

**Lecture 4:**

# **Efficiently Scheduling Image Processing Pipelines (in Halide)**

---

**Visual Computing Systems  
Stanford CS348K, Spring 2026**

# Today's themes

- **Techniques for efficiently mapping image processing applications (like those we've discussed in the past two classes) to multi-core CPUs and GPUs**
- **The design of programming abstractions that facilitate efficient image processing applications**
- **And we'll finish off our Frankencamera discussion from last time**

# Finishing off Frankencamera discussion

# Summary: F-cam scope

- **F-cam provides a set of abstractions that allow for manipulating configurable camera components**
  - **Timeline-based specification of actions (“do this at this time”)**
  - **Feed-forward system: no feedback loops**
  - **F-cam makes the frame-acquisition capabilities/operations of the camera programmable by software**
  
- **F-cam architecture performs image processing as part of outputting frames, but...**
  - **This functionality as presented by the architecture is not programmable**
  - **Hence, F-cam does not provide an image processing language**
  - **Other than work performed by the image processing stage, F-cam applications perform their own image processing (e.g., on smartphone/camera’s CPU or GPU resources) based on the frames and frame metadata they receive from the F-cam architectre**

# Reminder: key aspect in the design of any system

## Choosing the “right” representations for the job

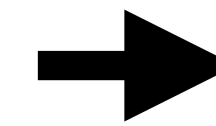
- **Good representations are productive to use:**
  - Embody the natural way of thinking about a problem
- **Good representations enable the system to provide **useful services**:**
  - Validating/providing certain guarantees (correctness, resource bounds, conversion of quantities, type checking)
  - Performance optimizations (parallelization, vectorization, use of specialized hardware)
  - Implementations of common, difficult-to-implement functionality (complex array indexing code, texture mapping in 3D graphics, auto-differentiation, etc.)

# Reminder: C++ code for a 3x3 “box blur”

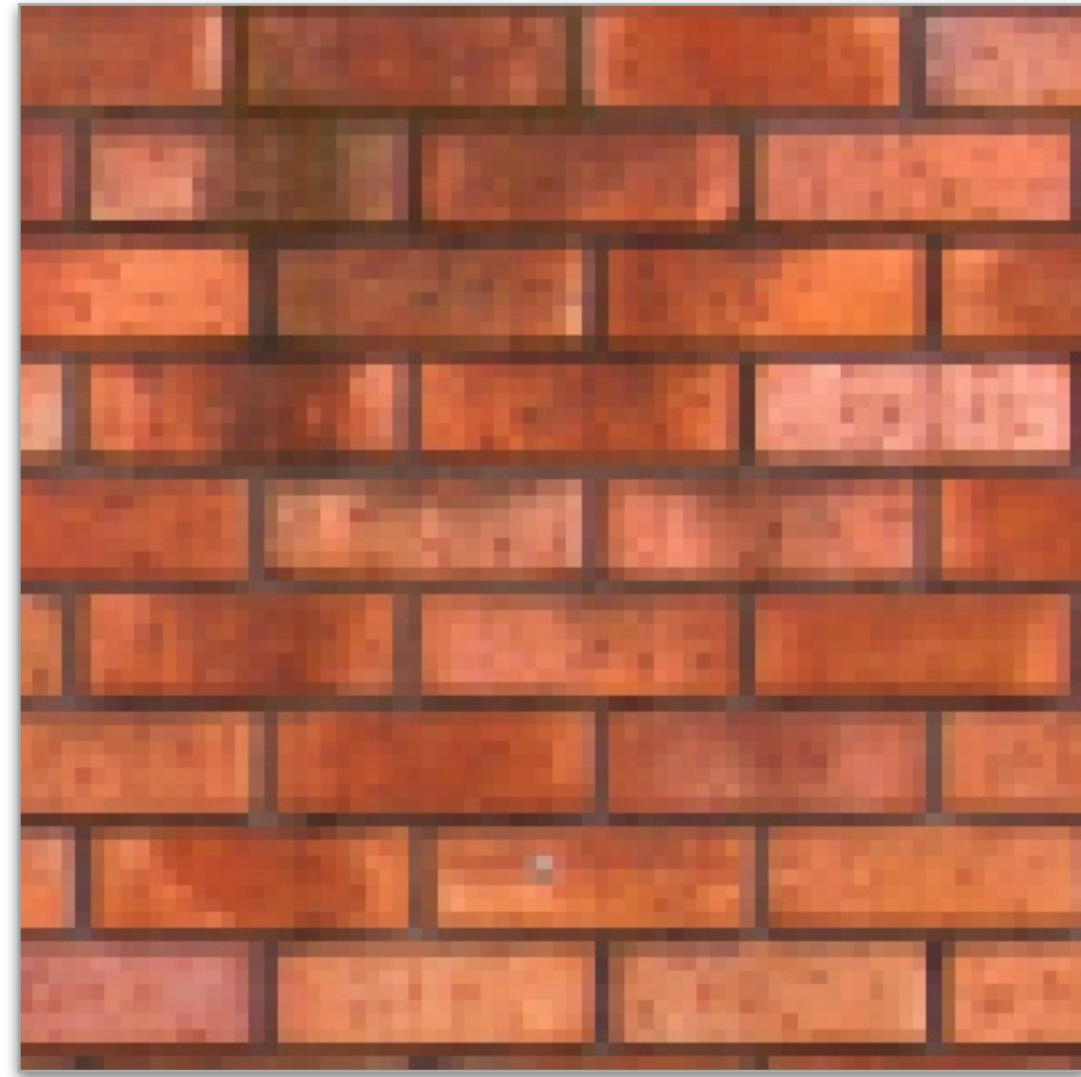
```
int WIDTH = 1024;
int HEIGHT = 1024;
float input[(WIDTH+2) * (HEIGHT+2)];
float output[WIDTH * HEIGHT];

float weights[] = {1.f/9, 1.f/9, 1.f/9,
                  1.f/9, 1.f/9, 1.f/9,
                  1.f/9, 1.f/9, 1.f/9};

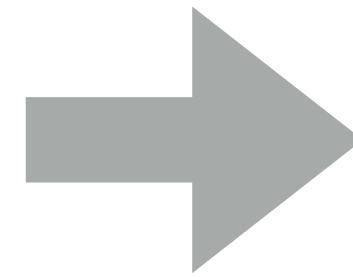
for (int j=0; j<HEIGHT; j++) {
  for (int i=0; i<WIDTH; i++) {
    float tmp = 0.f;
    for (int jj=0; jj<3; jj++)
      for (int ii=0; ii<3; ii++)
        tmp += input[(j+jj)*(WIDTH+2) + (i+ii)] * weights[jj*3 + ii];
    output[j*WIDTH + i] = tmp;
  }
}
```



# Consider a new task: sharpening an image



**Input**



**Output**

**Question: imagine you were asked to design a system for executing sharpen as efficiently as possible on a variety of parallel processors (CPUs, GPUs, etc.)**

**What would the interface to your system be?**

# Four different representations of sharpen

```
Image input;  
Image output = sharpen(input);
```

1

$$F = \begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$$

2

```
Image input;  
Image output = convolve(input, F);
```

```
Image input;  
Image output;  
output[i][j]
```

3

```
    = F[0][0] * input[i-1][j-1] +  
      F[0][1] * input[i-1][j]   +  
      F[0][2] * input[i-1][j+1] +  
      F[1][0] * input[i][j-1]   +  
      F[1][1] * input[i][j]     +
```

...

```
float input[(WIDTH+2) * (HEIGHT+2)];  
float output[WIDTH * HEIGHT];
```

4

```
float weights[] = {0., -1., 0.,  
                  -1., 5, -1.,  
                  0., -1., 0.};
```

```
for (int j=0; j<HEIGHT; j++) {  
    for (int i=0; i<WIDTH; i++) {  
        float tmp = 0.f;  
        for (int jj=0; jj<3; jj++)  
            for (int ii=0; ii<3; ii++)  
                tmp += input[(j+jj)*(WIDTH+2) + (i+ii)]  
                    * weights[jj*3 + ii];  
        output[j*WIDTH + i] = tmp;  
    }  
}
```

# Diversity of tasks: image processing tasks from previous lectures

## Sobel Edge Detection

$$G_x = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} * I$$

$$G_y = \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix} * I$$

$$G = \sqrt{G_x^2 + G_y^2}$$

## 3x3 Gaussian blur

$$F = \begin{bmatrix} .075 & .124 & .075 \\ .124 & .204 & .124 \\ .075 & .124 & .075 \end{bmatrix}$$

## 2x2 downsample (via averaging)

$$\text{output}[x][y] = (\text{input}[2x][2y] + \text{input}[2x+1][2y] + \text{input}[2x][2y+1] + \text{input}[2x+1][2y+1]) / 4.f;$$

## Gamma Correction

$$\text{output}[x][y] = \text{pow}(\text{input}[x][y], 0.5f);$$

## LUT-based correction

$$\text{output}[x][y] = \text{lookup\_table}[\text{input}[x][y]];$$

## Local Pixel Clamp

```
float f(image input) {
    float min_value = min( min(input[x-1][y], input[x+1][y]),
                          min(input[x][y-1], input[x][y+1]) );
    float max_value = max( max(input[x-1][y], input[x+1][y]),
                          max(input[x][y-1], input[x][y+1]) );
    output[x][y] = clamp(min_value, max_value, input[x][y]);
    output[x][y] = f(input);
}
```

## Histogram

```
bin[input[x][y]]++;
```

# Image processing workload characteristics

- **Structure: sequences (more precisely: DAGs) of operations on images**
- **Natural to think about algorithms in terms of their local, per-pixel behavior: e.g., output at pixel  $(x,y)$  is function of input image pixels in the neighborhood around  $(x,y)$**
- **Common case: computing value of output pixel  $(x,y)$  depends on access to a bounded local “window” of input image pixels around  $(x,y)$ ... (e.g. convolution, but also true of median filter, bilateral filter, etc.)**
- **Some algorithms require data-dependent data access (e.g., data-dependent access to lookup tables)**
- **Upsampling/downsampling (e.g., to create image pyramids)**
- **Computations that reduce information across many pixels (e.g., computing maximum brightness pixel in an image, building a histogram)**
- **FFTs on small patches of an image (to convert from pixel domain to frequency domain)**

# Halide language for image processing

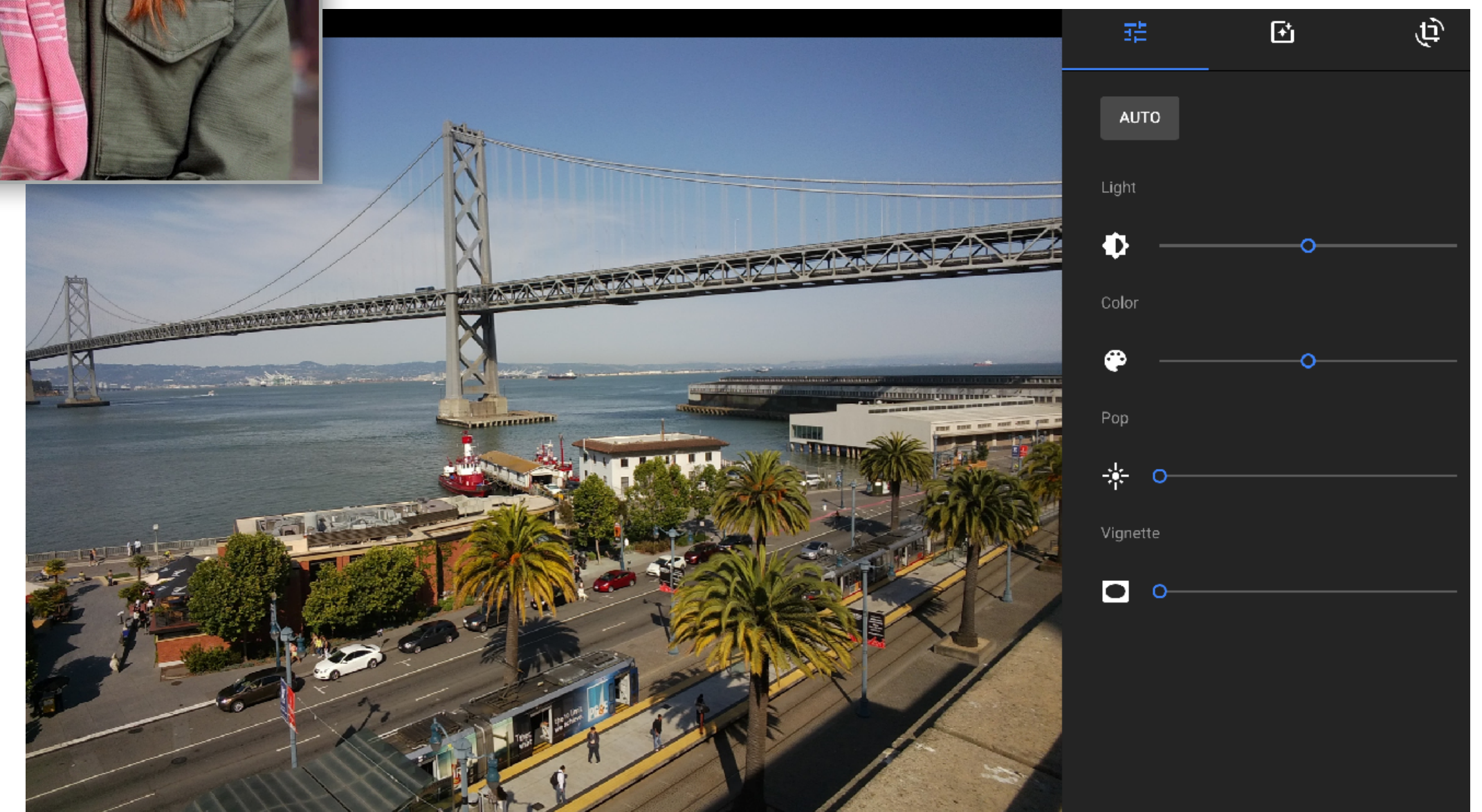
[Ragan-Kelley / Adams 2012]

# Halide goals

- **Expressive: facilitate intuitive expression of a broad class of image processing applications**
  - **e.g., all the components of a modern camera RAW pipeline**
  
- **High performance: want to generate code that efficiently utilizes the multi-core and SIMD processing resources of modern CPUs and GPUs, and is memory bandwidth efficient**

# Halide used in practice

- Halide used to implement camera processing pipelines on Google phones
  - HDR+, aspects of portrait mode, etc...
- Industry usage at Instagram, Adobe, etc.



# C++ code for a 3x3 “box blur”

```
float input[(WIDTH+2) * (HEIGHT+2)];  
float output[WIDTH * HEIGHT];
```

```
float weights[] = {1./9, 1./9, 1./9,  
                  1./9, 1./9, 1./9,  
                  1./9, 1./9, 1./9};
```

```
for (int j=0; j<HEIGHT; j++) {  
    for (int i=0; i<WIDTH; i++) {  
        float tmp = 0.f;  
        for (int jj=0; jj<3; jj++)  
            for (int ii=0; ii<3; ii++)  
                tmp += input[(j+jj)*(WIDTH+2) + (i+ii)] * weights[jj*3 + ii];  
        output[j*WIDTH + i] = tmp;  
    }  
}
```

**Total work per output image =  
9 x WIDTH x HEIGHT**

**For NxN filter:  $N^2$  x WIDTH x HEIGHT**

**For now: ignore boundary pixels and assume output image is smaller than input (makes convolution loop bounds much simpler to write)**

# 3x3 box blur in Halide

```
Var x, y;  
Func blurx, out;  
Image<uint8_t> in = load_image("myimage.jpg");
```

Functions map integer coordinates to values  
(e.g., colors of corresponding pixels)



```
// expression for computing convolution result for one output pixel
```


```
out(x,y) = 1/9.f * (in(x-1,y-1) + in(x,y-1) + in(x+1,y-1) +  
                  in(x-1,y)   + in(x,y)   + in(x+1,y) +  
                  in(x-1,y+1) + in(x,y+1) + in(x+1,y+1) );
```

```
// execute pipeline on domain of size 1024x1024
```

```
Image<uint8_t> result = out.realize(1024, 1024);
```

Total work per output image =  
9 x WIDTH x HEIGHT

For NxN filter:  $N^2$  x WIDTH x HEIGHT



Value of blurx at coordinate (x,y) is given by expression  
that accesses three values of in

**Halide function:** an infinite (but discrete) set of values defined on N-D domain

**Halide expression:** a side-effect free expression that describes how to compute a function's value at a point in its domain in terms of the values of other functions.

# An example application: two-pass blur

A 2D separable filter (such as a box filter) can be evaluated via two 1D filtering operations



**Input**



**Horizontal Blur**



**Vertical Blur**

**Note: I've exaggerated the blur for illustration (the end result is actually a 30x30 blur, not 3x3)**

# Two-pass 3x3 blur in C++

```
int WIDTH = 1024;
int HEIGHT = 1024;
float input[(WIDTH+2) * (HEIGHT+2)];
float tmp_buf[WIDTH * (HEIGHT+2)];
float output[WIDTH * HEIGHT];

float weights[] = {1.f/3, 1.f/3, 1.f/3};

for (int j=0; j<(HEIGHT+2); j++)
  for (int i=0; i<WIDTH; i++) {
    float tmp = 0.f;
    for (int ii=0; ii<3; ii++)
      tmp += input[j*(WIDTH+2) + i+ii] * weights[ii];
    tmp_buf[j*WIDTH + i] = tmp;
  }

for (int j=0; j<HEIGHT; j++) {
  for (int i=0; i<WIDTH; i++) {
    float tmp = 0.f;
    for (int jj=0; jj<3; jj++)
      tmp += tmp_buf[(j+jj)*WIDTH + i] * weights[jj];
    output[j*WIDTH + i] = tmp;
  }
}
```

1D horizontal blur

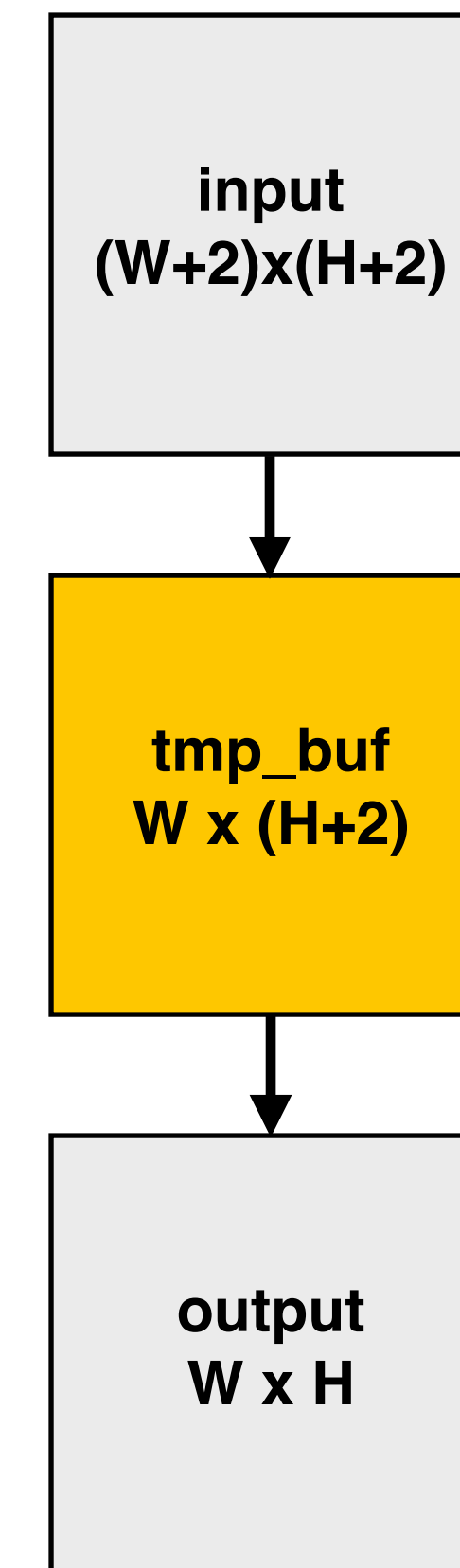
1D vertical blur

Total work per image = 6 x WIDTH x HEIGHT

For NxN filter: 2N x WIDTH x HEIGHT

WIDTH x HEIGHT extra storage

2x lower arithmetic intensity than 2D blur. Why?



# Two pass blur in Halide

Simple domain-specific language embedded in C++ for describing sequences of image processing operations

```
Var x, y;  
Func blurx, out;  
Halide::Buffer<uint8_t> in = load_image("myimage.jpg");
```


Functions map integer coordinates to values  
(e.g., colors of corresponding pixels)



```
// perform 3x3 box blur in two-passes
```

```
blurx(x,y) = 1/3.f * (in(x-1,y) + in(x,y) + in(x+1,y));  
out(x,y) = 1/3.f * (blurx(x,y-1) + blurx(x,y) + blurx(x,y+1));
```

Value of blurx at coordinate (x,y) is given by expression  
that accesses three values of in



**Halide function:** an infinite (but discrete) set of values defined on N-D domain

**Halide expression:** a side-effect free expression that describes how to compute a function's value at a point in its domain in terms of the values of other functions.

# A more complicated Halide program

Simple domain-specific language embedded in C++ for describing sequences of image processing operations

```
Var x, y;
Func blurx, blurry, bright, out;
Halide::Buffer<uint8_t> in = load_image("myimage.jpg");
Halide::Buffer<uint8_t> lookup = load_image("s_curve.jpg"); // 255-pixel 1D image

// perform 3x3 box blur in two-passes
blurx(x,y) = 1/3.f * (in(x-1,y) + in(x,y) + in(x+1,y));
blurry(x,y) = 1/3.f * (blurx(x,y-1) + blurx(x,y) + blurx(x,y+1));

// brighten blurred result by 25%, then clamp
bright(x,y) = min(blurry(x,y) * 1.25f, 255);

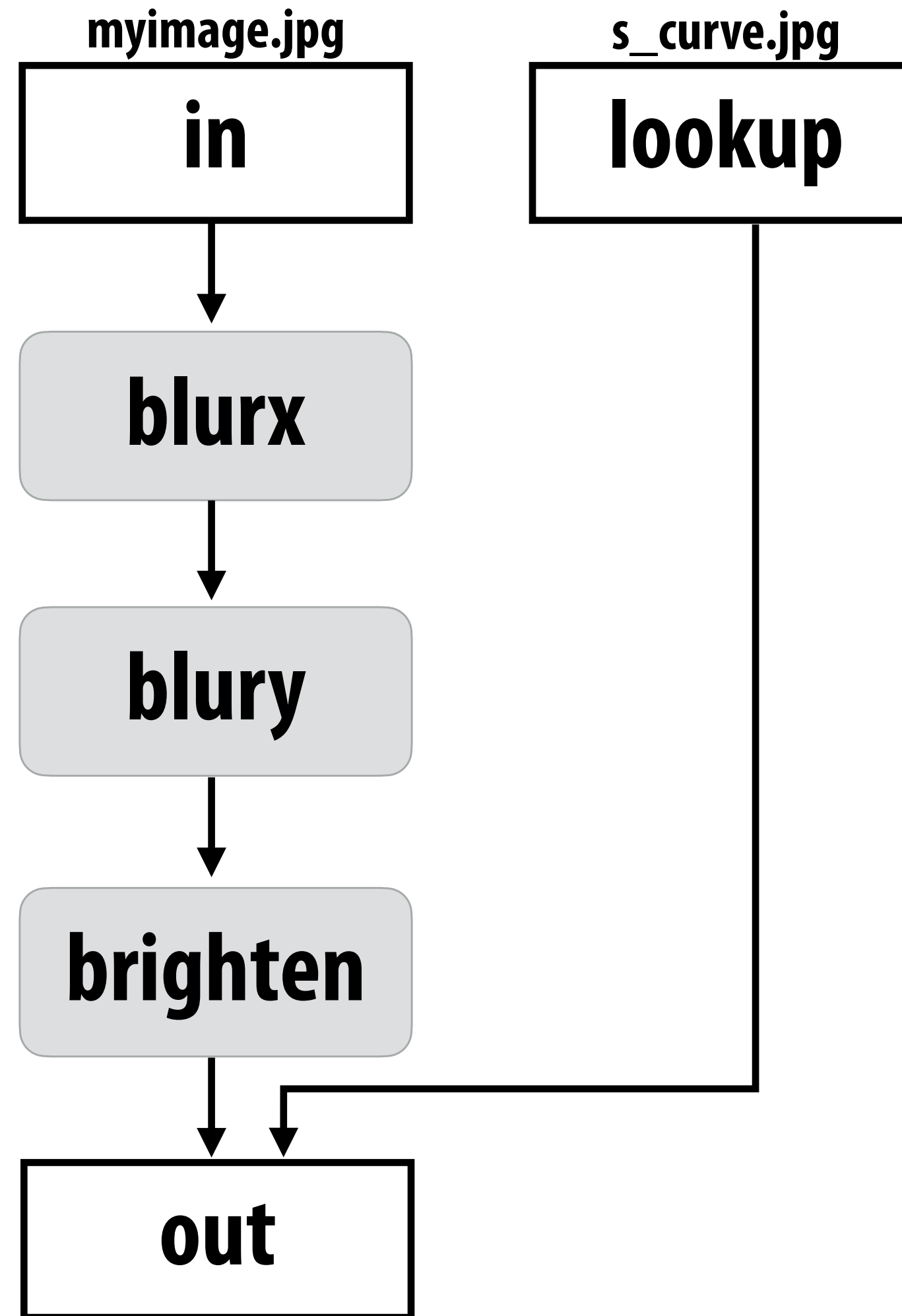
// access lookup table to contrast enhance
out(x,y) = lookup(bright(x,y));

// execute pipeline to materialize values of out in range (0:800,0:600)
Halide::Buffer<uint8_t> result = out.realize(800, 600);
```

Functions map integer coordinates to values (e.g., colors of corresponding pixels)

Value of blurx at coordinate (x,y) is given by expression accessing three values of in

# Image processing as a DAG



# Image processing pipelines feature complex DAGs of functions

Benchmark	Number of Halide functions
Two-pass blur	2
Unsharp mask	9
Harris Corner detection	13
Camera RAW processing	30
Non-local means denoising	13
Max-brightness filter	9
Multi-scale interpolation	52
Local-laplacian filter	103
Synthetic depth-of-field	74
Bilateral filter	8
Histogram equalization	7
VGG-16 deep network eval	64

**Real-world production applications may features hundreds to thousands of functions!**

**Google HDR+ pipeline: over 2000 Halide functions.**

# Key aspects of representation

## ■ Intuitive expression:

- Adopts local “point wise” view of expressing algorithms
- Halide language is declarative. It does not define order of iteration over elements in a domain, or even what values in domain are stored!
  - **It only defines what operations are needed to compute these values.**
  - **Iteration over domain points is implicit (no explicit loops)**

```
Var x, y;  
Func blurx, out;  
Halide::Buffer<uint8_t> in = load_image("myimage.jpg");  
  
// perform 3x3 box blur in two-passes  
blurx(x,y) = 1/3.f * (in(x-1,y)    + in(x,y)      + in(x+1,y));  
out(x,y) = 1/3.f * (blurx(x,y-1) + blurx(x,y) + blurx(x,y+1));  
  
// execute pipeline on domain of size 800x600  
Halide::Buffer<uint8_t> result = out.realize(800, 600);
```

# **Efficiently executing Halide programs**

## **(The interesting part!)**

# Two-pass 3x3 blur in C++

```
int WIDTH = 1024;
int HEIGHT = 1024;
float input[(WIDTH+2) * (HEIGHT+2)];
float tmp_buf[WIDTH * (HEIGHT+2)];
float output[WIDTH * HEIGHT];

float weights[] = {1.f/3, 1.f/3, 1.f/3};

for (int j=0; j<(HEIGHT+2); j++)
  for (int i=0; i<WIDTH; i++) {
    float tmp = 0.f;
    for (int ii=0; ii<3; ii++)
      tmp += input[j*(WIDTH+2) + i+ii] * weights[ii];
    tmp_buf[j*WIDTH + i] = tmp;
  }

for (int j=0; j<HEIGHT; j++) {
  for (int i=0; i<WIDTH; i++) {
    float tmp = 0.f;
    for (int jj=0; jj<3; jj++)
      tmp += tmp_buf[(j+jj)*WIDTH + i] * weights[jj];
    output[j*WIDTH + i] = tmp;
  }
}
```

1D horizontal blur

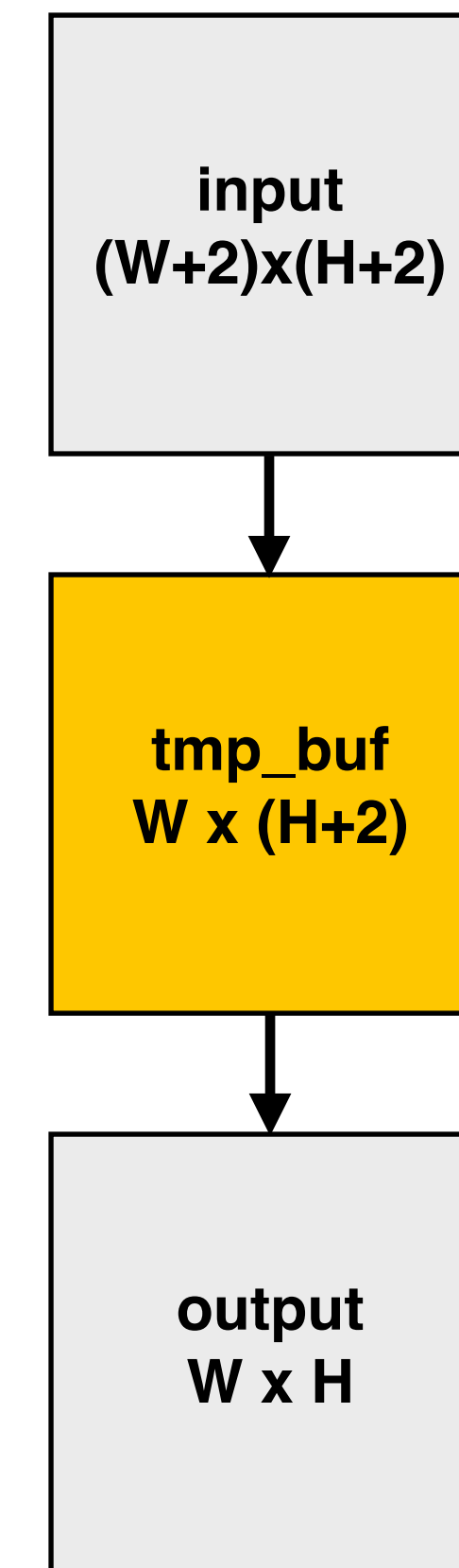
1D vertical blur

Total work per image =  $6 \times \text{WIDTH} \times \text{HEIGHT}$

For  $N \times N$  filter:  $2N \times \text{WIDTH} \times \text{HEIGHT}$

$\text{WIDTH} \times \text{HEIGHT}$  extra storage

2x lower arithmetic intensity than 2D blur. Why?



# Two-pass image blur, “chunked” (version 1)

```
int WIDTH = 1024;
int HEIGHT = 1024;
float input[(WIDTH+2) * (HEIGHT+2)];
float tmp_buf[WIDTH * 3];
float output[WIDTH * HEIGHT];

float weights[] = {1.f/3, 1.f/3, 1.f/3};

for (int j=0; j<HEIGHT; j++) {
    for (int j2=0; j2<3; j2++)
        for (int i=0; i<WIDTH; i++) {
            float tmp = 0.f;
            for (int ii=0; ii<3; ii++)
                tmp += input[(j+j2)*(WIDTH+2) + i+ii] * weights[ii];
            tmp_buf[j2*WIDTH + i] = tmp;
        }
    for (int i=0; i<WIDTH; i++) {
        float tmp = 0.f;
        for (int jj=0; jj<3; jj++)
            tmp += tmp_buf[jj*WIDTH + i] * weights[jj];
        output[j*WIDTH + i] = tmp;
    }
}
```

Only 3 rows of intermediate buffer need to be allocated

Produce 3 rows of tmp\_buf (only what's needed for one row of output)

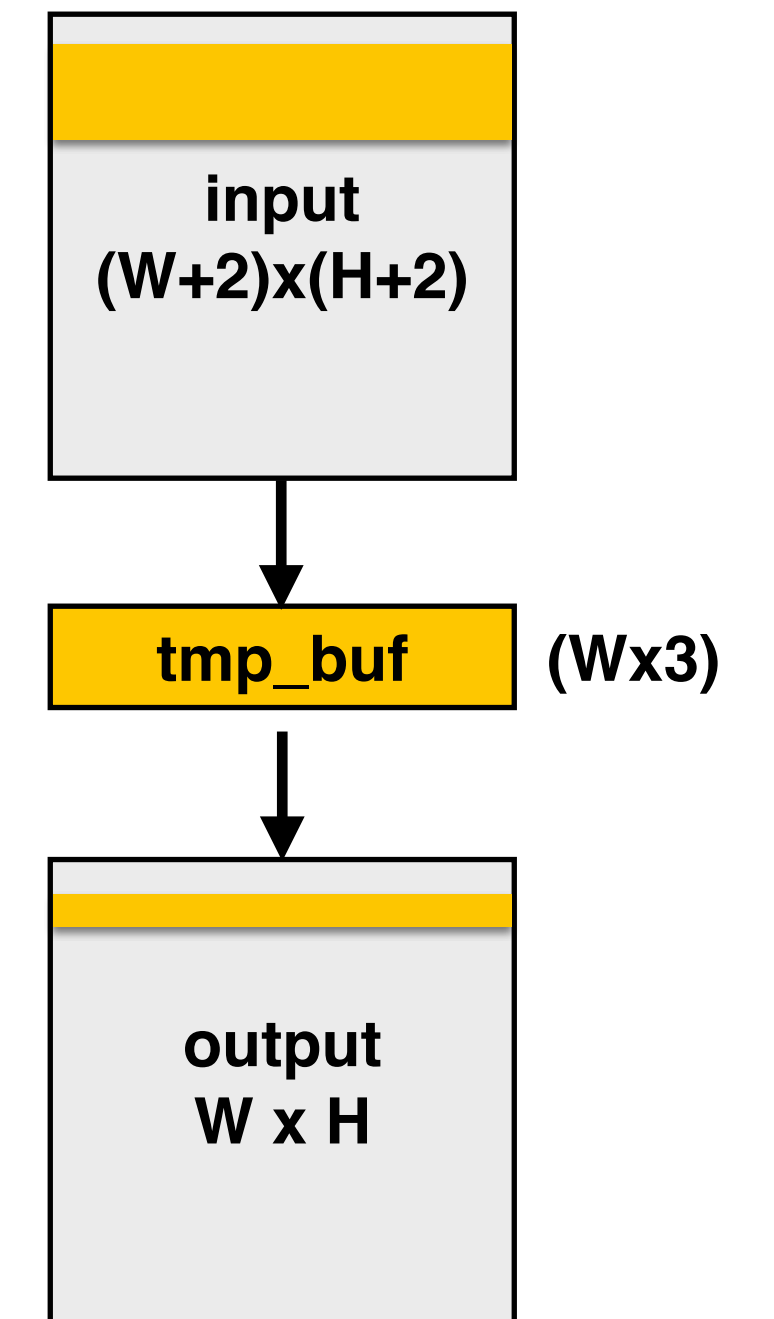
Combine them together to get one row of output

Total work per row of output:

- step 1: 3 x 3 x WIDTH work
- step 2: 3 x WIDTH work

Total work per image = 12 x WIDTH x HEIGHT ????

Loads from tmp\_buffer are cached (assuming tmp\_buffer fits in cache)



# Two-pass image blur, "chunked" (version 2)

```
int WIDTH = 1024;
int HEIGHT = 1024;
float input[(WIDTH+2) * (HEIGHT+2)];
float tmp_buf[WIDTH * (CHUNK_SIZE+2)];
float output[WIDTH * HEIGHT];

float weights[] = {1.f/3, 1.f/3, 1.f/3};

for (int j=0; j<HEIGHT; j+=CHUNK_SIZE) {
    for (int j2=0; j2<CHUNK_SIZE+2; j2++)
        for (int i=0; i<WIDTH; i++) {
            float tmp = 0.f;
            for (int ii=0; ii<3; ii++)
                tmp += input[(j+j2)*(WIDTH+2) + i+ii] * weights[ii];
            tmp_buf[j2*WIDTH + i] = tmp;
        }
    for (int j2=0; j2<CHUNK_SIZE; j2++)
        for (int i=0; i<WIDTH; i++) {
            float tmp = 0.f;
            for (int jj=0; jj<3; jj++)
                tmp += tmp_buf[(j2+jj)*WIDTH + i] * weights[jj];
            output[(j+j2)*WIDTH + i] = tmp;
        }
}
```

Sized so entire buffer fits in cache  
(capture all producer-consumer locality)

Produce enough rows of tmp\_buf to  
produce a CHUNK\_SIZE number of rows  
of output

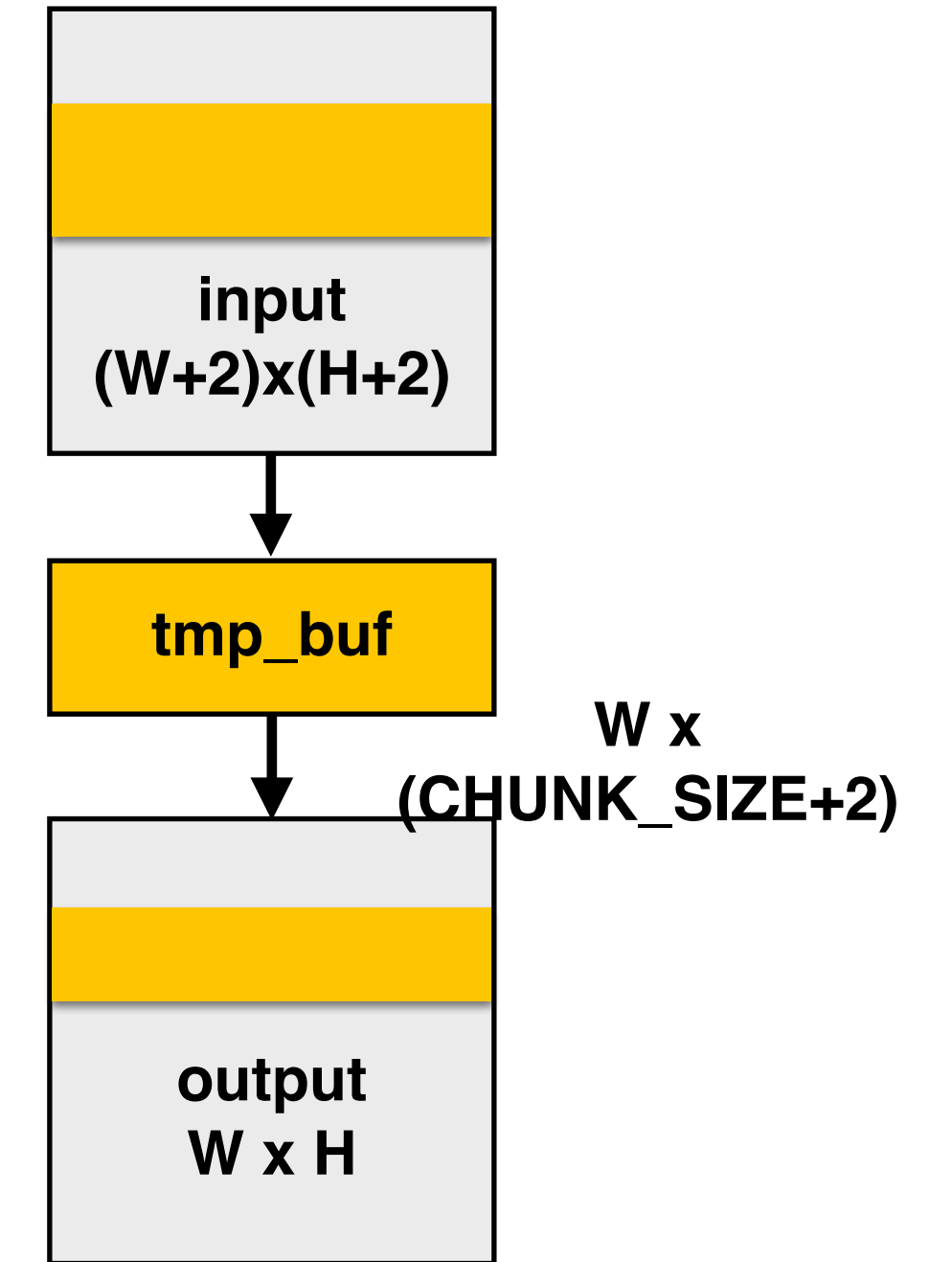
Produce CHUNK\_SIZE rows of output

Total work per chunk of output: (assume CHUNK\_SIZE = 16)

- Step 1: 18 x 3 x WIDTH work
- Step 2: 16 x 3 x WIDTH work

Total work per image:  $(34/16) \times 3 \times \text{WIDTH} \times \text{HEIGHT}$   
 $= 6.4 \times \text{WIDTH} \times \text{HEIGHT}$

Trends to ideal value of 6 x WIDTH x HEIGHT as CHUNK\_SIZE is increased!



# Still not done

- **We have not parallelized loops for multi-core execution**
- **We have not used SIMD vector instructions to execute loop bodies**
- **Other common performance optimizations: loop unrolling, etc...**

# Optimized implementation of 3x3 box blur in x86 SSE intrinsics

Good: ~10x faster on a quad-core CPU than my original two-pass code

Bad: specific to SSE (not AVX2), CPU-code only, hard to tell what is going on at all!

```
void fast_blur(const Image &in, Image &blurred) {
    _m128i one_third = _mm_set1_epi16(21846);
    #pragma omp parallel for
    for (int yTile = 0; yTile < in.height(); yTile += 32) {
        _m128i a, b, c, sum, avg;
        _m128i tmp[(256/8)*(32+2)];
        for (int xTile = 0; xTile < in.width(); xTile += 256) {
            _m128i *tmpPtr = tmp;
            for (int y = -1; y < 32+1; y++) {
                const uint16_t *inPtr = &(in(xTile, yTile+y));
                for (int x = 0; x < 256; x += 8) {
                    a = _mm_loadu_si128((__m128i*)(inPtr-1));
                    b = _mm_loadu_si128((__m128i*)(inPtr+1));
                    c = _mm_load_si128((__m128i*)(inPtr));
                    sum = _mm_add_epi16(_mm_add_epi16(a, b), c);
                    avg = _mm_mulhi_epi16(sum, one_third);
                    _mm_store_si128(tmpPtr++, avg);
                    inPtr += 8;
                }
            }
            tmpPtr = tmp;
            for (int y = 0; y < 32; y++) {
                _m128i *outPtr = (_m128i *)&(blurred(xTile, yTile+y));
                for (int x = 0; x < 256; x += 8) {
                    a = _mm_load_si128(tmpPtr+(2*256)/8);
                    b = _mm_load_si128(tmpPtr+256/8);
                    c = _mm_load_si128(tmpPtr++);
                    sum = _mm_add_epi16(_mm_add_epi16(a, b), c);
                    avg = _mm_mulhi_epi16(sum, one_third);
                    _mm_store_si128(outPtr++, avg);
                }
            }
        }
    }
}
```

Multi-core execution  
(partition image vertically)

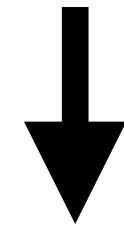
Modified iteration order:  
256x32 tiled iteration (to  
maximize cache hit rate)

use of SIMD vector  
intrinsics

two passes fused into one:  
tmp data read from cache

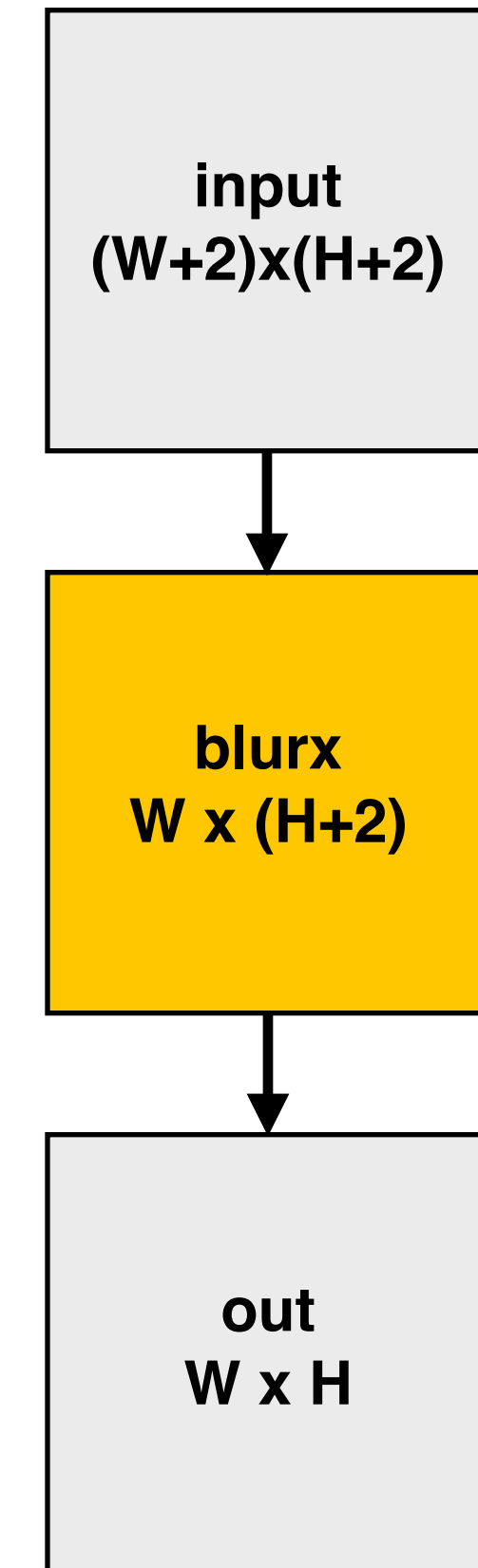
# One (serial) implementation of Halide

```
Func blurx, out;  
Var x, y, xi, yi;  
Halide::Buffer<uint8_t> in = load_image("myimage.jpg");  
  
// the "algorithm description" (declaration of what to do)  
blurx(x,y) = (in(x-1, y) + in(x,y) + in(x+1,y)) / 3.0f;  
out(x,y)   = (blurx(x,y-1) + blurx(x,y) + blurx(x,y+1)) / 3.0f;  
  
// execute pipeline on domain of size 1024x1024  
Halide::Buffer<uint8_t> result = out.realize(1024, 1024);
```



Equivalent "C-style" loop nest:

```
allocate in(1024+2, 1024+2); // (width,height)... initialize from image  
allocate blurx(1024,1024+2); // (width,height)  
allocate out(1024,1024); // (width,height)  
  
for y=0 to 1024:  
  for x=0 to 1024+2:  
    blurx(x,y) = ... compute from in  
  
for y=0 to 1024:  
  for x=0 to 1024:  
    out(x,y) = ... compute from blurx
```



# Key aspect in the design of any system:

Choosing the “right” representations for the job

- **Good representations are productive to use:**
  - Embody the natural way of thinking about a problem
- **Good representations enable the system to provide the application useful services:**
  - Validating/providing certain guarantees (correctness, resource bounds, type checking)
  - Performance (parallelization, vectorization, use of specialized hardware)

**Now the job is not expressing an image processing computation, but generating an efficient implementation of a specific Halide program.**

# A second set of representations for “scheduling”

```
Func blurx, out;  
Var x, y, xi, yi;  
Halide::Buffer<uint8_t> in = load_image("myimage.jpg");  
  
// the “algorithm description” (declaration of what to do)  
blurx(x,y) = (in(x-1, y) + in(x,y) + in(x+1,y)) / 3.0f;  
out(x,y)   = (blurx(x,y-1) + blurx(x,y) + blurx(x,y+1)) / 3.0f;
```

```
// “the schedule” (how to do it)  
out.tile(x, y, xi, yi, 256, 32).vectorize(xi,8).parallel(y);
```

```
blurx.compute_at(x).vectorize(x, 8);
```

Produce elements `blurx` on demand for each tile of output.

Vectorize the `x` (innermost) loop

When evaluating `out`, use 2D tiling order (loops named by `x, y, xi, yi`).  
Use tile size 256 x 32.

Vectorize the `xi` loop (8-wide)

Use threads to parallelize the `y` loop

“Schedule”

```
// execute pipeline on domain of size 1024x1024  
Halide::Buffer<uint8_t> result = out.realize(1024, 1024);
```

Scheduling primitives allow the programmer to specify a high-level “sketch” of how to schedule the algorithm onto a parallel machine, but leave the details of emitting the low-level platform-specific code to the Halide compiler

# Primitives for iteration order

1	2	3	4	5	6
7	8	9	10	11	12
13	14	15	16	17	18
19	20	21	22	23	24
25	26	27	28	29	30
31	32	33	34	35	36

serial y, serial x

1	7	13	19	25	31
2	8	14	20	26	32
3	9	15	21	27	33
4	10	16	22	28	34
5	11	17	23	29	35
6	12	18	24	30	36

serial x, serial y

Specify both order and how to parallelize  
(multi-thread, vectorize via SIMD instr)

1	2
3	4
5	6
7	8
9	10
11	12

serial y  
vectorized x

1	2
1	2
1	2
1	2
1	2
1	2
1	2

parallel y  
vectorized x

t0  
t1

1	2	5	6	9	10
3	4	7	8	11	12
13	14	17	18	21	22
15	16	19	20	23	24
25	26	29	30	33	34
27	28	31	32	35	36

split x into  $2x_0+x_i$ ,  
split y into  $2y_0+y_i$ ,  
serial  $y_0, x_0, y_i, x_i$

2D blocked iteration order

(In diagram, numbers indicate sequential order of processing within a thread)

# Ordering Halide loop nests

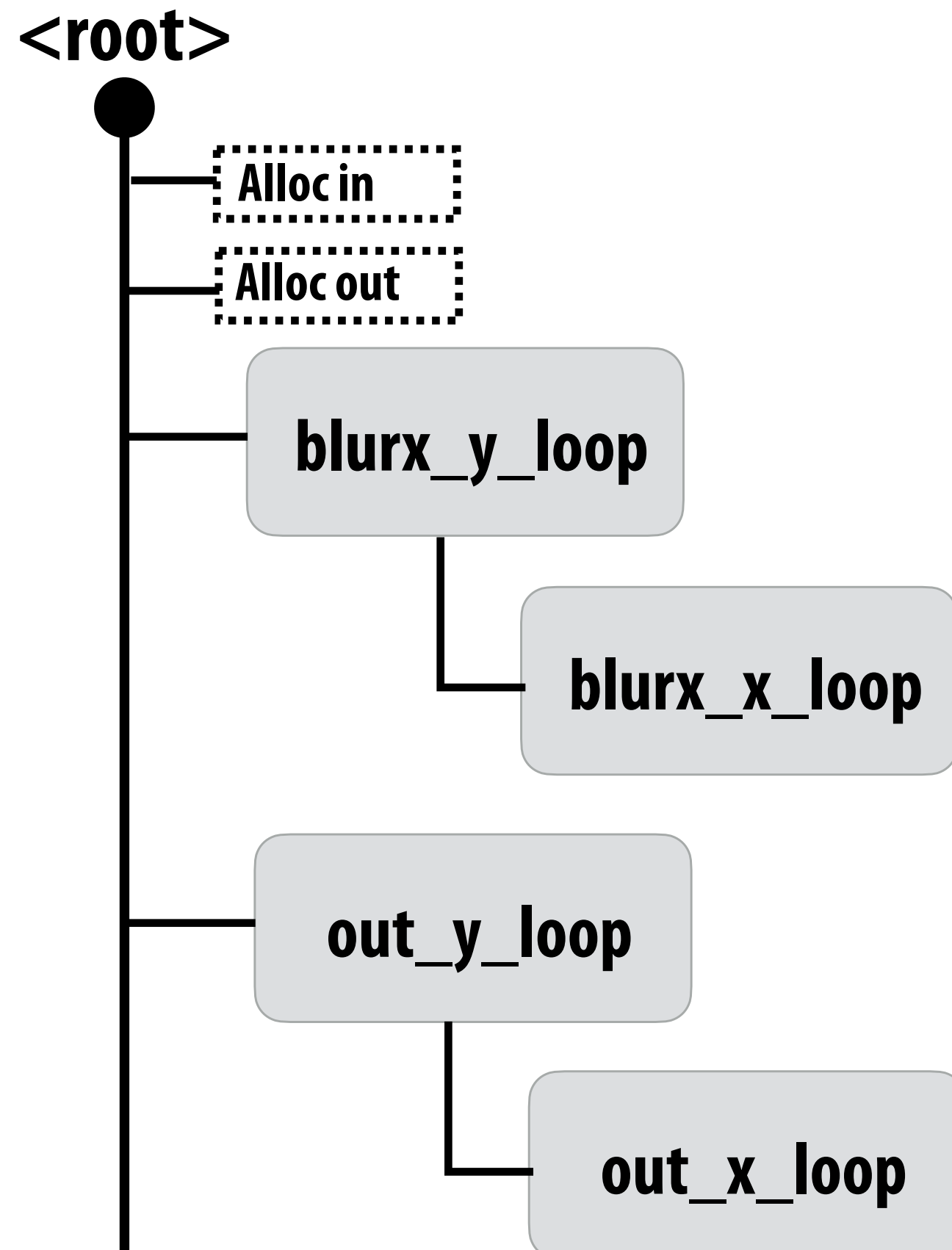
## Halide algorithm:

```
blurx(x,y) = (in(x-1, y) + in(x,y) + in(x+1,y)) / 3.0f;  
out(x,y)   = (blurx(x,y-1) + blurx(x,y) + blurx(x,y+1)) / 3.0f;  
Halide::Buffer<uint8_t> result = out.realize(1024, 1024);
```

## Halide schedule:

```
blurx.compute_root();
```

## Loop nest diagram of implementation:



## C-code equivalent:

```
allocate in(1024+2, 1024+2); // (width,height)... initialize from image  
allocate blurx(1024,1024+2); // (width,height)  
allocate out(1024,1024); // (width,height)
```

```
for y=0 to 1024:  
  for x=0 to 1024+2:  
    blurx(x,y) = ... compute from in
```

┌  
└ Loops for computing values of blurx

```
for y=0 to 1024:  
  for x=0 to 1024:  
    out(x,y) = ... compute from blurx
```

┌  
└ Loops for computing values of out

# Ordering Halide loop nests

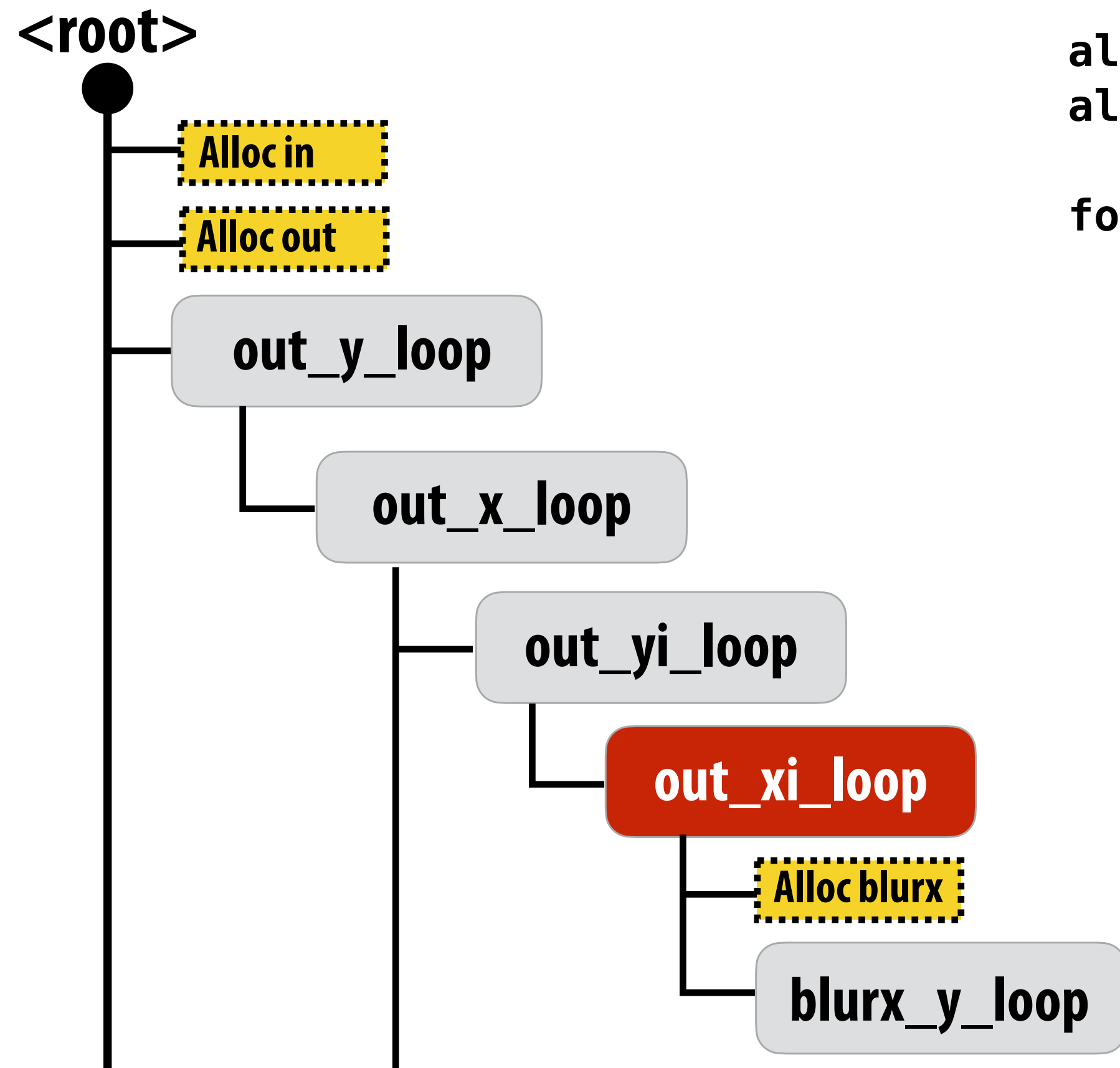
## Halide algorithm:

```
blurx(x,y) = (in(x-1, y) + in(x,y) + in(x+1,y)) / 3.0f;  
out(x,y)   = (blurx(x,y-1) + blurx(x,y) + blurx(x,y+1)) / 3.0f;  
Halide::Buffer<uint8_t> result = out.realize(1024, 1024);
```

## Halide schedule:

```
out.tile(x, y, xi, yi, 256, 32);  
blurx.compute_at(out, xi);
```

## Loop nest diagram of implementation:



## Another possible implementation:

```
allocate in(1024+2, 1024+2); // (width,height)... initialize from image  
allocate out(1024,1024);    // (width,height)
```

```
for y=0 to num_tiles_y:  
  for x=0 to num_tiles_x: Outer loops over tiles of out
```

```
  for yi=0 to 32: Inner loops for computing values of out  
    for xi=0 to 256:
```

```
      idx_x = x*256+xi;  
      idx_y = y*32+yi
```

```
      allocate blurx(1,3)
```

**Only allocate 3 elements of blurx**

```
      // compute 3 elements of blurx needed for out(idx_x, idx_y) here
```

```
      for blurx_y=0 to 3:  
        blurx(0, blurx_y) = ... // compute blurx from in
```

```
      out(idx_x, idx_y) = ... // compute out from blurx
```

# Ordering Halide loop nests

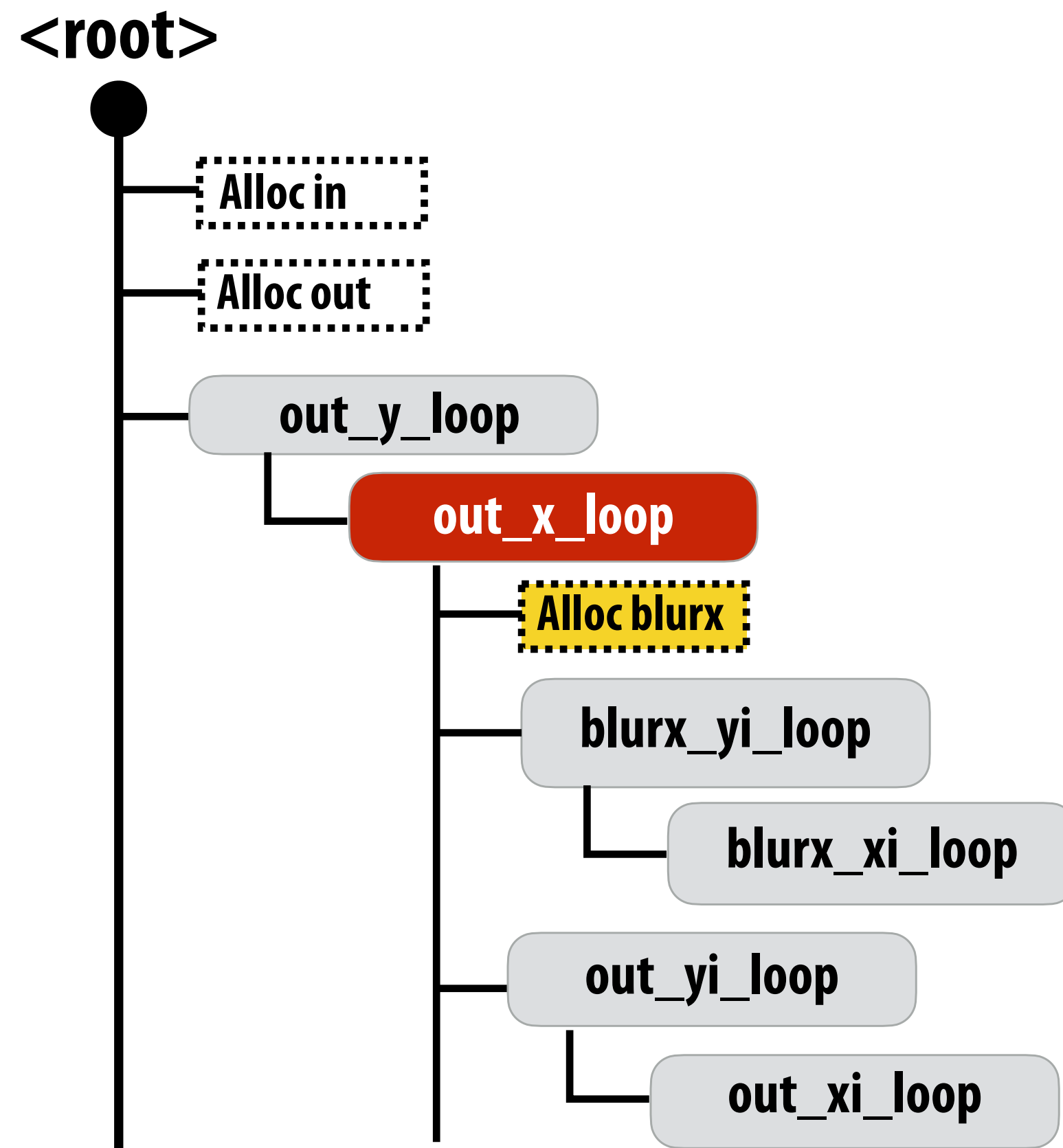
## Halide algorithm:

```
blurx(x,y) = (in(x-1, y) + in(x,y) + in(x+1,y)) / 3.0f;  
out(x,y)   = (blurx(x,y-1) + blurx(x,y) + blurx(x,y+1)) / 3.0f;  
Halide::Buffer<uint8_t> result = out.realize(1024, 1024);
```

## Halide schedule:

```
out.tile(x, y, xi, yi, 256, 32);  
blurx.compute_at(out, x);
```

## Loop nest diagram of implementation:



## C-code equivalent:

```
allocate in(1024+2, 1024+2); // (width,height)... initialize from image  
allocate out(1024,1024); // (width,height)  
  
for y=0 to num_tiles_y:  
  for x=0 to num_tiles_x: Outer loops over tiles of out  
  
    allocate blurx(256, 34) ← Only allocate a tile of blurx  
    for yi=0 to 32+2: Loops for computing values of blurx  
      for xi=0 to 256: compute blurx from in  
        blurx(xi,yi) = // compute blurx from in  
  
    for yi=0 to 32: Inner loops for computing values of out (loops over elements)  
      for xi=0 to 256:  
        idx_x = x*256+xi;  
        idx_y = y*32+yi  
        out(idx_x, idx_y) = // compute out from blurx
```

# An interesting Halide schedule

*“Compute necessary elements of blurx within out’s xi loop nest, but fill in tile-sized buffer allocated at x loop nest.”*

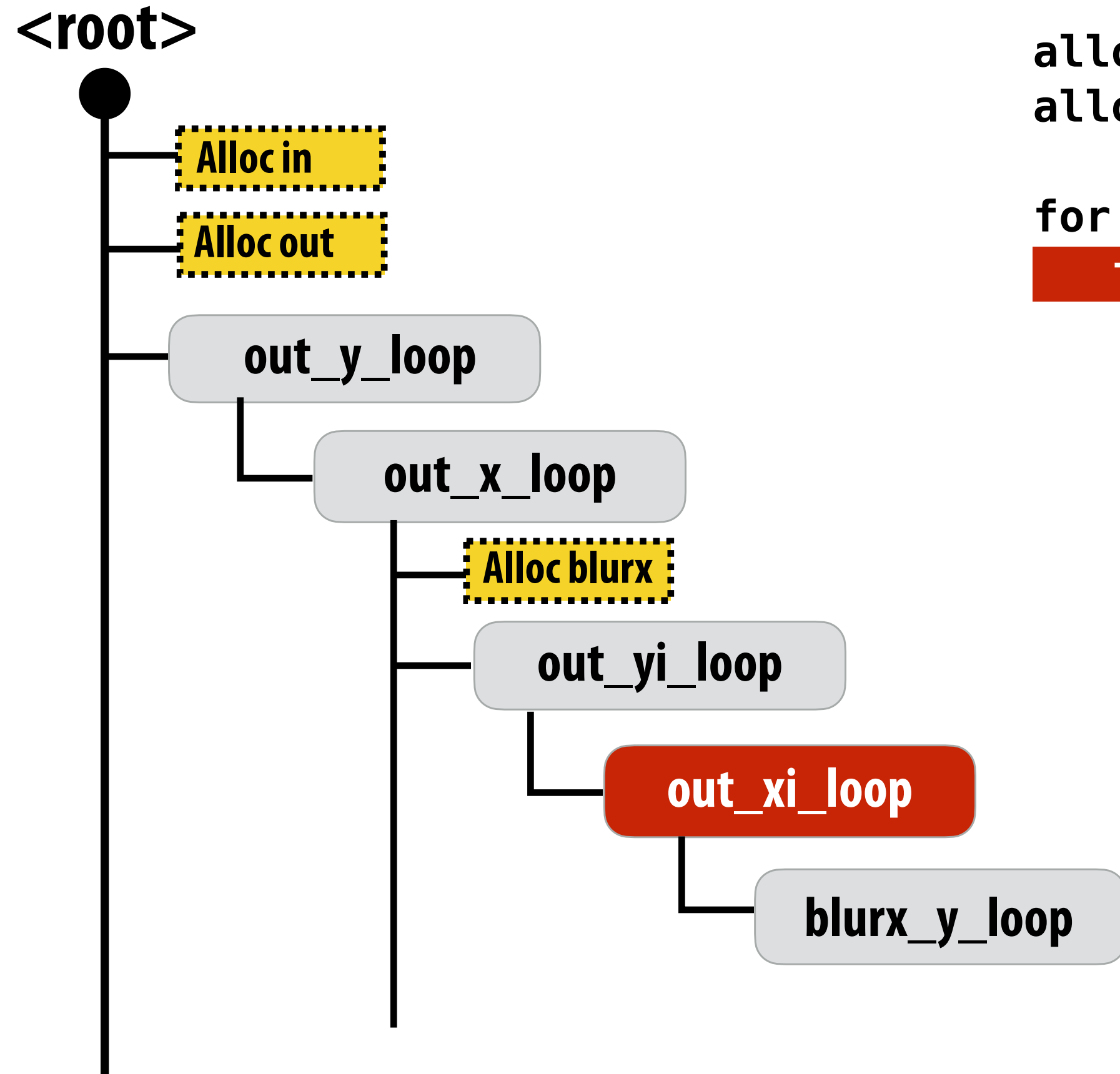
## Halide algorithm:

```
blurx(x,y) = (in(x-1, y) + in(x,y) + in(x+1,y)) / 3.0f;  
out(x,y)   = (blurx(x,y-1) + blurx(x,y) + blurx(x,y+1)) / 3.0f;  
Halide::Buffer<uint8_t> result = out.realize(1024, 1024);
```

## Halide schedule:

```
out.tile(x, y, xi, yi, 256, 32);  
blurx.store_at(out, x);  
blurx.compute_at(out, xi);
```

## Loop nest diagram of implementation:



## C-code equivalent:

```
allocate in(1024+2, 1024+2); // (width,height)... initialize from image  
allocate out(1024,1024);    // (width,height)  
  
for y=0 to num_tiles_y:  
  for x=0 to num_tiles_x:  Outer loops over tiles of out  
    allocate blurx(256, 34) Allocate a tile of blurx  
    for yi=0 to 32: Loops for computing values of out  
      for xi=0 to 256:  
        idx_x = x*256+xi;  
        idx_y = y*32+yi;  
        // compute 3 elements of blurx needed for out(idx_x, idx_y) here  
        for blurx_y=0 to 3: This recomputes values. Can compiler be smarter?  
          blurx(xi, yi + blurx_y) = ... // compute blurx from in  
        out(idx_x, idx_y) = ... // compute out from blurx
```

# “Sliding optimization” (reduces redundant computation)

*“Compute necessary elements of blurx within out’s xi loop nest, but fill in tile-sized buffer allocated at x loop nest.”*

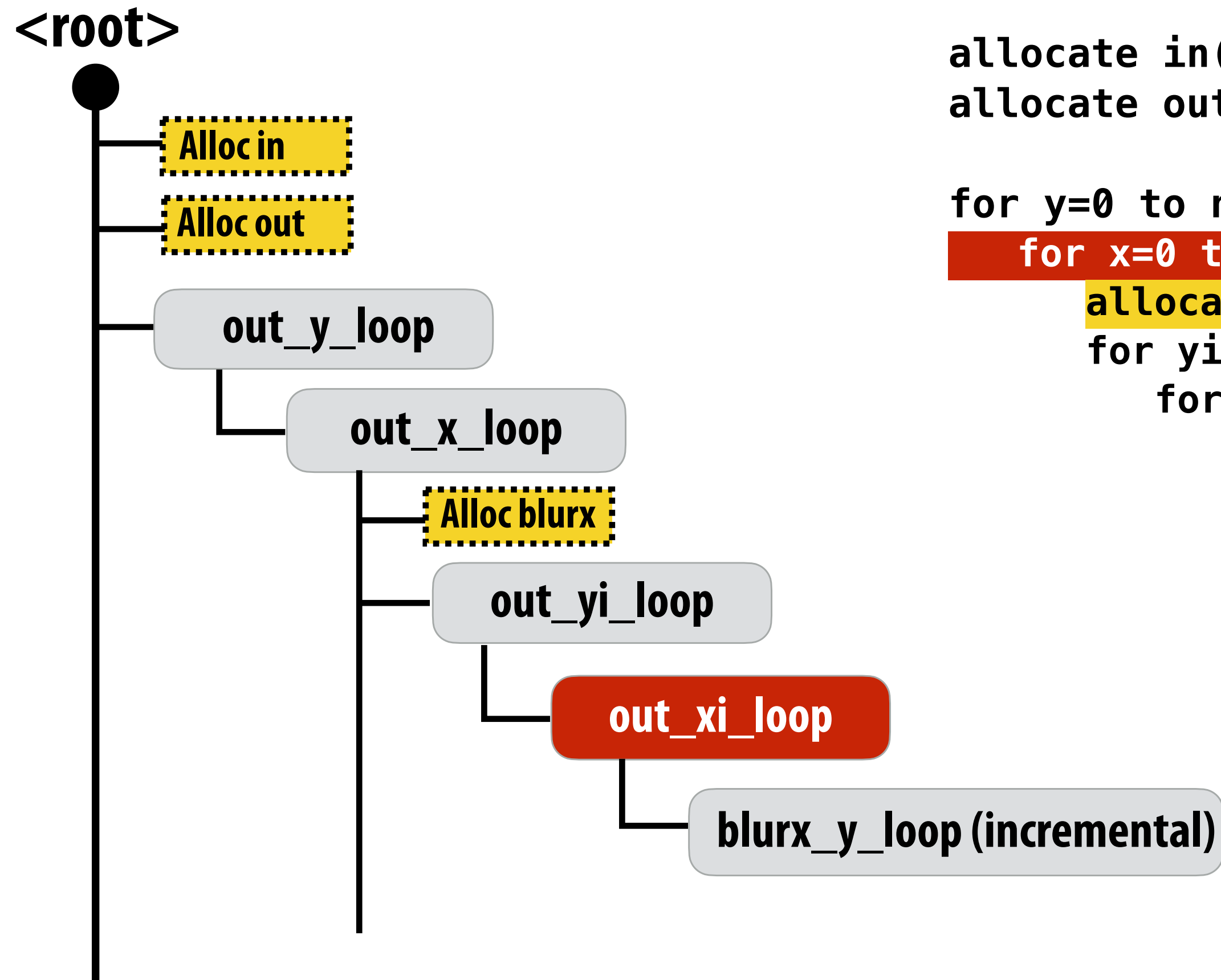
## Halide algorithm:

```
blurx(x,y) = (in(x-1, y) + in(x,y) + in(x+1,y)) / 3.0f;  
out(x,y)   = (blurx(x,y-1) + blurx(x,y) + blurx(x,y+1)) / 3.0f;  
Halide::Buffer<uint8_t> result = out.realize(1024, 1024);
```

## Halide schedule:

```
out.tile(x, y, xi, yi, 256, 32);  
blurx.store_at(out, x);  
blurx.compute_at(out, xi);
```

## Loop nest diagram of implementation:



## C-code equivalent:

```
allocate in(1024+2, 1024+2); // (width,height)... initialize from image
allocate out(1024,1024);    // (width,height)

for y=0 to num_tiles_y:
  for x=0 to num_tiles_x:
    allocate blurx(256, 34)
    for yi=0 to 32:
      for xi=0 to 256:
        idx_x = x*256+xi;
        idx_y = y*32+yi;

        if (yi=0) {
          // compute 3 elements of blurx needed for out(idx_x, idx_y) here
          for blurx_y=0 to 3:
            blurx(xi, yi + blurx_y) = ... // compute blurx from in
        } else
          blurx(xi, yi + 2) = ... // compute one additional element of blurx

        out(idx_x, idx_y) = ... // compute out from blurx
```

Outer loops over tiles of out  
Allocate a tile of blurx  
Loops for computing values of out

Steady state: only one new element of blurx needs to be computed per output

# “Folding optimization” (reduces intermediate storage)

“Compute necessary elements of blurx within out’s xi loop nest, but fill in tile-sized buffer allocated at x loop nest.”

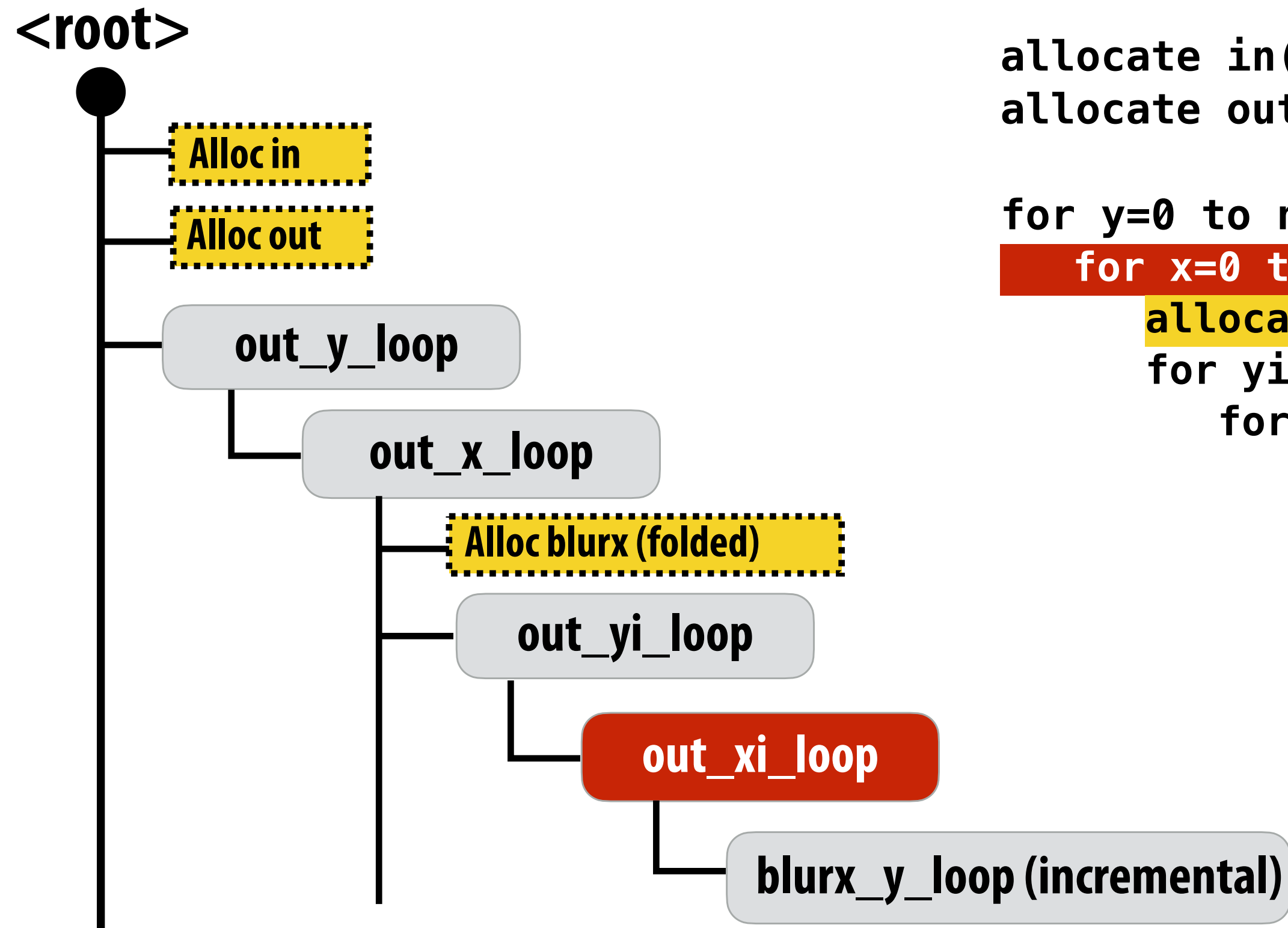
## Halide algorithm:

```
blurx(x,y) = (in(x-1, y) + in(x,y) + in(x+1,y)) / 3.0f;  
out(x,y)   = (blurx(x,y-1) + blurx(x,y) + blurx(x,y+1)) / 3.0f;  
Halide::Buffer<uint8_t> result = out.realize(1024, 1024);
```

## Halide schedule:

```
out.tile(x, y, xi, yi, 256, 32);  
blurx.store_at(out, x);  
blurx.compute_at(out, xi);
```

## Loop nest diagram of implementation:



## C-code equivalent:

```
allocate in(1024+2, 1024+2); // (width,height)... initialize from image  
allocate out(1024,1024); // (width,height)  
  
for y=0 to num_tiles_y:  
  for x=0 to num_tiles_x: Outer loops over tiles of out  
    allocate blurx(256, 3) Allocate only 3 rows of blurx tile  
    for yi=0 to 32: Loops for computing values of out  
      for xi=0 to 256:  
        idx_x = x*256+xi;  
        idx_y = y*32+yi;  
  
        if (yi=0) {  
          // compute 3 elements of blurx needed for out(idx_x, idx_y) here  
          for blurx_y=0 to 3:  
            blurx(xi, yi + blurx_y) = ... // compute blurx from in  
          } else  
            blurx(xi, (yi + 2) % 3) = ... // compute one additional element of  
            // blurx  
        }  
        out(idx_x, idx_y) = ... // compute out from blurx
```

Accesses to blurx modified to access appropriate row of circular buffer:

# Summary of scheduling the 3x3 box blur

```
// the "algorithm description" (declaration of what to do)
blurx(x,y) = (in(x-1, y) + in(x,y) + in(x+1,y)) / 3.0f;
out(x,y)   = (blurx(x,y-1) + blurx(x,y) + blurx(x,y+1)) / 3.0f;
```

```
// "the schedule" (how to do it)
out.tile(x, y, xi, yi, 256, 32).vectorize(xi,8).parallel(y);
blurx.compute_at(out, x).vectorize(x, 8);
```

## Equivalent parallel loop nest:

---

```
for y=0 to num_tiles_y: // iters of this loop are parallelized using threads
  for x=0 to num_tiles_x:
    allocate blurx(256, 34)
    for yi=0 to 32+2:
      for xi=0 to 256+2 by 8:
        blurx(xi,yi) = ... // compute blurx from in using 8-wide
                           // SIMD instructions here
                           // compiler generates boundary conditions
                           // since 256+2 isn't evenly divided by 8
    for yi=0 to 32:
      for xi=0 to 256 by 8:
        idx_x = x*256+xi;
        idx_y = y*32+yi
        out(idx_x, idx_y) = ... // compute out from blurx using 8-wide
                               // SIMD instructions here
```

# What is the philosophy of Halide

- **Programmer** is responsible for describing an image processing algorithm
- **Programmer** has knowledge to schedule application efficiently on machine (but it's slow and tedious), so give programmer another language to express their high-level scheduling decisions
  - Loop structure of code
  - Unrolling / vectorization / multi-core parallelization
- **The system** (Halide compiler) is not smart, it provides the service of mechanically carrying out the nitty gritty details of implementing the schedule using mechanisms available on the target machine (pthreads, AVX intrinsics, CUDA code, etc.)
  - There are deviations from this philosophy in Halide? What are they?

# Constraints on language

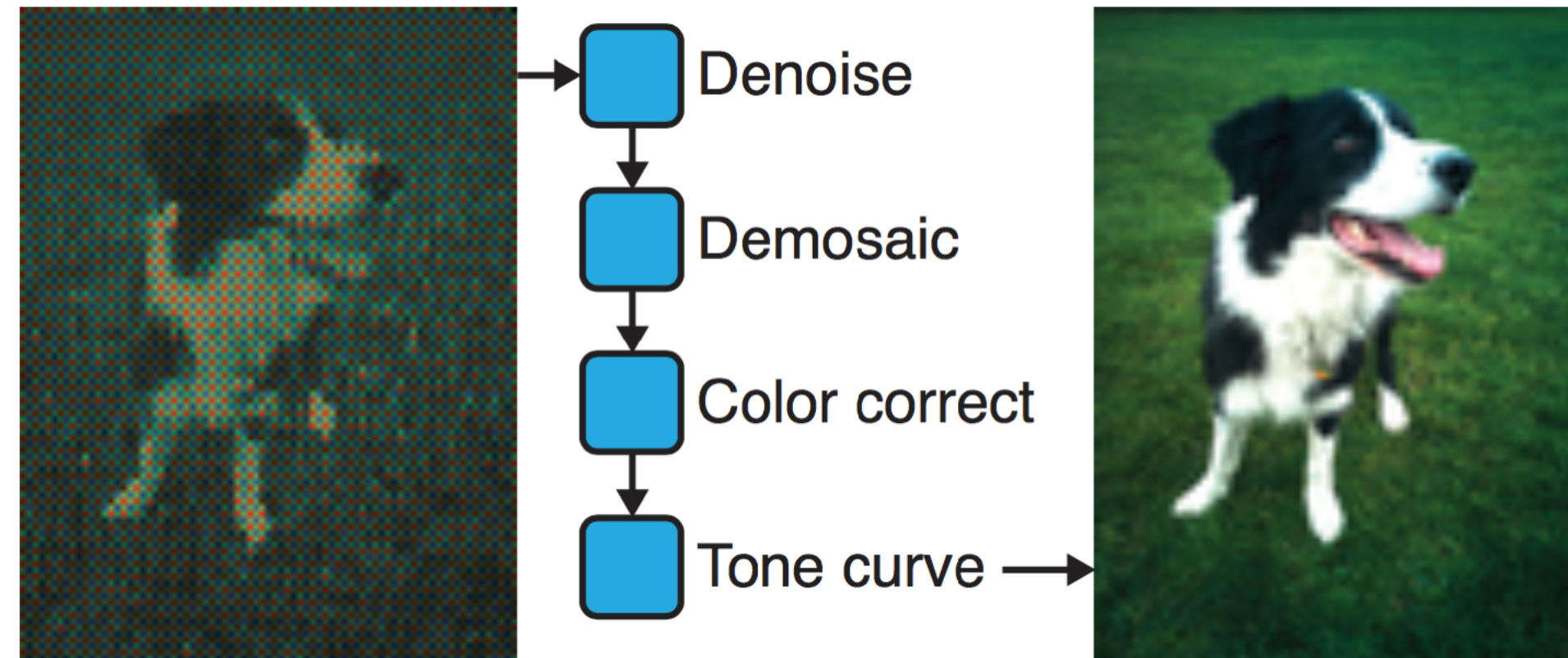
(to enable compiler to provide desired services)

- **Application domain scope: computation on regular N-D domains**
- **Only feed-forward pipelines (includes special support for reductions and fixed recursion depth)**
- **All dependencies inferable by compiler**

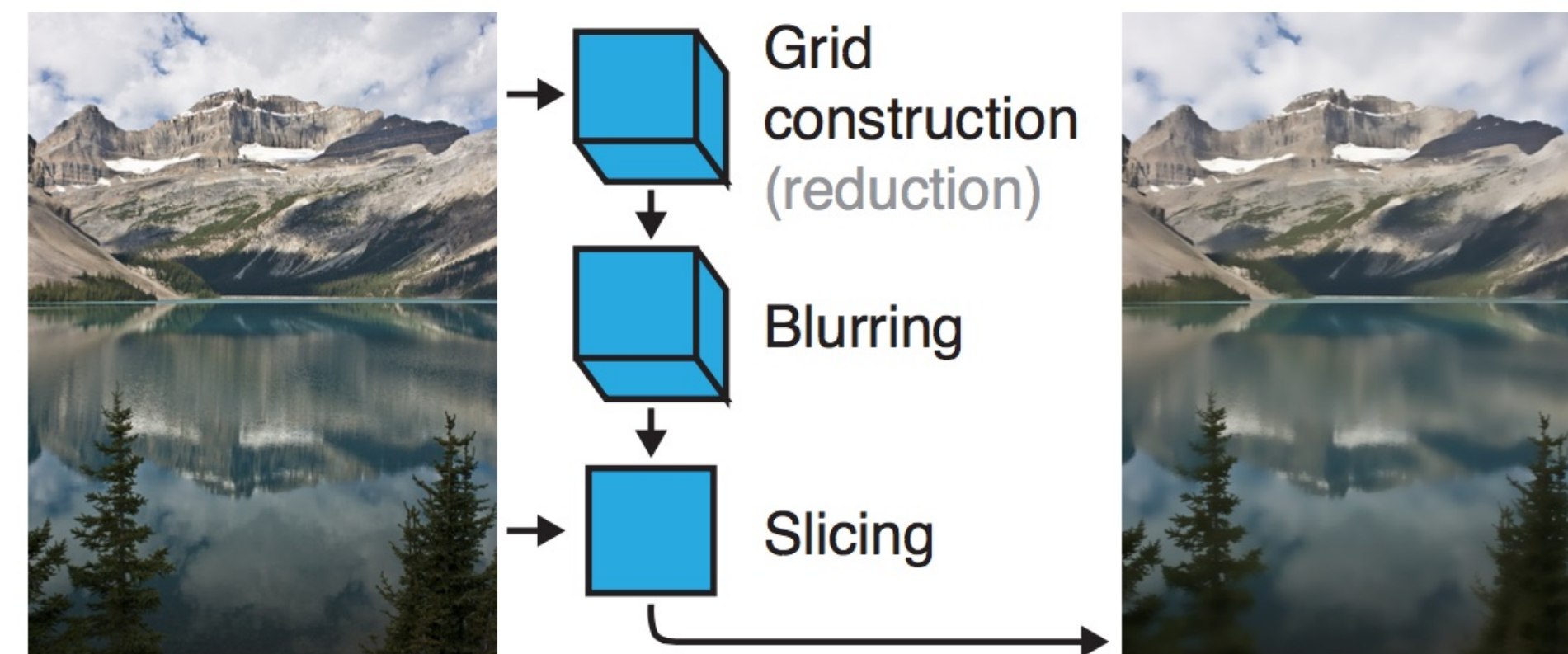
# Initial academic Halide results

[Ragan-Kelley 2012]

- **Application 1: camera RAW processing pipeline**  
(Convert RAW sensor data to RGB image)
  - **Original: 463 lines of hand-tuned ARM NEON assembly**
  - **Halide: 2.75x less code, 5% faster**



- **Application 2: bilateral filter**  
(Common image filtering operation used in many applications)
  - **Original 122 lines of C++**
  - **Halide: 34 lines algorithm + 6 lines schedule**
    - **CPU implementation: 5.9x faster**
    - **GPU implementation: 2x faster than hand-written CUDA**



# Stepping back: what is Halide?

- **Halide is a DSL for helping expert developers optimize image processing code more rapidly**
  - **Halide does not decide how to optimize a program for a novice programmer (ignoring the auto scheduler, see tonight's reading)**
  - **Halide provides a small number of primitives for a programmer that has strong knowledge of code optimization to rapidly express what optimizations the system should apply**
    - **parallel, vector, unroll, split, reorder, store\_at, compute\_at...**
  - **Halide compiler carries out the mapping of that strategy to a machine**

# Automatically generating Halide schedules

- **Problem: it turned out that very few programmers have the technical ability to write good Halide schedules**
  - Circa 2017... 80+ programmers at Google write Halide
  - Very small number trusted to write schedules
- **Recent work: Halide compiler analyzes the Halide program to automatically generate efficient schedules for the programmer [Mullapudi 2016, Adams 2019]**
  - As of Adams 2019, you'd have to work hard to manually author a schedule that is better than the schedule generated by the Halide autoscheduler for a complex image processing pipeline

# Influence on code generation for ML applications

## Example: Apache TVM



## Apache TVM

An End to End Machine Learning Compiler Framework for CPUs, GPUs and accelerators

### ☰ Schedule Primitives in TVM

split

tile

fuse

reorder

bind

compute\_at

compute\_inline

compute\_root

Summary

Reduction

## Tuning Parameters of Thread Numbers

How to schedule the workload, say, 32x32 among the threads of one cuda block? Intuitively, it should be like 1

```
num_thread_y = 8
num_thread_x = 8
thread_y = tvm.thread_axis((0, num_thread_y), "threadIdx.y")
thread_x = tvm.thread_axis((0, num_thread_x), "threadIdx.x")
ty, yi = s[Output].split(h_dim, nparts=num_thread_y)
tx, xi = s[Output].split(w_dim, nparts=num_thread_x)
s[Output].reorder(ty, tx, yi, xi)
s[Output].bind(ty, thread_y)
s[Output].bind(tx, thread_x)
```

There are two parameters in the schedule: `num_thread_y` and `num_thread_x`. How to determine the optimal  
Below is the result with Filter = [256, 1, 3, 3] and stride = [1, 1]:

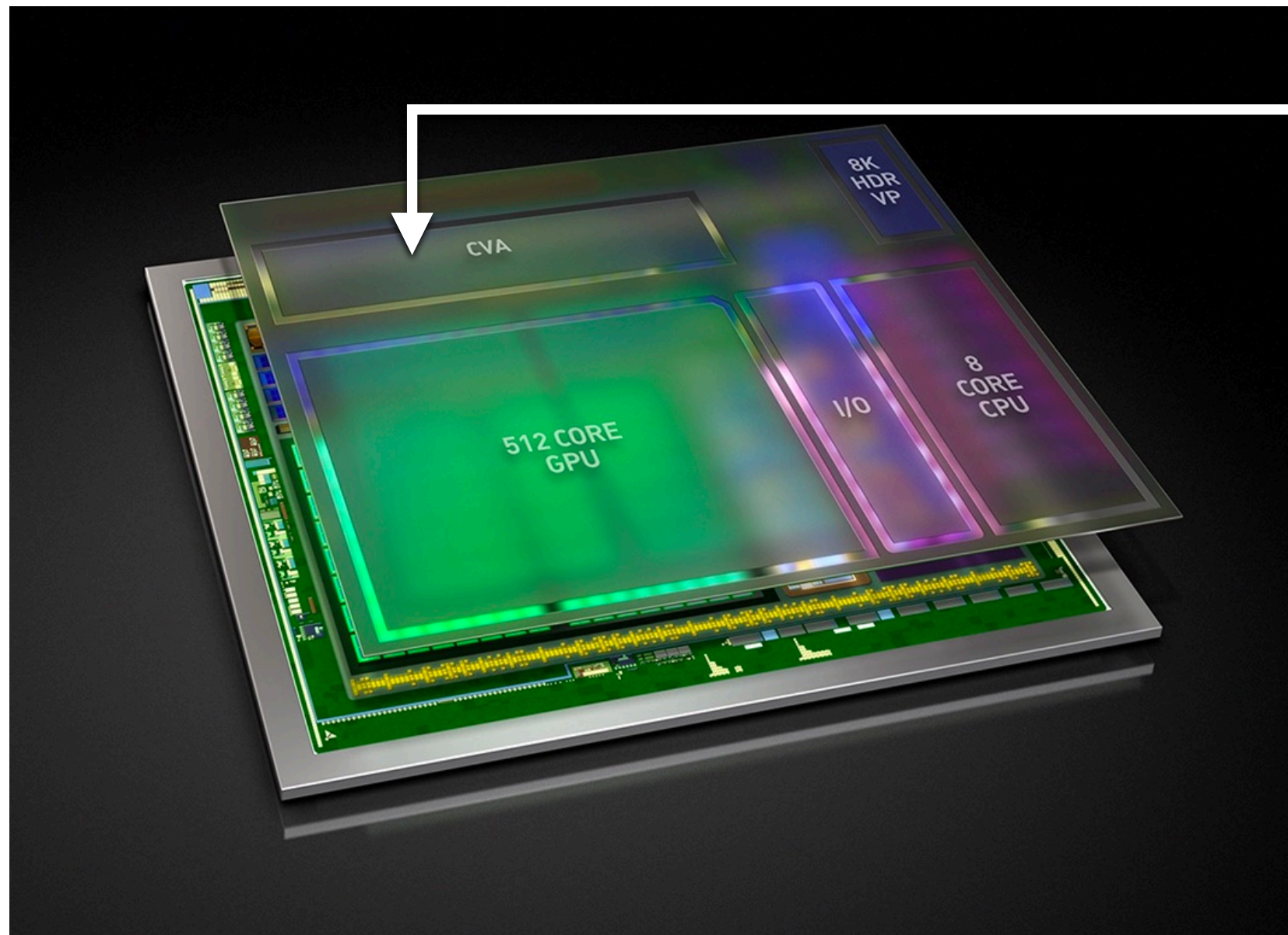
Case	Input	num_thread_y	num_thread_x
1	[1, 256, 32, 32]	8	32
2	[1, 256, 32, 32]	4	32
3	[1, 256, 32, 32]	1	32
4	[1, 256, 32, 32]	32	1

Many interesting observations from above results:

# AI Accelerators

# Efficiency estimates \*

- **Estimated overhead of programmability (instruction stream, control, etc.)**
  - **Half-precision FMA (fused multiply-add) 2000%**
  - **Half-precision DP4 (vec4 dot product) 500%**
  - **Half-precision 4x4 MMA (matrix-matrix multiply + accumulate) 27%**



**NVIDIA Xavier (SoC for automotive domain)**

**Features a Computer Vision Accelerator (CVA), a custom module for deep learning acceleration (large matrix multiply unit)**

**~ 2x more efficient than NVIDIA V100 MMA instruction despite being highly specialized component. (includes optimization of gating multipliers if either operand is zero)**

\* Estimates by Bill Dally using academic numbers, SysML talk, Feb 2018

# Ampere GPU SM (A100)

Each SM core has:

64 fp32 ALUs (mul-add)

32 int32 ALUs

4 “tensor cores”

Execute  $8 \times 4 \times 4 \times 8$  matrix mul-add instr

$A \times B + C$  for matrices A,B,C

A, B stored as fp16, accumulation with fp32 C

There are 108 SM cores in the GA100 GPU:

6,912 fp32 mul-add ALUs

432 tensor cores

1.4 GHz max clock

= 19.5 TFLOPs fp32

+ 312 TFLOPs (fp16/32 mixed) in tensor cores



Single instruction to perform  
 $2 \times 8 \times 4 \times 8$  FP16 +  $8 \times 8$  TF32 ops



The NVIDIA tensor core approach is an evolutionary design: add DNN-specific instructions to a traditional programmable processor (“evolve, don’t replace”)

# Numerical data formats



Range

$10^{-38} - 10^{38}$

Accuracy

.000006%

**Reminder:**

$$-1^S \times (1 + (M \times 2^{-23})) \times 2^{(E-127)}$$



$6 \times 10^{-5} - 6 \times 10^4$

.05%



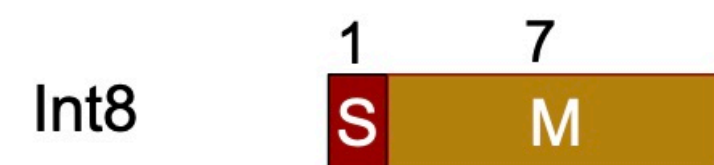
$0 - 2 \times 10^9$

33%



$0 - 6 \times 10^4$

33%

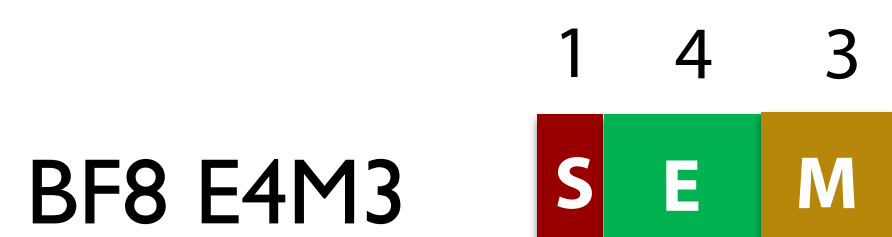


$0 - 127$

33%



**BF16: Same range as FP32, but lower accuracy**



$0 - 448$



$0 - 57344$

# Generating efficient code for AI accelerators

# What are the right primitives?

- **Halide adopts low-level primitives (mathematical expressions on values) for expressing image processing applications**
- **DNN operations often built on higher-level primitives like matrix multiplication ops, movement of matrices between locations in memory**
- **Could a Halide-like scheduling approach apply to different sets of primitives?**

## **ThunderKittens: A Simple Embedded DSL for AI kernels**

Benjamin Spector, Aaryan Singhal, Simran Arora, Chris Re

(This post is an extended introduction to a [longer post](#) we're releasing concurrently.)



# Modern look at “kernel generation”

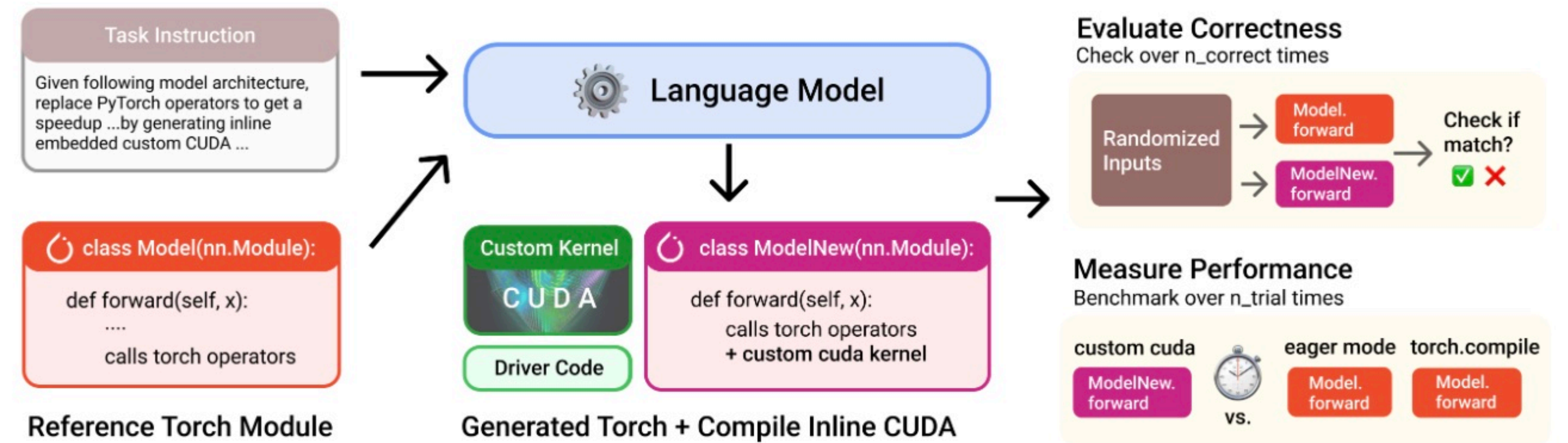
- Given interest in usage of LLMs for code generation, rekindled interest in automation of DNN kernels
- I'd like to hear the class' thoughts

## KernelBench: Can LLMs Write GPU Kernels?

 Anne Ouyang\*  
Stanford

 Simon Guo\*  
Stanford

 Azalia Mirhoseini  
Stanford



A benchmark designed to evaluate the ability of LLMs to generate efficient GPU kernels for optimizing neural network performance

# **A few comments on projects**

# Projects

- **Team size**
  - **Teams of 2 (no questions asked)**
  - **Teams of 3 (possible, but need permission and a clear description of roles)**
  
- **Topics: must meet two requirements**
  - **You are really excited to work a lot on it**
  - **Forces you to perform work that applies the intellectual themes of the course**
    - **Clearly defining a job (defining definition of “good”)**
    - **Choosing the right algorithms for the job, or designing abstractions/primitives to do the job**
    - **Measuring how your choices are responsible for meeting the definition of good**
  
- **Presentations will be 10 minute presentations on the last class day of the quarter — times TBD**

# Ideas from class topics so far

- **Implement as much as the Google HDR+ pipeline as you can**
  - **I have sample assignments and can point you at input data**
  - **Or implement and extend another recent computational photography feature**
- **For an application space that you know well, attempt to design an API for authoring programs for this application space**
  - **Or for a space of scheduling choices...**