# Halide and Clockwork Compiler
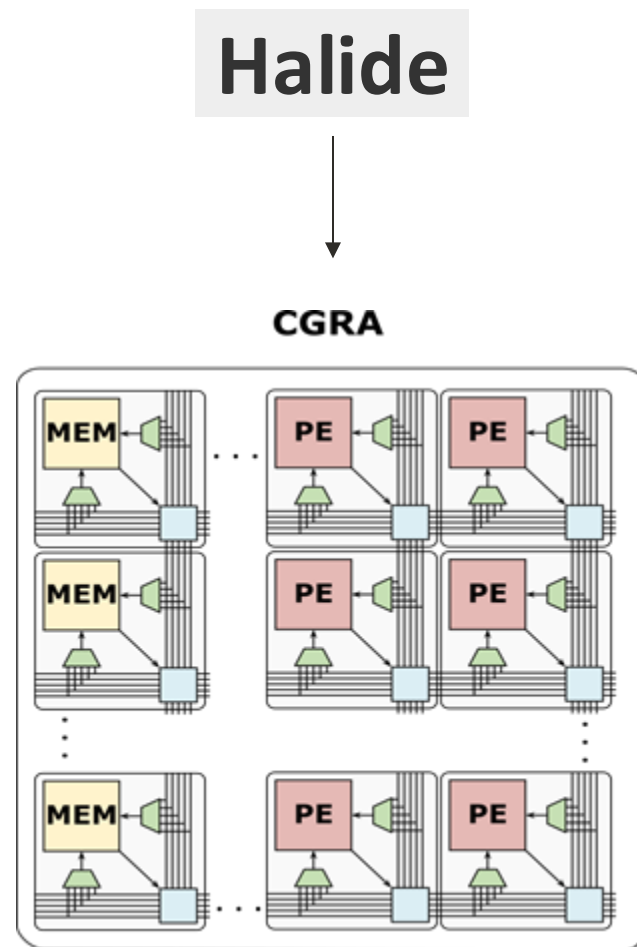
Jeff Setter

# Overview of our System

**Halide**

Applications: in Halide DSL (stencils + DNNs)

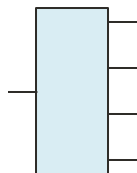Hardware: CGRAs

- Coarse-Grained Reconfigurable Array

CGRA

# Talk Overview:

1) **Halide Scheduling**

```
app += algorithm;
app.scheduling();
```

2) Clockwork to map
Unified Buffers

# Halide

Motivation: need a concise way to describe specializations for each hardware target

=> Separate algorithm from schedule

<u>Algorithm</u>: description of the computation (*what* output values)

- Mathematical operations

<u>Schedule</u>: optimization decisions (*how* to compute)

- Loop optimizations to run fast

Jonathan Ragan-Kelley, Connelly Barnes, Andrew Adams, Sylvain Paris, Fredo Durand, and Saman Amarasinghe.
*Halide: A language and compiler for optimizing parallelism, locality, and recomputation in image processing pipelines.* In PLDI 2013, page 519-530.
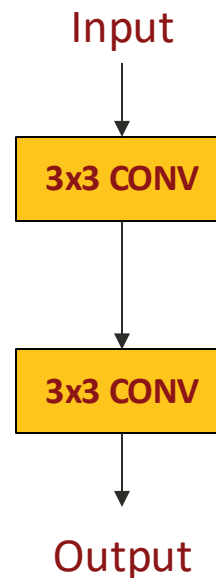
# Example Algorithm: cascade for CPU

```
// Algorithm:
// Define 3x3 window of kernel weights
RDom r(0, 3, 0, 3); // Reduction
kernel(0,0) = 1;  kernel(0,1) = 2;  ... kernel(2,2) = 1;

// First convolution
conv1(x, y) += kernel(r.x, r.y) * input(x + r.x, y + r.y);
conv1_norm(x, y) = conv1(x,y) / 16;

// Second convolution
conv2(x, y) += kernel(r.x, r.y) * conv1_norm(x + r.x, y + r.y);
conv2_norm(x, y) = conv2(x,y) / 16;

output(x, y) = conv2_norm(x, y);
```

Input

3x3 CONV

3x3 CONV

Output

# Halide Scheduling: parallelism and memory

Efficient and fast application execution relies on **parallelism** and good **memory management**.

Optimizations for CPU

- Loop parallelism: unrolling loops, vectorization, threading
- Memory: tiling, memory granularity (fusion)

# Halide Scheduling: parallelism and memory

Efficient and fast application execution relies on **parallelism** and good **memory management**.

Optimizations for CPU

- Loop parallelism: unrolling loops, vectorization, threading
- Memory: tiling, memory granularity (fusion)

Scheduling for **CGRAs**

- Hardware parallelism: PEs simultaneously executing
- Memory: tiling, memory hierarchy, specialized for streaming

Jeff Setter. *Compiling Image Processing and Machine Learning Applications to Reconfigurable Accelerators*. In Stanford Doctoral Dissertation, 2023.

# Halide Scheduling to CGRA

Reused the existing Halide compiler toolchain as much as possible without introducing new primitives.

Hardware scheduling:

- <u>Accelerator</u>: input, output, image size
- <u>Loops</u>: reorder loops, tiling
- <u>Memory</u>: memory temporaries, hierarchy
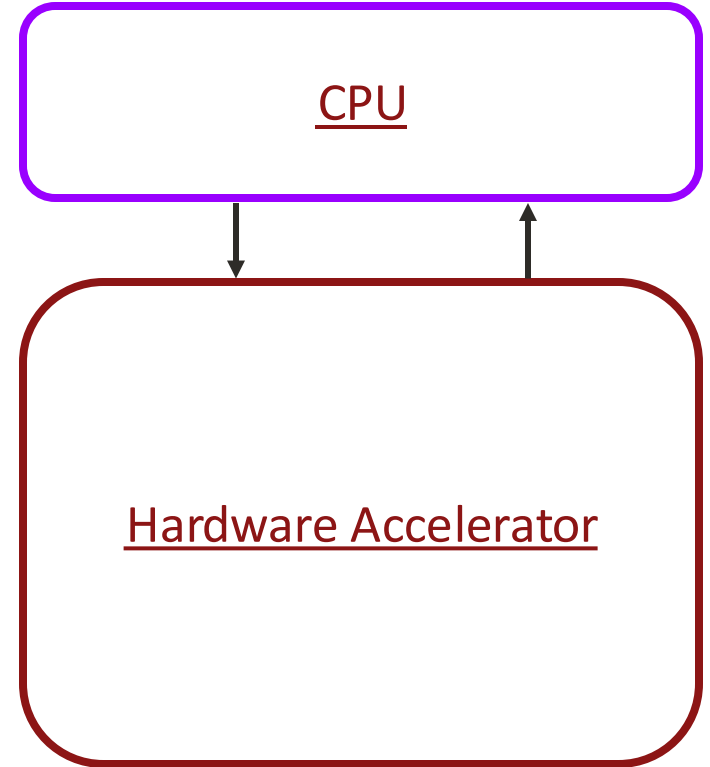- <u>Compute</u>: duplication

# Define Accelerator's Scope

`hw_accelerate(xi, xo)`

- Defines the accelerator output
- Argument specifies outermost loop

`stream_to_accelerator()`

- Defines an accelerator input



CPU

Hardware Accelerator

# Define Loops: strip mining, reordering, tiling

`split(x, xo, xi, 64)`
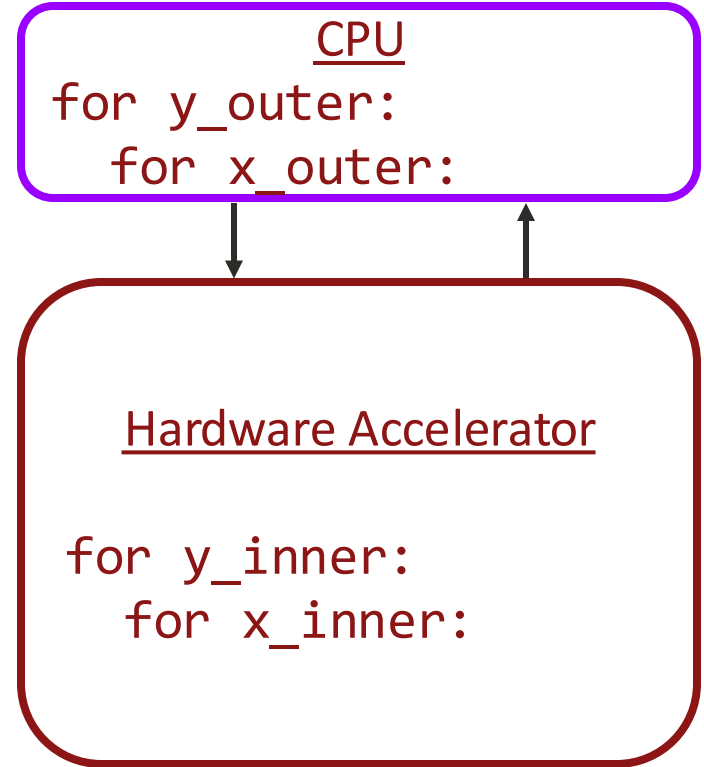
- create nested loops
- specified size for inner loop

`reorder(xi, yi, xo, yo)`

- new loop order (from inner to outer)

`tile(x,y, xo,yo, xi,yi, 64,64)`

- split two loops and reorder
- common for images

CPU
```
for y_outer:
    for x_outer:
```

Hardware Accelerator

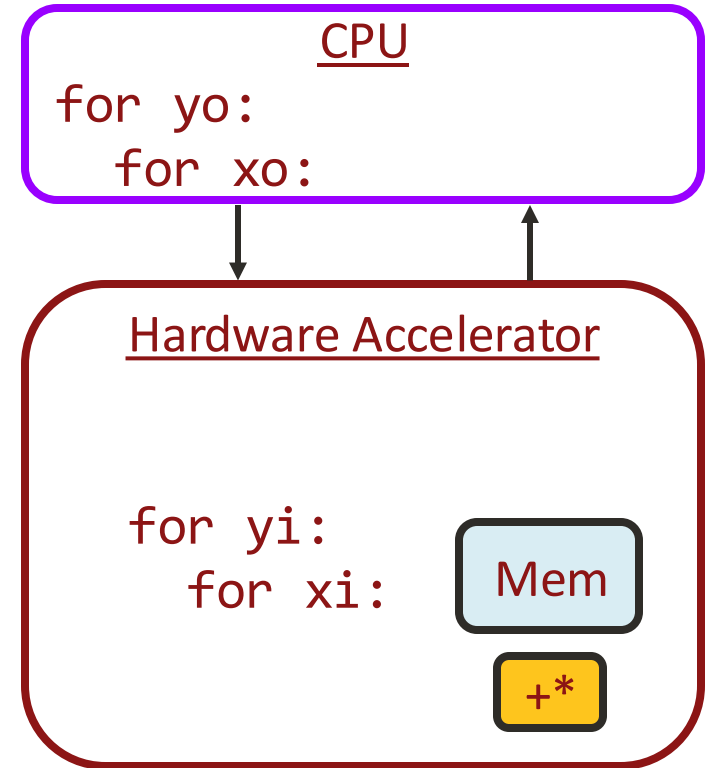```
for y_inner:
    for x_inner:
```

# Define Memories: temporaries

`store_at(consumer, xo)`

- create buffers between compute
- at which loop level to create buffer

`compute_at(consumer, xo)`

- for CGRA, we compute at tile level and leave loop fusion for Clockwork
- <u>not</u>: choosing line buffer or double buffer



CPU
```
for yo:
    for xo:
```

Hardware Accelerator
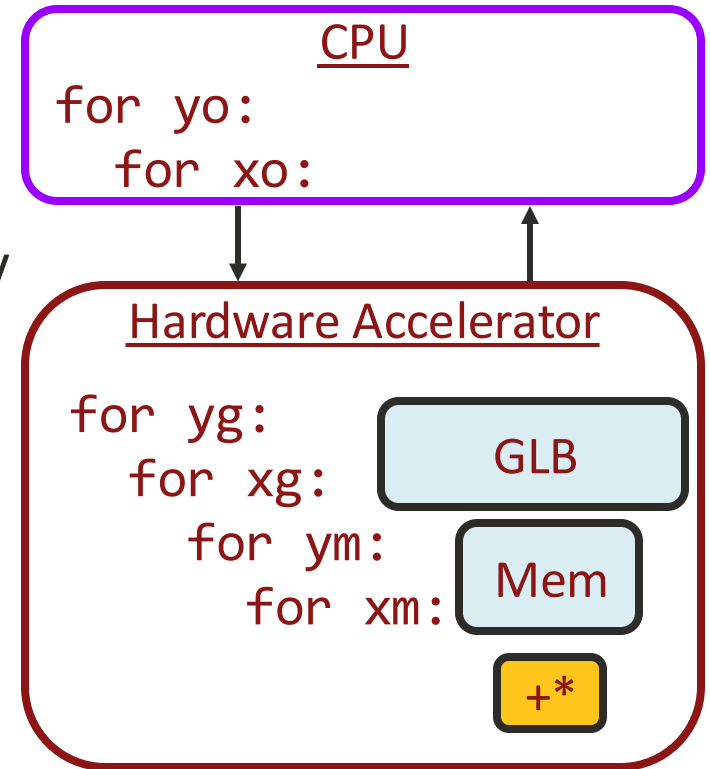
```
for yi:
    for xi:
```
Mem

+*

# Define Memories: hierarchy

`in()`

- creates a copy (identity function)
- Needed to move data between memory hierarchy levels with `tile` and `compute_at`

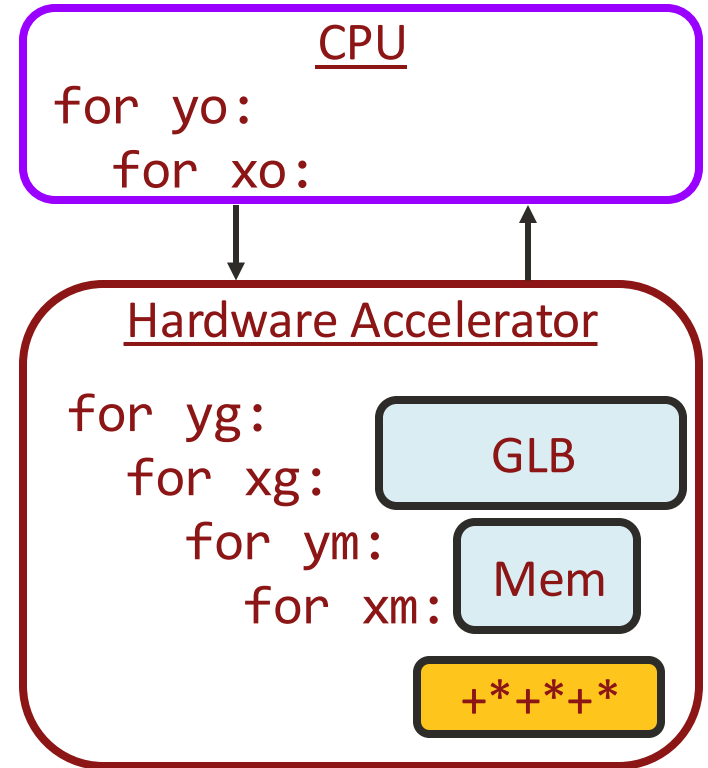`store_in(MemoryType::GLB)`

- specify what storage type to use
- Examples: Host, GLB, Memory, Pond

CPU
```
for yo:
  for xo:
```

Hardware Accelerator
```
for yg:
  for xg:          GLB
    for ym:
      for xm:      Mem

                   +*
```

# Define Compute: duplication and sharing

`unroll(r.y, 3)`

- increases hardware parallelism
- in Halide IR duplicates loop, because multiple statements = utilizes more existing hardware

```
CPU
for yo:
  for xo:
```

```
Hardware Accelerator
for yg:
  for xg:          GLB
    for ym:
      for xm:      Mem

              +*+*+*
```

```
1  /* CGRA Schedule */
2  Var xii, xio, yii, yio, xi, xo, yi, yo;
3  output.bound(x, 0, outImgSizeX)
4         .bound(y, 0, outImgSizeY);
5
6  // Accelerate 360x360 tiles in GLB
7  hw_output.in().compute_root()
8      .tile(x,y, xo,yo, xi,yi, 360,360)
9      .hw_accelerate(xi, xo);
10
11 // Send 60x60 tiles to CGRA fabric
12 hw_output
13     .tile(x,y, xio,yio, xii,yii, 60,60)
14     .compute_at(hw_output.in(), xo)
15     .store_in(MemoryType::GLB);
16
17 // conv2 kernel with unrolled reduction
18 conv2_norm.compute_at(hw_output, xio);
19 conv2.update()
20     .unroll(r2.x).unroll(r2.y);
21
22 // conv1 kernel with unrolled reduction
23 conv1_norm.compute_at(hw_output, xio);
24 conv1.update()
25     .unroll(r.x).unroll(r.y);
26
27 // MemoryTile and GLB for input stream
28 hw_input.in().in()
29     .compute_at(hw_output, xio);
30 hw_input.in()
31     .compute_at(hw_output.in(), xo)
32     .store_in(MemoryType::GLB);
33 hw_input.accelerator_input();
34
35 kernel.compute_at(hw_output, xio);
```

CPU - output

Hardware Accelerator

GLB

conv2

Mem

conv1

Mem

GLB

CPU - input

# Halide Compiler and Codegen

- Compiler: performs loop transformations

- Codegen to Clockwork: generates files of **compute kernels** and buffers defined as **sequences of memory accesses**
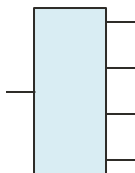
```
1   // Cascade has a unified buffer named "conv2_stencil" storing 16-bit values
2   prg.buffer_port_widths["conv2_stencil"] = 16;
3
4   // These loops correspond to the iteration domain
5   auto y = prg.add_loop("y", 0, 99);
6   auto x = y->add_loop("x", 0, 600);
7
8   // The function indicates which compute kernel to use
9   auto hcompute_conv2_stencil_1 = x->add_op("op_hcompute_conv2_stencil_1");
10  hcompute_conv2_stencil_1->add_function("hcompute_conv2_stencil_1");
11
12  // These loads are access maps for outputs of the "conv1_stencil" buffer
13  hcompute_conv2_stencil_1->add_load("conv1_stencil", "y", "x");
14  hcompute_conv2_stencil_1->add_load("conv1_stencil", "y", "(x + 1)");
15  hcompute_conv2_stencil_1->add_load("conv1_stencil", "y", "(x + 2)");
16  hcompute_conv2_stencil_1->add_load("conv1_stencil", "(y + 1)", "x");
17  hcompute_conv2_stencil_1->add_load("conv1_stencil", "(y + 1)", "(x + 1)");
18  hcompute_conv2_stencil_1->add_load("conv1_stencil", "(y + 1)", "(x + 2)");
19  hcompute_conv2_stencil_1->add_load("conv1_stencil", "(y + 2)", "x");
20  hcompute_conv2_stencil_1->add_load("conv1_stencil", "(y + 2)", "(x + 2)");
21  hcompute_conv2_stencil_1->add_load("conv1_stencil", "(y + 2)", "(x + 1)");
22  hcompute_conv2_stencil_1->add_load("conv2_stencil", "y", "x");
23
24  // This store is an access map for the input of the "conv2_stencil" buffer
25  hcompute_conv2_stencil_1->add_store("conv2_stencil", "y", "x");
```

# Talk Overview:

1) Halide Scheduling

```
app += algorithm;
app.scheduling();
```

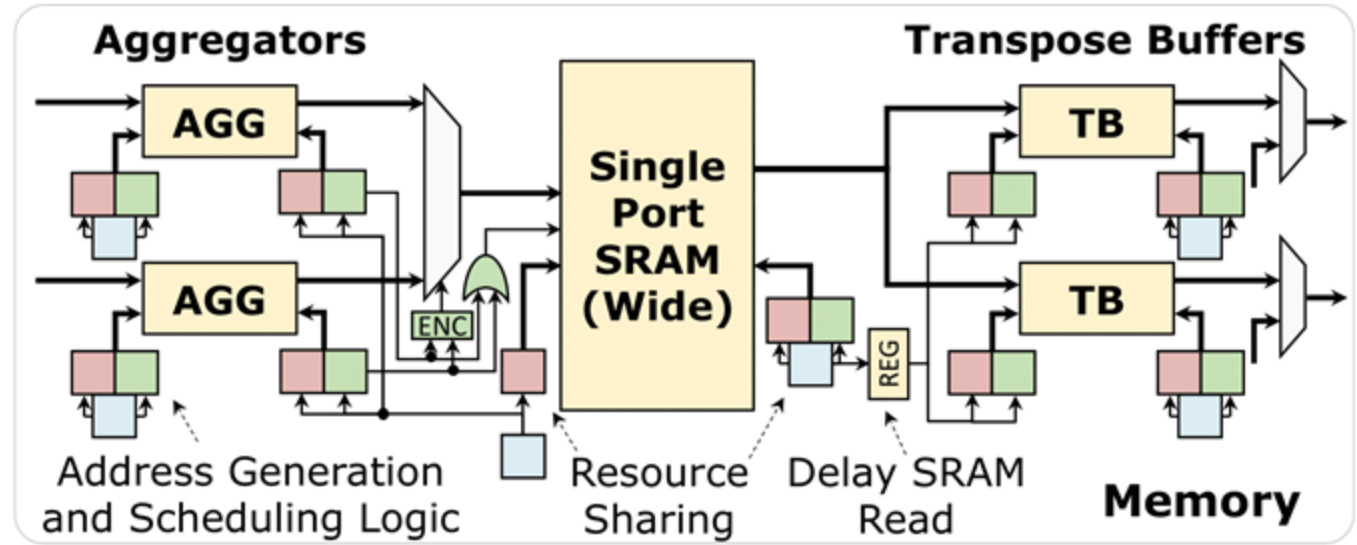2) **Clockwork to map Unified Buffers**

# Unified Buffer Motivation

**Motivation**:

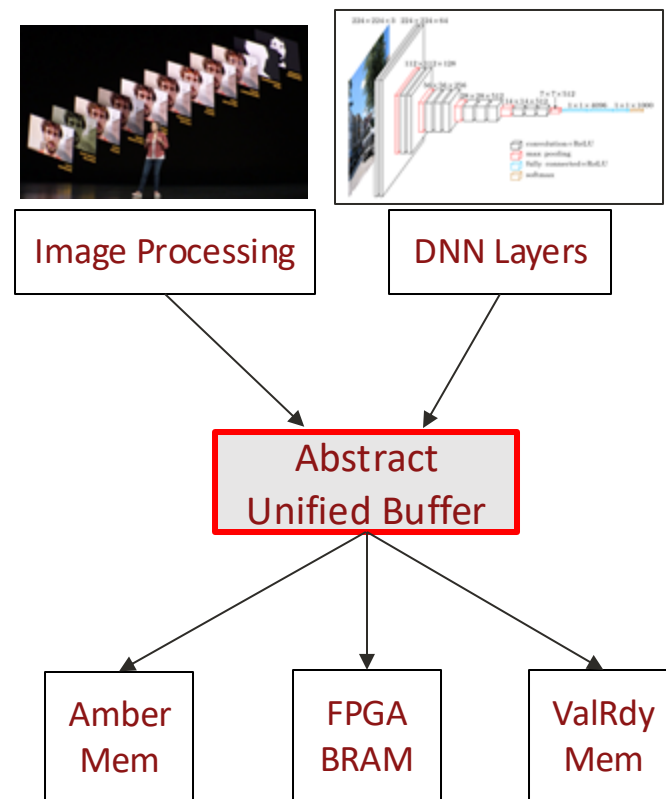need abstraction for memories in streaming accelerators



Push memories are **efficient**, but also **complex**

# Unified Buffer Abstraction

Create an abstraction to **interface** between the Halide **application** and the **hardware** motifs
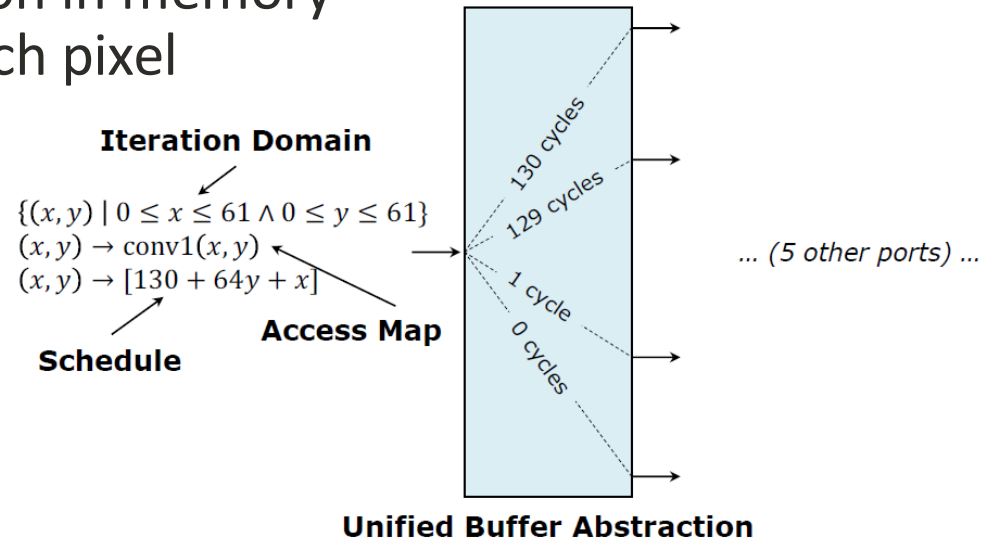
- <u>Domain</u>: fully describes streaming applications (stencils, DNNs)

- <u>Hardware Mapping</u>:
  - Sufficiently general to map to dataflow accelerators classes (FPGA, CGRA; statically scheduled, ready-valid)
  - with efficiency of common hardware stream memories (linebuffers, double buffers)



Image Processing | DNN Layers

Abstract Unified Buffer

Amber Mem | FPGA BRAM | ValRdy Mem

Qiaoyi Liu, Jeff Setter, Dillon Hu, Maxwell Strange, Kathleen Feng, Mark Horowitz, Priyanka Raina, and Fredrik Kjolstad. *Unified buffer: Compiling image processing and machine learning applications to push-memory accelerators*. TACO 2023.

# Unified Buffer Properties

Properties for each memory port:

- <u>dependencies</u>: relation to other memories
- <u>iteration domain</u>: valid indices
- <u>access map</u>: relative location in memory
- <u>op schedule</u>: cycle time each pixel executes

**Iteration Domain**

$$\{(x, y) \mid 0 \leq x \leq 61 \land 0 \leq y \leq 61\}$$
$$(x, y) \rightarrow \text{conv1}(x, y)$$
$$(x, y) \rightarrow [130 + 64y + x]$$

**Schedule**

**Access Map**

130 cycles
129 cycles
1 cycle
0 cycles

... (5 other ports) ...

**Unified Buffer Abstraction**

# Task for Memory Mapper: Clockwork

Optimize

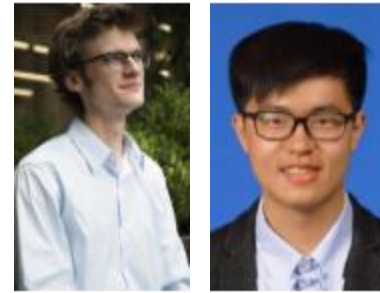- Improve schedule: reduce latency as much as possible using loop fusion and loop pipelining

Map to memories

- Occurs after schedule optimization
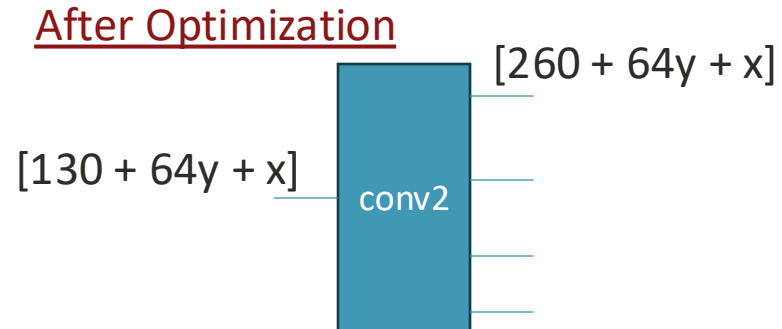- Account for hardware constraints
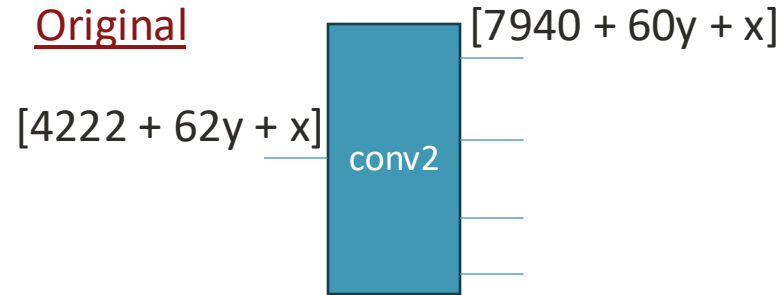
Implemented by Qiaoyi

Qiaoyi Liu. *Compiling Applications to Reconfigurable Push Memory Accelerators*. In Stanford Doctoral Dissertation, 2023.
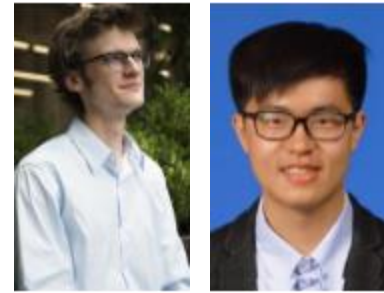
# Unified Buffer: Clockwork scheduling

Clockwork scheduling:

- Polyhedral analysis in Clockwork, created by Dillon + Qiaoyi
- <u>Purpose</u>: statically determine cycle times when all stores/loads occur
- <u>Calculation</u>: determined by dependencies, then minimize latency
- <u>Buffer size</u>: Number of cycles between store and last load for each pixel

<u>Original</u>

$[4222 + 62y + x]$ → conv2 → $[7940 + 60y + x]$

<u>After Optimization</u>

$[130 + 64y + x]$ → conv2 → $[260 + 64y + x]$

Dillon Hu, Steve Dai, and Pat Hanrahan.
*Clockwork: Resource-efficient static scheduling for multi-rate image processing applications on FPGAs*. In FCCM 2021.

# Unified Buffer: Clockwork scheduling

Clockwork scheduling:

- Polyhedral analysis in Clockwork, created by Dillon + Qiaoyi
- Purpose: statically determine cycle times when all stores/loads occur
- Calculation: determined by dependencies, then minimize latency
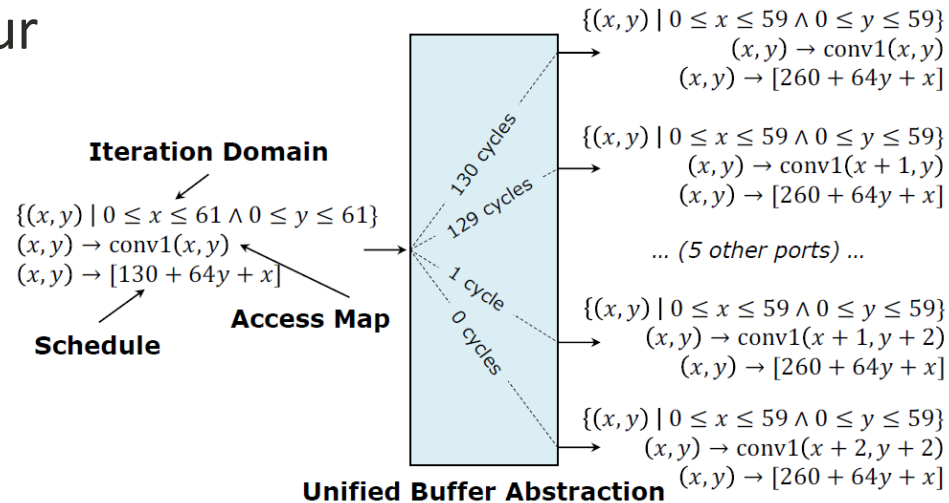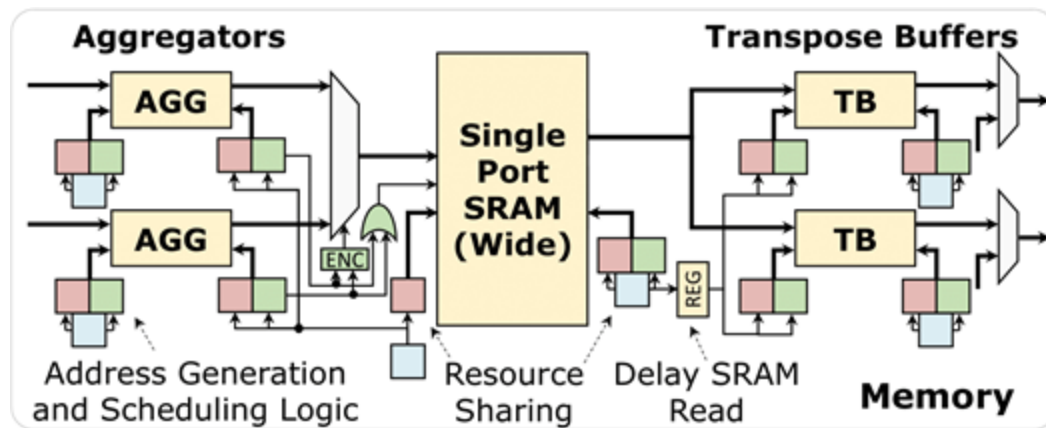- Buffer size: Number of cycles between store and last load for each pixel



**Iteration Domain**

$\{(x, y) \mid 0 \le x \le 61 \wedge 0 \le y \le 61\}$
$(x, y) \to \text{conv1}(x, y)$
$(x, y) \to [130 + 64y + x]$

**Schedule**

**Access Map**

130 cycles
129 cycles
1 cycle
0 cycles

$\{(x, y) \mid 0 \le x \le 59 \wedge 0 \le y \le 59\}$
$(x, y) \to \text{conv1}(x, y)$
$(x, y) \to [260 + 64y + x]$

$\{(x, y) \mid 0 \le x \le 59 \wedge 0 \le y \le 59\}$
$(x, y) \to \text{conv1}(x + 1, y)$
$(x, y) \to [260 + 64y + x]$

… (5 other ports) …

$\{(x, y) \mid 0 \le x \le 59 \wedge 0 \le y \le 59\}$
$(x, y) \to \text{conv1}(x + 1, y + 2)$
$(x, y) \to [260 + 64y + x]$

$\{(x, y) \mid 0 \le x \le 59 \wedge 0 \le y \le 59\}$
$(x, y) \to \text{conv1}(x + 2, y + 2)$
$(x, y) \to [260 + 64y + x]$

**Unified Buffer Abstraction**

# Unified Buffer: mapping to MEM tiles

Mapping to hardware

- Occurs after unified buffer abstraction and scheduling
- **One unified buffer != one memory tile**
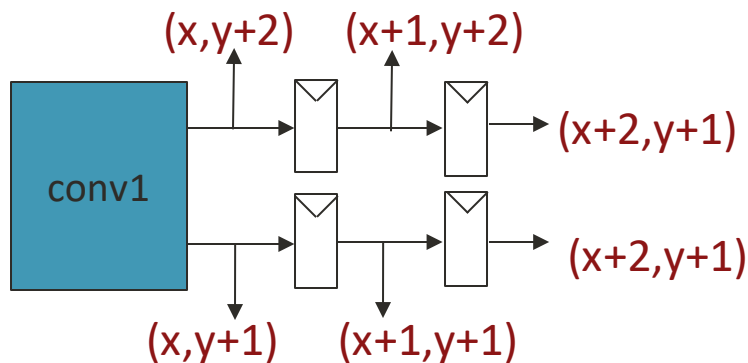
Amber CGRA Memory Tile:

- 2048 B capacity
- 2 inputs, 2 outputs
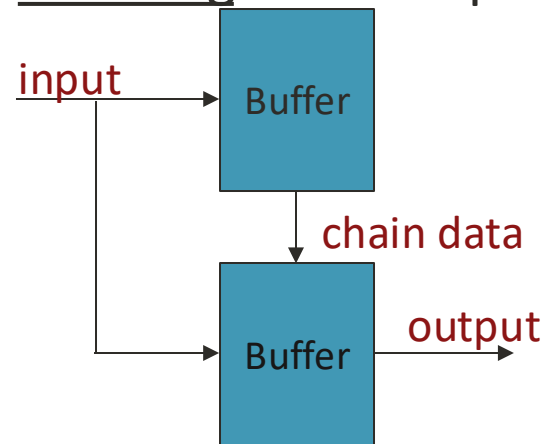- 6-level affine loopnest
- 4-wide SRAM reads/writes

# Clockwork: mapping steps

Key steps:

- <u>Wide fetch</u>: align with 4-wide single-port SRAM
- <u>Banking</u>: more IOs

- <u>Shift register optimization</u>: identify locality opportunities

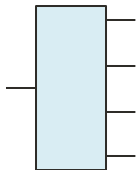- <u>Chaining</u>: more capacity

# Summary

```
app += algorithm;
app.scheduling();
```

Used **Halide** to specify an application using an algorithm, and applied a CGRA schedule to create an efficient application execution on hardware.

**Clockwork** extracted Unified Buffers to enable us to map to the memories on the CGRA.

Thank You

Thank You

# Thank You

# Demo notes

Commands

- aha map --apps apps/cascade

- Generated Log
  - Full application (CPU enclosed around accelerator = _hls_target.hw_output)
  - HalideIR printout of cascade for CGRA
  - Runs CPU and clockwork versions and compares generated output images
  - Runs clockwork memory compilation (with logging in mem_cout)
  - Runs MetaMapper

# Demo notes

Important files

- Halide generator: algorithm+schedule
    - /aha/Halide-to-Hardware/apps/hardware_benchmarks/apps/cascade/cascade_generator.cpp
- Generated files from compiler
    - /aha/Halide-to-Hardware/apps/hardware_benchmarks/apps/cascade/bin/*
- Input to Clockworks
    - /aha/Halide-to-Hardware/apps/hardware_benchmarks/apps/cascade/bin/cascade_memory.cpp
    - /aha/Halide-to-Hardware/apps/hardware_benchmarks/apps/cascade/bin/cascade_compute.json
- Output of Clockwork and MetaMapper: mapped design with memories and compute
    - /aha/Halide-to-Hardware/apps/hardware_benchmarks/apps/cascade/bin/design_top.json