

# High-Performance Image Processing

*Last lecture on Halide*

Frédo Durand  
MIT CSAIL

# High-Performance Image Processing

*Last lecture on Halide*

Frédo Durand  
MIT CSAIL

# Halide reductions

Loops are implicit !!

Reduction domain: all the location that will be aggregated

Multidimensional

RDom(baseX, extentX, baseY, extentY,...)

reduction domain

Initialize your Func

myFunc(Var1, Var2)=initialValue;

→ implicit loop  
i.e. you specify end → start  
unlike regular loops

Update equation

myFunc(Expr, Expr) = f ( myFunc(Expr, Expr), rdom );

will be called for each RDom location

arbitrary Expr of the Func, its Var, and the RDom. The RDom can be on the left and right

# Sum in Halide

```
Image input = Image(Float(32), im);
```

```
Var x, y, c;
```

```
Func mySum;
```

```
RDom r(0, input.width(), 0, input.height()); reduction domain  
mySum(c)=0.0; init  
r.y → first dim of an RDom  
is always called  $\approx$ 
```

```
mySum(c) +=input(r.x, r.y, c); update  
→ implicit loop
```

```
output = mySum.realize(input.channels());
```

# Sum pseudocode

```
for c in (0..input.channels()): }  
    out[c]=0.0  
  
for ry in (0..input.height()): } Reduction loop  
    for rx in (0..input.width()): }  
        for c in (0..input.channels()): } Output loop  
            out[c] += input[rx, ry, c] } update
```

Note that the reduction loops are the outer loop

# In general, updates

**Funcs must always have a first pure-function initialization**

similar to what we've done so far

e.g.  $f(x,y)=\text{input}(x,y)^*2$ ; or  $f(x,y)=0.0$ ;

**Then they can have multiple updates**

either at specific locations:  $f(32,46) = 1.0$ ;

or sets of locations:  $f(x, 17)=x+1$ ; (applies to all  $x$ , implicit loop)

or with RDoms (implicit loops, more flexibility for arithmetics)

**In general:**

`myFunc(Expr,Expr) = f ( myFunc(Expr,Expr), Rdom );`

will be called for each RDom location

**Unlike Func definitions, updates can have Expr on LHS**

# Sum pseudocode

```
for c in (0..input.channels()): }  
    out[c]=0.0  
  
for ry in (0..input.height()): } Reduction loop  
    for rx in (0..input.width()): }  
        for c in (0..input.channels()): } Output loop  
            out[c] += input[rx, ry, c] } update
```

Note that the reduction loops are the outer loop

# Sum in Halide

```
Image input = Image(Float(32), im);
```

```
Var x, y, c;
```

```
Func mySum;
```

```
RDom r(0, input.width(), 0, input.height()); reduction domain  
mySum(c)=0.0; init  
r.y → first dim of an RDom  
is always called  $\approx$ 
```

```
mySum(c) +=input(r.x, r.y, c); update  
→ implicit loop
```

```
output = mySum.realize(input.channels());  
3
```

# Sum pseudocode

```
for c in (0..input.channels()): }  
    out[c]=0.0  
  
for ry in (0..input.height()): } Reduction loop  
    for rx in (0..input.width()): }  
        for c in (0..input.channels()): } Output loop  
            out[c] += input[rx, ry, c] } update
```

Note that the reduction loops are the outer loop

# In general, updates

**Funcs must always have a first pure-function initialization**

similar to what we've done so far

e.g.  $f(x,y)=\text{input}(x,y)^*2$ ; or  $f(x,y)=0.0$ ;

**Then they can have multiple updates**

either at specific locations:  $f(32,46) = 1.0$ ;

or sets of locations:  $f(x, 17)=x+1$ ; (applies to all  $x$ , implicit loop)

or with RDoms (implicit loops, more flexibility for arithmetics)

**In general:**

`myFunc(Expr,Expr) = f ( myFunc(Expr,Expr), Rdom );`

will be called for each RDom location

**Unlike Func definitions, updates can have Expr on LHS**

# Pure Var vs. RDom

**Pure Var update:**  $f(x, 17) = x + 6$ .

only for output pixel coordinates

pure Var ( $x$ ) can only appear unmodified on LHS,  
and at the same location as in pure definition.

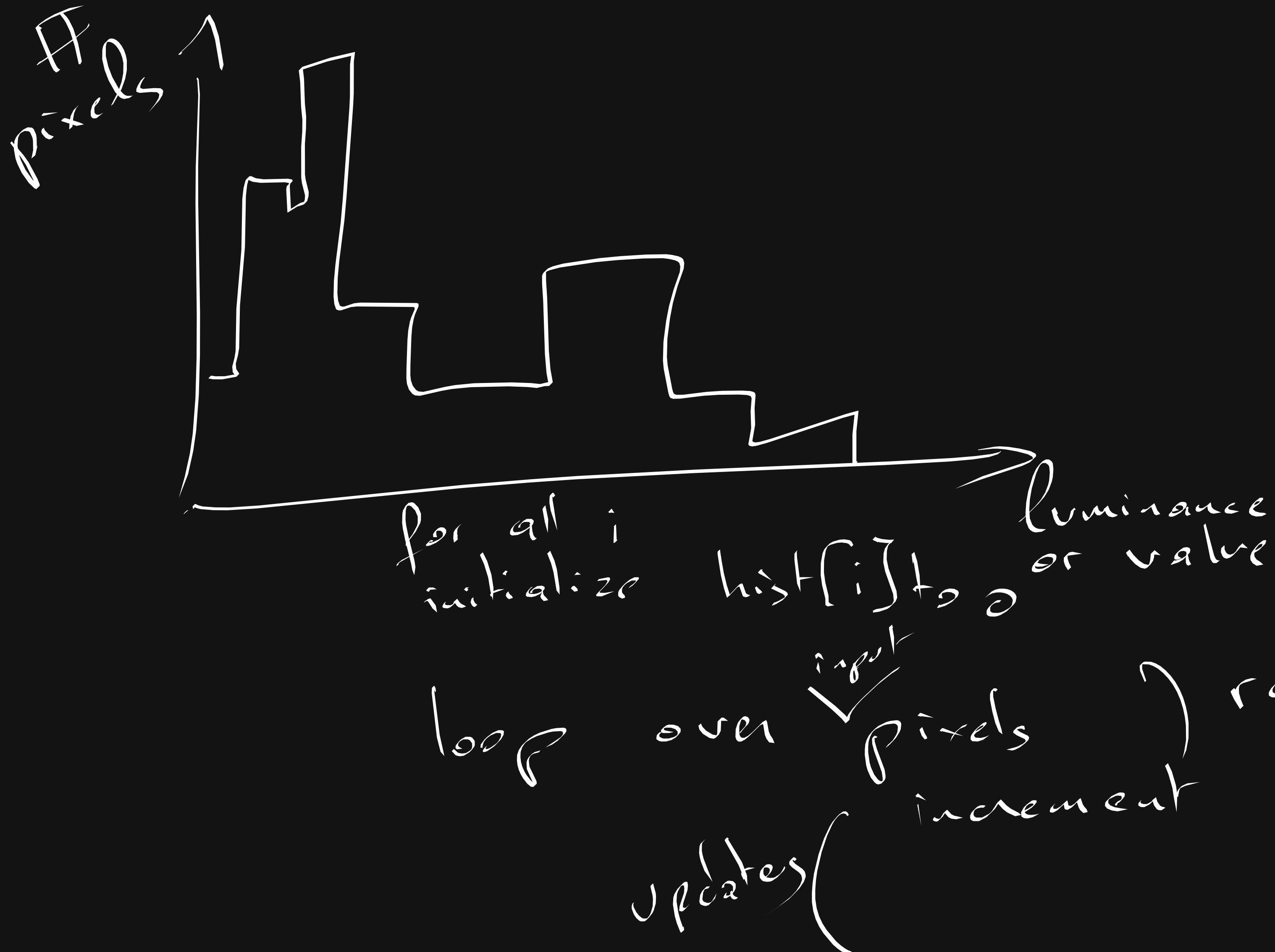
e.g.  $f(x+1) = x * 7$ ; or  $f(y, x) = 16$ ; are invalid (assuming it's  $f(x, y)$ )

**RDom**s can be used more arbitrarily

typically iterate over input pixels



# Histograms (for 8 bit images)



# Histogram in C (for 8 bit images)

```
int c_result[256];      output array
for (int x = 0; x < 256; x++) {           init.
    c_result[x] = 0;
}
for (int r_y = 0; r_y < input.height(); r_y++) {           reduction
    for (int r_x = 0; r_x < input.width(); r_x++) {
        c_result[input(r_x, r_y)] += 1;           } update
    }
}
```

# Histogram in Halide (assuming 8 bit images)

assume input, call it  $i$

```
Func histogram("histogram");          (output)

Var i ("i");
histogram(i) = 0;
Random r (0, in.width(), 0, in.height(), "r");

histogram(in(r.x, r.y)) += 1           update
```

- realize -- (256)

# Histogram in Halide

```
Func histogram("histogram");
Var x ("x") =  
    histogram(x) = 0;           init  
  
RDom r(0, input.width(), 0, input.height());      implicit loop  
histogram(input(r.x, r.y)) += 1;                  update  
  
Image<int> halide_result = histogram.realize(256);
```

# C equivalent

```
int c_result[256];
for (int x = 0; x < 256; x++) {
    c_result[x] = 0;
}
for (int r_y = 0; r_y < input.height(); r_y++) {
    for (int r_x = 0; r_x < input.width(); r_x++) {
        c_result[input(r_x, r_y)] += 1;
    }
}
```

redacted  
domain

Note: no loop over output indices  
because no Var was used  
in update equation  
(unlike sum)



# Convolution as reduction

O  
( for each pure var output pixel at  $x, y$   
init output( $x, y$ ) = O  
loop over kernel RDom cx cy  
add pixel ( $x - cx, y - cy$ )  
+ maybe divide by normalization

# Convolution as reduction

Box blur

```
Func blur ("blur"), clamped("clamped");  
Var x("x"), y("y"), c("c");  
clamped(x, y, c) = input(clamp(x, 0, input.width()-1),  
                           clamp(y, 0, input.height()-1), c);
```

$$\text{blur}(x, y, c) = O_O ;$$

Random r(O, kernel\_size; O, kernel\_size);

$\text{blur}(x, y, c) += \text{clamped}\left(x + r \cdot x - \text{kernel\_size}/2,$   
 $y + r \cdot y - \text{kernel\_size}/2,$   
 $c\right) / \text{pow}(\text{kernel\_size}, 2);$   
 $* \text{kernel}(r \cdot x, r \cdot y)$  for general  
case

blur.realize(—)

# Convolution as reduction

```
Func blur ("blur"), clamped("clamped");
Var x("x"), y("y"), c("c");
clamped(x, y, c) = input(clamp(x, 0, input.width()-1),
                           clamp(y, 0, input.height()-1), c);

RDom r(0, kernel_width, 0, kernel_width, "r");

blur(x,y,c) = 0.0;
blur(x,y,c) += clamped(x+r.x-kernel_width/2,
                        y+r.y-kernel_width/2,
                        c) / pow(kernel_width,2);
```

# Note

Again, no loop!

Here, reduction over 2 of the 3 dimensions

# Convolution as reduction

```
Func blur ("blur"), clamped("clamped");  
Var x("x"), y("y"), c("c");  
clamped(x, y, c) = input(clamp(x, 0, input.width()-1),  
                           clamp(y, 0, input.height()-1), c);
```

```
RDom r(0, kernel_width, 0, kernel_width, "r");
```

```
blur(x,y,c) = 0.0;  
blur(x,y,c) += clamped(x+r.x-kernel_width/2,  
                        y+r.y-kernel_width/2,  
                        c) / pow(kernel_width,2);
```

*pure  
Var*

*RDom*

*update*

*→ 5 loops*

# Note

Again, no loop!

Here, reduction over 2 of the 3 dimensions

# Equivalent

```
for y in xrange(input.height()):  
    for x in xrange(input.width()):  
        for c in xrange(input.channels):  
            out[x,y,c]=0  
  
for ry in xrange(kernel_width):  
    for rx in xrange(kernel_width):  
        for y in xrange(input.height()):  
            for x in xrange(input.width()):  
                for c in xrange(input.channels):  
                    out[x,y,c]+=clampedInput[x+rx kernel_width/2,  
                                         y+ry-kernel_width/2,  
                                         c] / kernel_width**2
```

int

RD<sub>src</sub>

Pure Var

update

# Convolution as reduction - alternative

```
Func blur ("blur"), clamped("clamped");
Var x("x"), y("y"), c("c");
clamped(x, y, c) = input(clamp(x, 0, input.width()-1),
                           clamp(y, 0, input.height()-1), c);

RDom r(0, kernel_width, 0, kernel_width, "r");

blur(x,y,c) = 0.0;

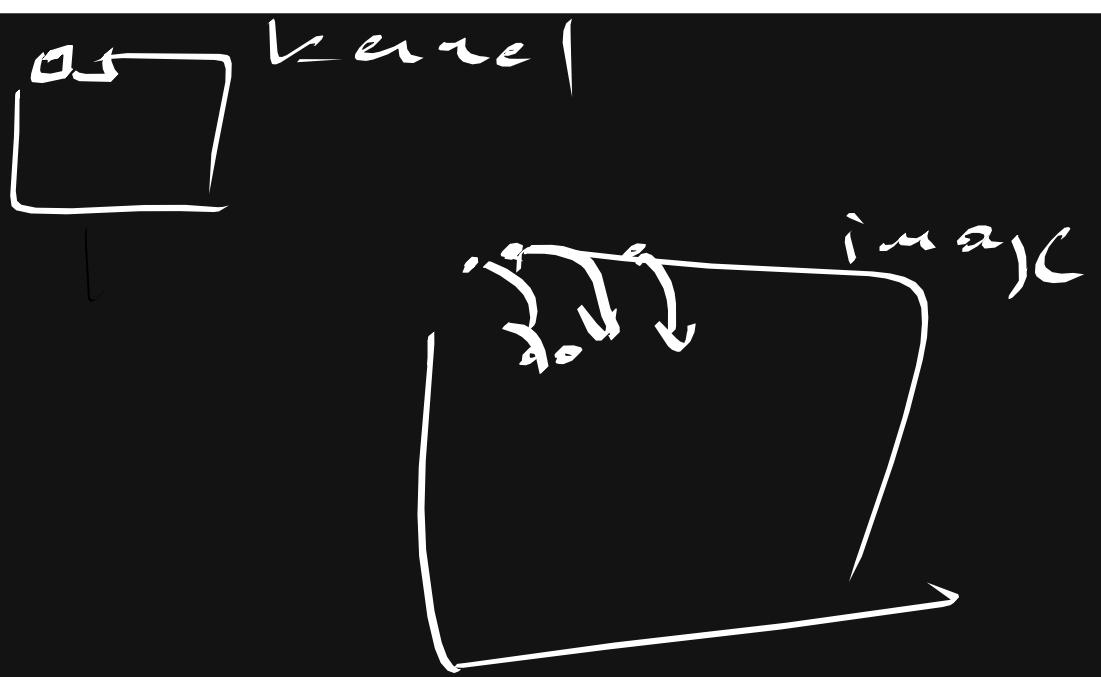
blur(x,y,c) += clamped(x+r.x-kernel_width/2,
                        y+r.y-kernel_width/2,
                        c);
blur(x,y,c) /= pow(kernel_width,2);
```



# Loop order

```
for ry in xrange(kernel_width):  
    for rx in xrange(kernel_width):  
        for y in xrange(input.height()):  
            for x in xrange(input.width()):  
                for c in xrange(input.channels):  
                    out[x,y,c] += clampedInput[x+rx, y+ry, c] / kernel_width**2
```

RPN  
Post  
order



# Problem, bad order

```
for ry in xrange(kernel_width):  
    for rx in xrange(kernel_width):  
        for y in xrange(input.height()):  
            for x in xrange(input.width()):  
                for c in xrange(input.channels):  
                    out[x,y,c] += clampedInput[x+rx kernel_width/2,  
                                         y+ry-kernel_width/2,  
                                         c] / kernel_width**2
```

**Reduction loops are always outside**

Because otherwise the update semantics might be changed

**Can result in very bad locality for convolution**

**Cannot always be reordered directly (semantics should not change)**

But there is a trick

# **Inline/helper trick**

**Trick to make Halide do the right thing with reductions**

**I'll show you how it works**

**Then I'll show you a sugar that hides it from you**



# The helper/inline trick

same code as before

```
Func blur ("blur"), clamped("clamped");  
Var x("x"), y("y"), c("c");  
clamped(x, y, c) = input(clamp(x, 0, input.width()-1),  
                           clamp(y, 0, input.height()-1), c);
```

```
RDom r(0, kernel_width, 0, kernel_width, "r");
```

```
blur(x,y,c) = 0.0;  
blur(x,y,c) += clamped(x+r.x-kernel_width/2,  
                           y+r.y-kernel_width/2,  
                           c) / pow(kernel_width,2);
```

```
Func superBlur("superBlur");  
superBlur(x,y,c)=blur(x,y,c)
```

done

# Trick (default) schedule equivalent

```
produce {  
    RDom r(0, kernel_width, 0, kernel_width, 'r');  
    blur(x,y,c) = 0.0;  
    blur(x,y,c) += clamped(x+r.x-kernel_width/2, y+r.y-kernel_width/2, c) /  
    pow(kernel_width,2);  
}  
end [ superBlur(x,y,c)=blur(x,y,c) ]  
          output final consumer → root  
          simplified version code → inline
```

for y  
for x  
for c  
 "they produce all need values of producer  
 // for 1 consumer at x,y,c  
 // i.e. one value of blur at x,y,c

```

produce {
    RDom r(0, kernel_width, 0, kernel_width, 'r');
    blur(x,y,c) = 0.0;
    blur(x,y,c) += clamped(x+r.x-kernel_width/2, y+r.y-kernel_width/2, c) / pow(kernel_width,2);
}
add [ superBlur(x,y,c)=blur(x,y,c) ] output final consumer → root

```

for y

for x

for c

// here produce all need values of produce  
// for 1 consumer at x,y,c  
// i.e. one value of blur at x,y,c

for r.y in (0..kernel\_w)  
 for r.x in (0..kernel\_w)

simp^1  
vector code  
of

inline

```

RDom r(0, kernel_width, 0, kernel_width, 'r');
blur(x,y,c) = 0.0;
blur(x,y,c) += clamped(x+r.x-kernel_width/2, y+r.y-kernel_width/2, c) / pow(kernel_width,2);
} else [ superBlur(x,y,c)=blur(x,y,c) ] → inline

```

→ output final consumer → root

for y  
 for x  
 for c

// user produce all need values of producer  
 // for 1 consumer at x,y,c  
 // i.e. one value of blur at x,y,c

for r.y in (0..kernel\_w)  
 for r.x (0..kernel\_w)  
 for y' in (y, y+kernel\_w)  
 for x' in (x, x+kernel\_w)

```
[pow(kernel_width, 2);
```

```
ad [ superBlur(x,y,c)=blur(x,y,c) ] output final consumer → root
```

for y

for x

for c

then produce all need values of producer

// for 1 consumer at x,y,c

// i.e. one value of blur at x,y,c

for r,y in (0 .. kernel\_w)

for r,x in (0 .. kernel\_w)

for y' in (y, u)

for x' in (x, v)

for c in (c, d)

```
ad [ superBlur(x,y,c)=blur(x,y,c)           output final consumer → root  
for y  
for x  
for c  
    blur produce all need values of produce  
    // for 1 consumer at x,y,c  
    // i.e. one value of blur at x,y,c
```

```
for r,y in (0 .. level-w)  
for r,x in (0 .. level-w)  
for y' in (y,w)  
for x' in (x,w)
```

for  $\lambda$

for  $\alpha$

for  $c$

then produce all need values of producer

// for 1 consumer at  $x, y, c$

// i.e. one value of blu at  $x, y, c$

for  $x, y$  in  $(0 .. land-w)$

for  $c, \alpha$  in  $(0 .. land-w)$

for  $\lambda$  in  $(y, u)$

for  $\beta$  in  $(x, z)$

for  $c$  in  $(c, d)$

for  $\lambda$

for  $\alpha$

for  $c$

"then produce all need values of producer

// for 1 consumer at  $x, y, c$

// i.e. one value of  $b_{\lambda}$  at  $x, y, c$

for  $\beta, \gamma$  in  $(0 .. \text{Lanc} - w)$

for  $\alpha$  in  $(0 .. \text{Lanc} - w)$

for  $\lambda$  in  $(\beta, \gamma)$

for  $c$  in  $(x, z)$

for  $y$

for  $x$

for  $c$

"then produce all need values of producer

// for 1 consumer at  $x, y, c$   
// i.e. one value of blu at  $x, y, c$

for  $y$  in  $Y_{\text{US}}$

for  $x$  in  $X_{\text{US}}$   
for  $c$  in  $C_{\text{US}}$   
 $\text{blu}(x, y, c) = 0$

for  $y$  in  $(0 \dots \text{land-w})$   
for  $x$  in  $(0 \dots \text{land-w})$

for  $y'$  in  $(y, y)$   
for  $x'$  in  $(x, x)$

accumulate

# Equivalent

```
for y in xrange(input.height()):  
    for x in xrange(input.width()):  
        for c in xrange(input.channels):
```

# Python equivalent

```
for y in xrange(input.height()):
    for x in xrange(input.width()):
        for c in xrange(input.channels):
            tmp=numpy.empty([1,1,1])
            for yi in xrange(1):
                for xi in xrange(1):
                    for ci in xrange(1):
                        tmp[xi,yi,ci]=0
            for ry in xrange(kernel_width):
                for rx in xrange(kernel_width):
                    for yi in xrange(1):
                        for xi in xrange(1):
                            for ci in xrange(1):
                                tmp[xi,yi,ci]+=clampedInput[x+rx kernel_width/2,
                                                               y+ry-kernel_width/2,
                                                               c] / kernel_width**2
                    superBlur[x,y,c]=tmp[0,0,0]
```

# Python equivalent without 1-iteration loops

```
superBlur=numpy.empty([input.width(), input.height(), input.channels()])
for y in xrange(input.height()):
    for x in xrange(input.width()):
        for c in xrange(input.channels():

            tmp=0
            for ry in xrange(kernel_width):
                for rx in xrange(kernel_width):

                    tmp +=clampedInput[x+rx kernel_width/2,
                                         y+ry-kernel_width/2,
                                         c] / kernel_width**2

            superBlur[x,y,c]=tmp
```

# **Rule of thumb**

**For stencil reduction**

i.e. each output pixel is a reduction over multiple input pixels

**Use the helper/inline trick**

**Encapsulate into an extra Func scheduled as default (inline)**

**In general, true for reduction where the free variables are independent of the reduction variables.**



# Sugar: sum (includes helper/inline trick)

```
Func blur("blur");
Var x("x"), y("y"), c("c");
Func clamped("clamped");
clamped(x, y, c) = input(clamp(x, 0, input.width()-1),
                           clamp(y, 0, input.height()-1), c);

RDom r(0, kernel_width, 0, kernel_width, "r");

blur(x,y,c) = sum(clamped(x+r*x-kernel_width/2, y, c)) ;  

                           ↓  

                           y+r*y-kernel_width/2
```

Under the hood: initializes, creates a helper/inline

# For hardcore Halide

(beyond the scope of this class & a little dirty)

`unroll_sum`

↳ conflates schedule  $\Sigma$  algorithm

does what it says



# Computing the kernel

for this class, I advise you do it in C++ and pass it to Halide as  
an image (possibly of size 1xkernel\_size)

Could be done in Halide  
how should it be scheduled?

*compute\_root()*



# More complex reductions

Bilateral grid

Output over grid

Reduction over input pixels

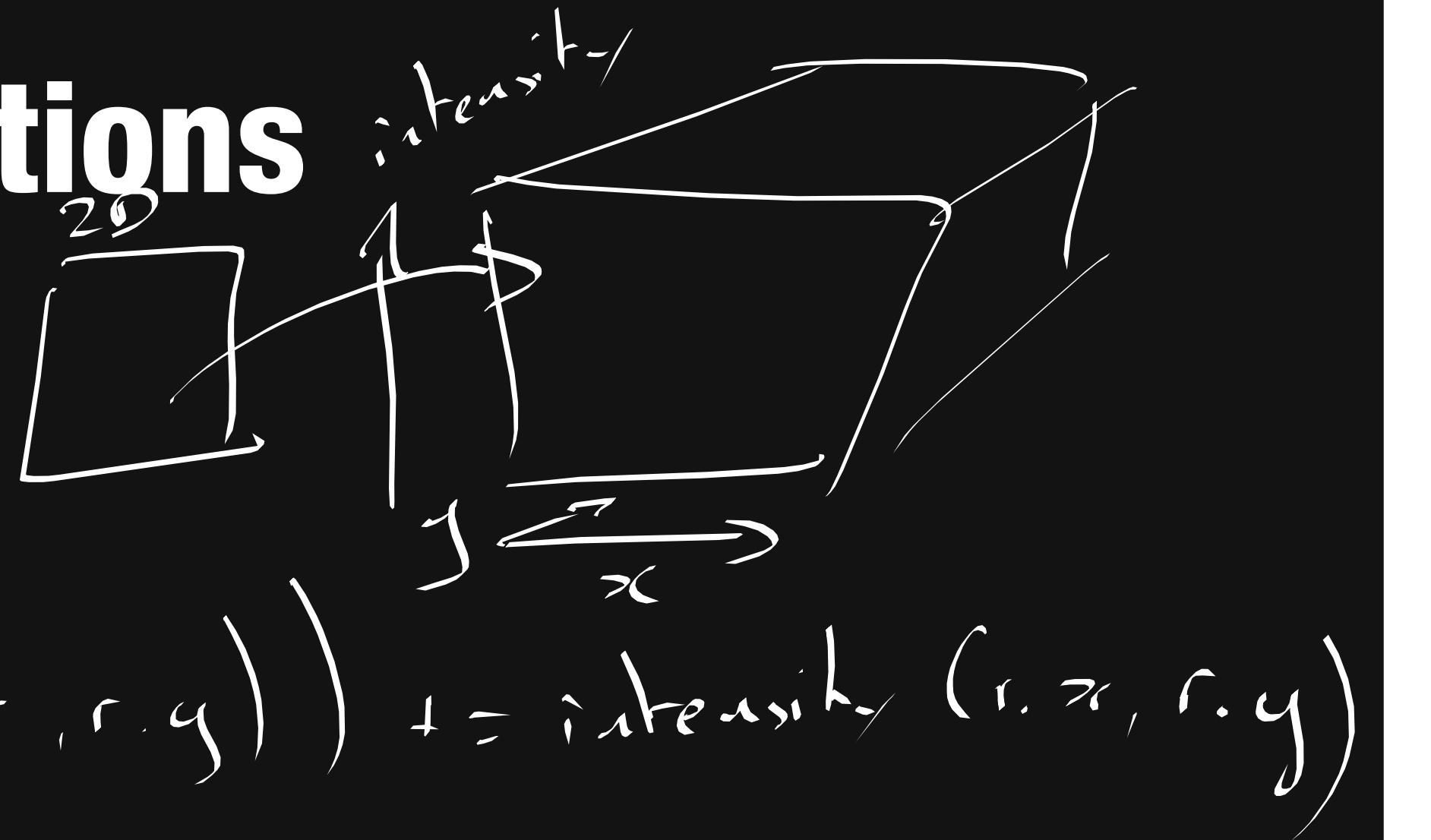
Update equation:

: grid construction

input pixels

grid ( $r_{x,y}$ ,  $c_{x,y}$ , intensity ( $r_{x,y}, c_{x,y}$ ))  $\leftarrow$  intensity ( $r_{x,y}, c_{x,y}$ )

+ probably  
downsampling



# More complex reductions

## Bilateral grid

Output over grid cells  $x, y, I$

Reduction over input image  $x, y$

Update equation:  $\text{grid}(rx, ry, \text{intensity}(rx, ry)) += \text{intensity}(rx, ry);$   
(plus another one for sum of weights)

**Note how the left-hand side is an Expr of the input image**

Only case where LHS can have something else than Var to index Funcs



# Compilation

We'll just JIT things (just in time)

your code calls the Halide library when it runs.

the Halide library includes a full compiler (LLVM) and generates the binary on the fly

**Alternative: traditional static compilation**

`compile_to_file`

useful to ship code, slightly more work to do the bindings .

# Other goodies

**store\_at**

**GPU specific (CUDA, OpenCL)**

gpu\_threads, gpu\_blocks, gpu\_tile

**tuples**

**extern**

**parameters**

