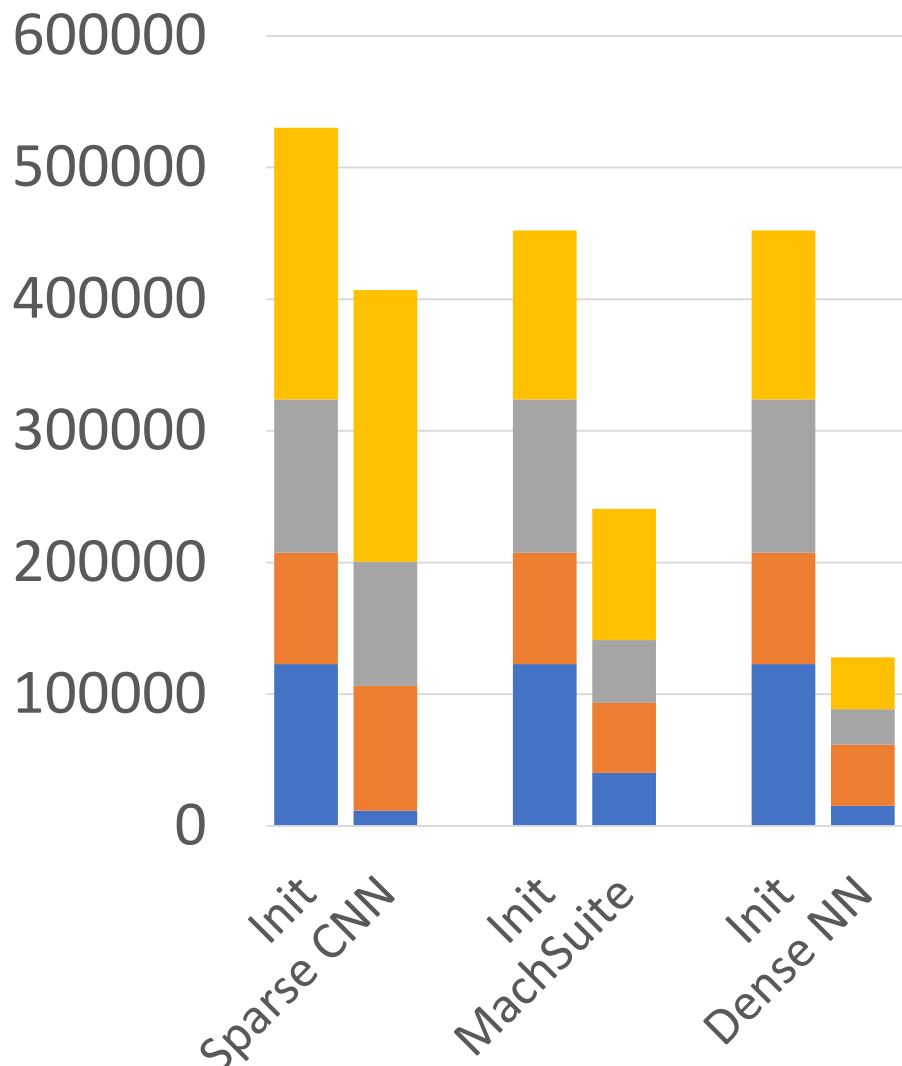  - Accelerator for neural network

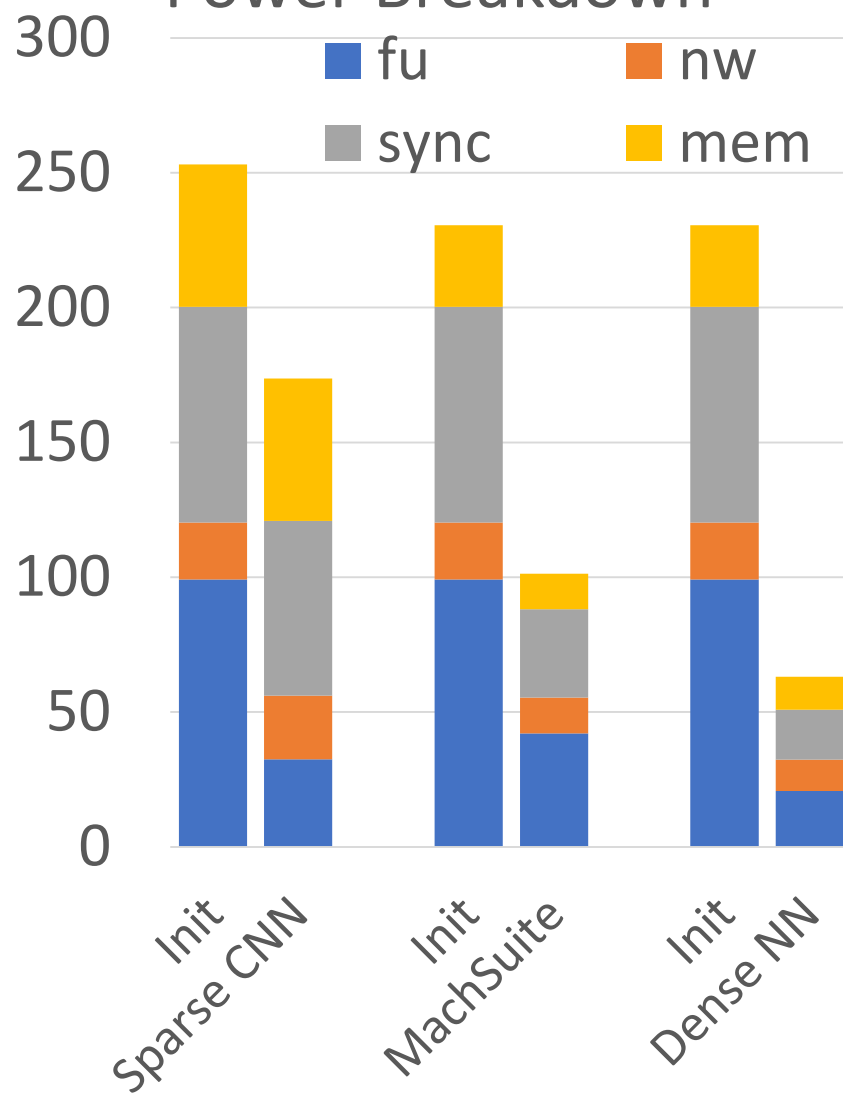Compiler Performance

# Design Space Explorer

- Workloads
  - Dense Neural Network
  - MachSuite
  - Sparse Convolutional Neural Network
- Initial Design
  - A 5x5 mesh with all capability (arithmetic, control, and indirect)
- Objective: $perf^2/mm^2$

# Design Space Explorer

## Area Breakdown



## Power Breakdown



Sparse CNN: 24h
MachSuite: 19.2h
Dense NN: 16h

# Conclusion

| | HLS | Manual | DSAGEN |
|---|---|---|---|
| Frontend | C+Pragma | DSL/Intrinsics, etc. | C+Pragma |
| Design Flow | Nearly Automated | Manual | Nearly Automated |
| Input | A **Single** Application | Multiple Target Applications | **Multiple** Target Applications |
| Output | **Application-Specific** Accel. | ASIC/Programmable Accel. | A **Programmable** Accelerator |
| Design Space | **Limited** | **Rich** | **Rich** |

# Q&A

- Our framework is working in progress at: https://github.com/PolyArch/dsa-framework
- All the questions and comments are welcomed