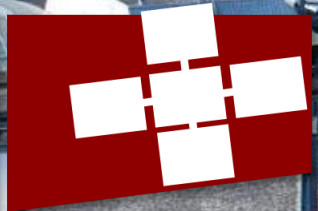**ETH**zürich

**D**INFK

J. DE FINE LICHT AND T. HOEFLER

# Productive Parallel Programming on FPGA with High-Level Synthesis

# Based on material from:

*Transformations of High-Level Synthesis Codes for High-Performance Computing*
https://arxiv.org/abs/1805.08288

# Code examples found at:

https://github.com/spcl/hls_tutorial_examples

## Virtual machine for emulation

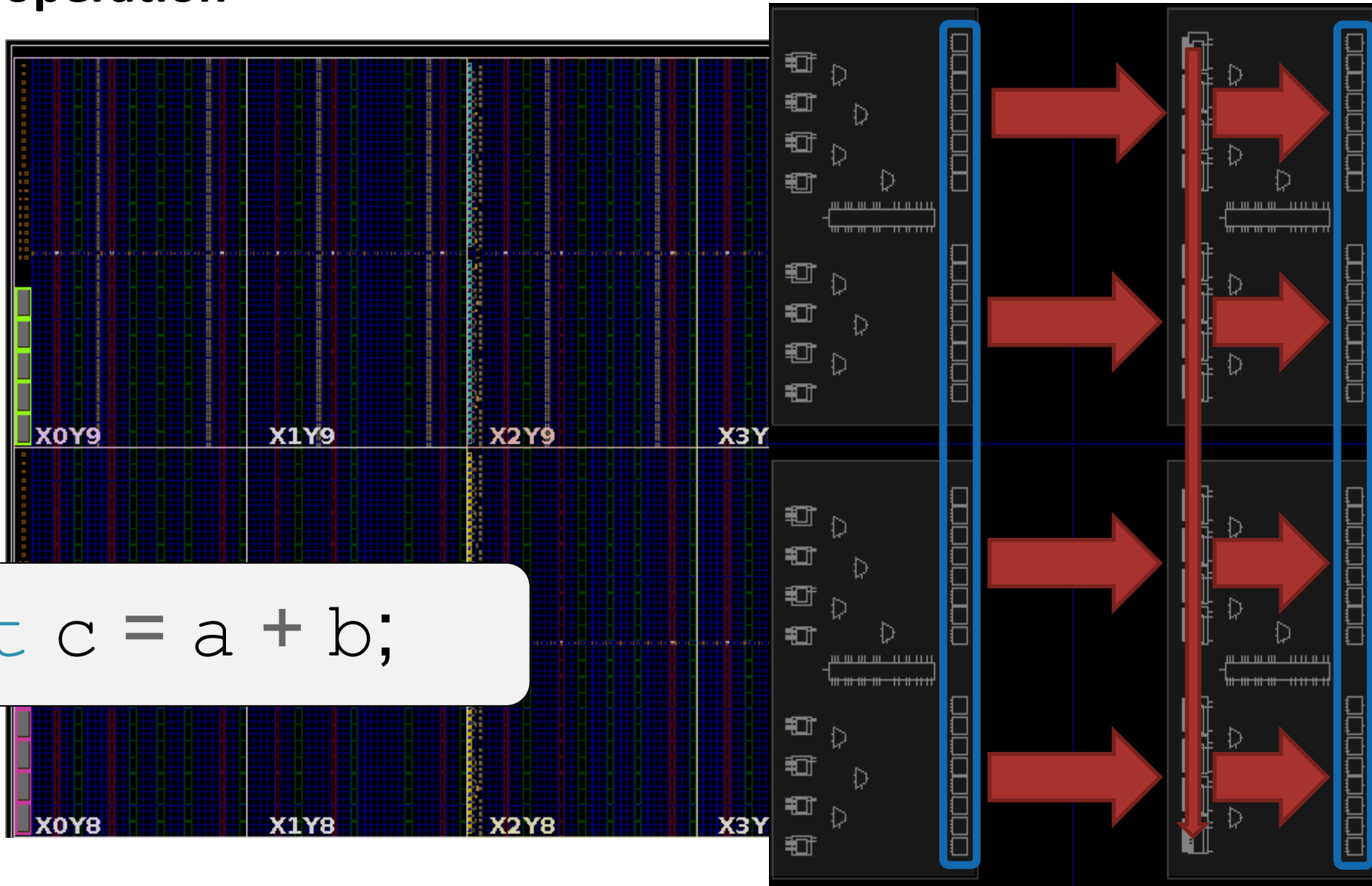http://spcl.inf.ethz.ch/~definelj/HLS_Tutorial.7z

## Nimbix Alveo Trial

https://www.nimbix.net/alveotrial
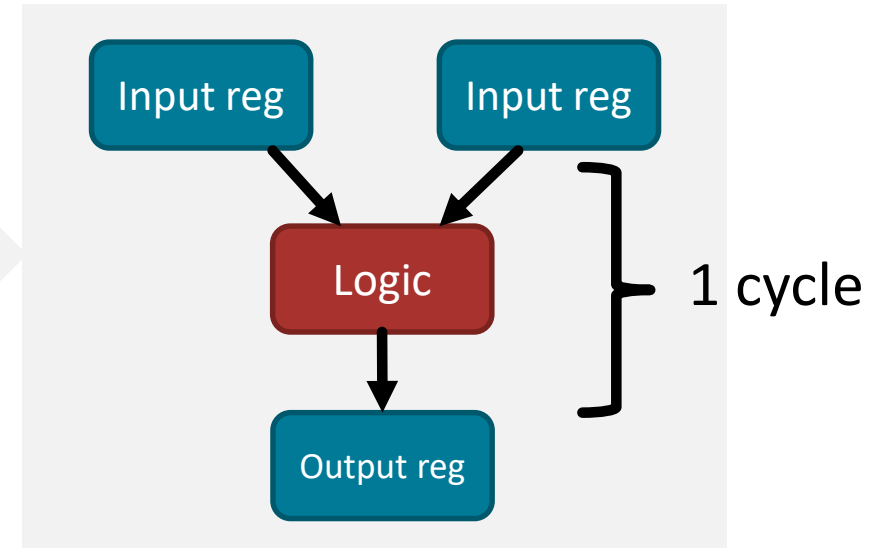
# Our goal



Perform compute in every piece of available logic – every cycle!!

# A single operation

```
int c = a + b;
```

X0Y9    X1Y9    X2Y9    X3Y

X0Y8    X1Y8    X2Y8    X3Y

# Register transfer level

```
always @(posedge clk)
    if (start) begin
        out <= in + 1;
    end
```



1 cycle

```
int c = a + b;
```

# Single <u>floating point operation</u>



Reg    Reg

Logic

1 cycle

*Transient register*

$L$ = 8 cycles

Logic
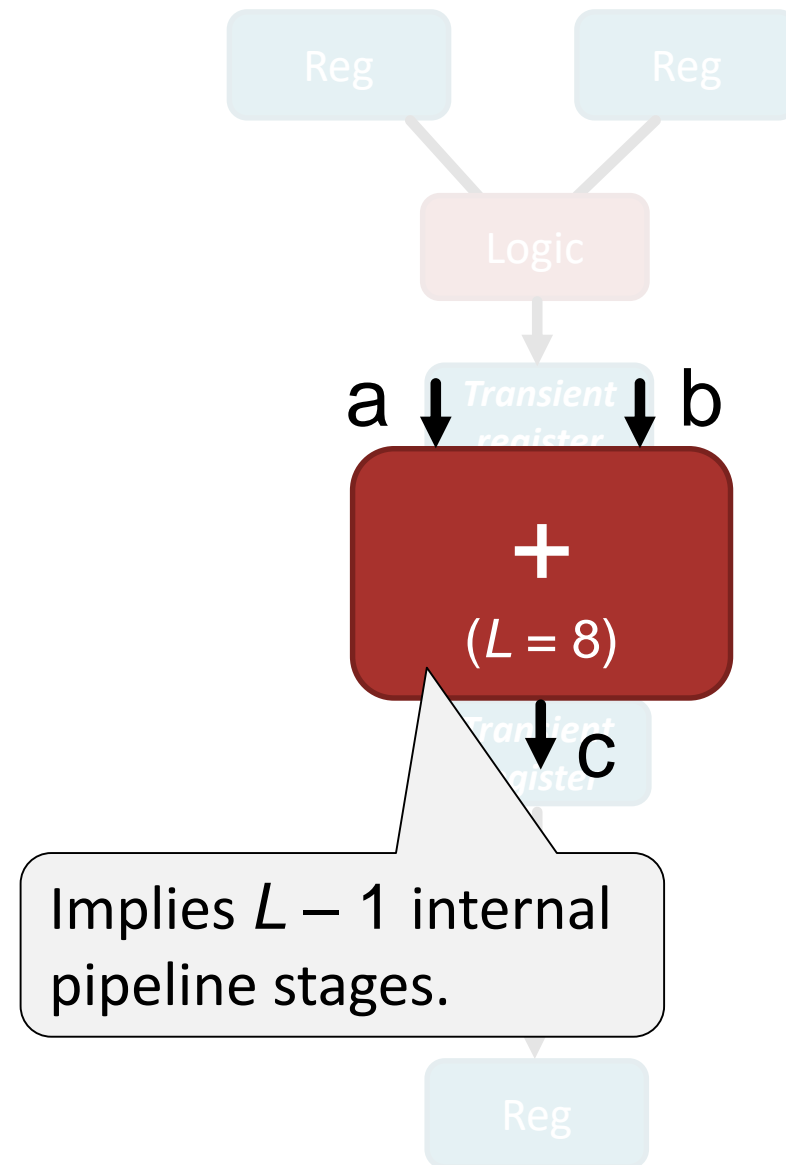
Logic too deep!

1 cycle

*Transient register*

...

Reg

float c = a + b;

Nakayama, T. *Hardware arrangement for floating-point addition and subtraction,*
1993, US Patent 5,197,023.

6

# Our point of view
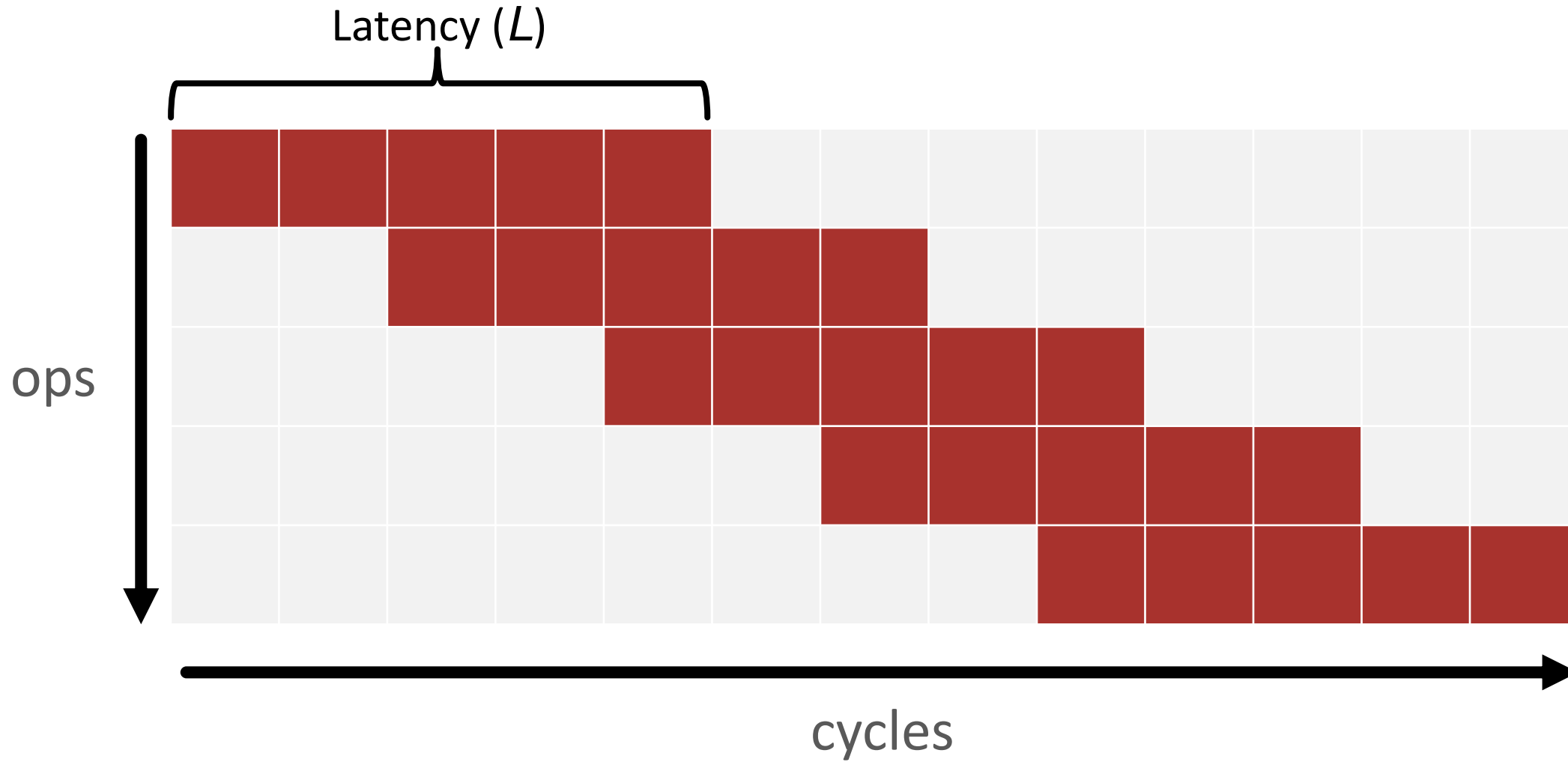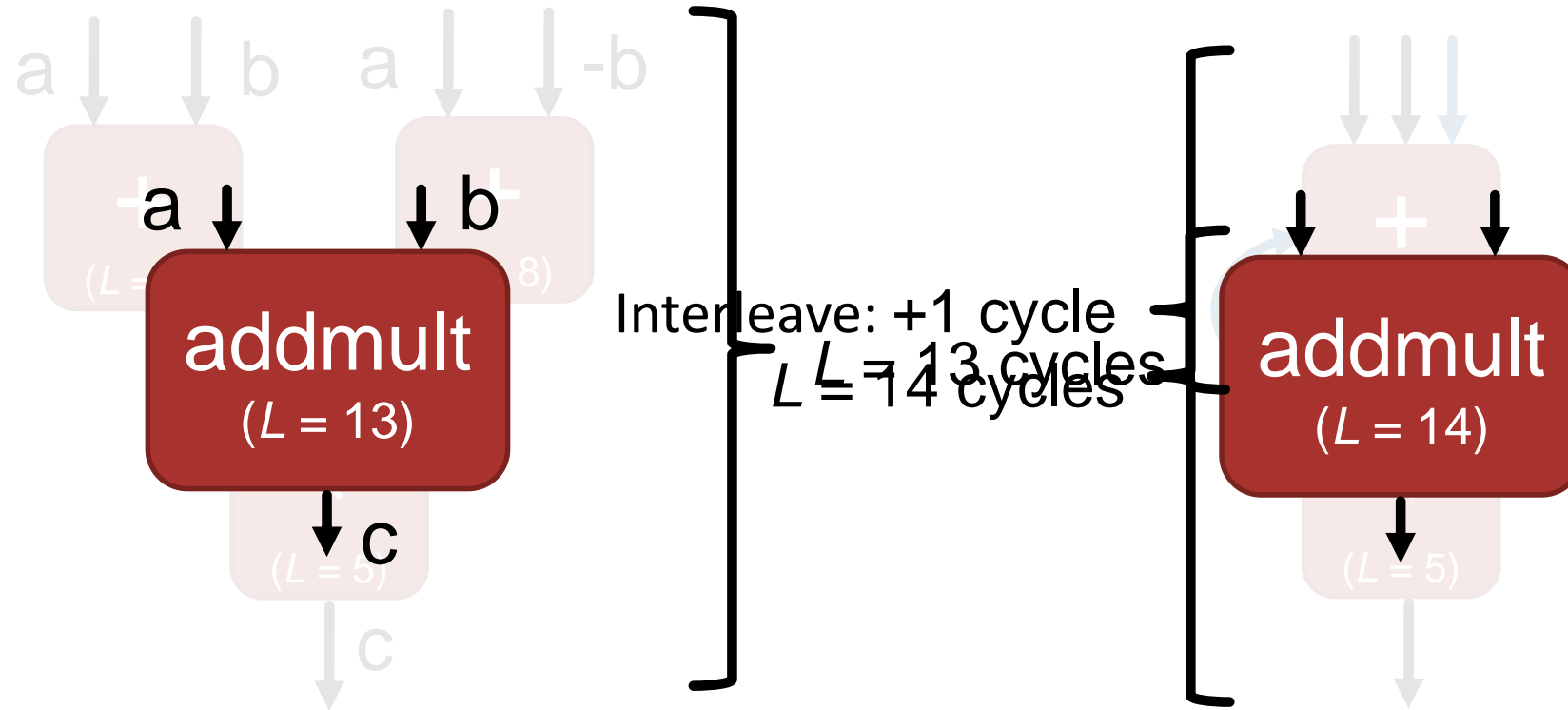
In HLS, we treat **pipelines**.

```
float c = a + b;
```

a Transient register b

**+**

*(L = 8)*

c

Implies *L* – 1 internal pipeline stages.

Reg    Reg

Logic

Reg

# Pipelines

# Multiple floating point operations

a ↓    ↓ b    a ↓    ↓ -b

a ↓    ↓ b

**addmult**
(*L* = 13)

↓ c

Interleave: +1 cycle
*L* = 13 cycles
*L* = 14 cycles

**addmult**
(*L* = 14)

float c = (a + b) * (a - b);

≥ Two ways to implement this
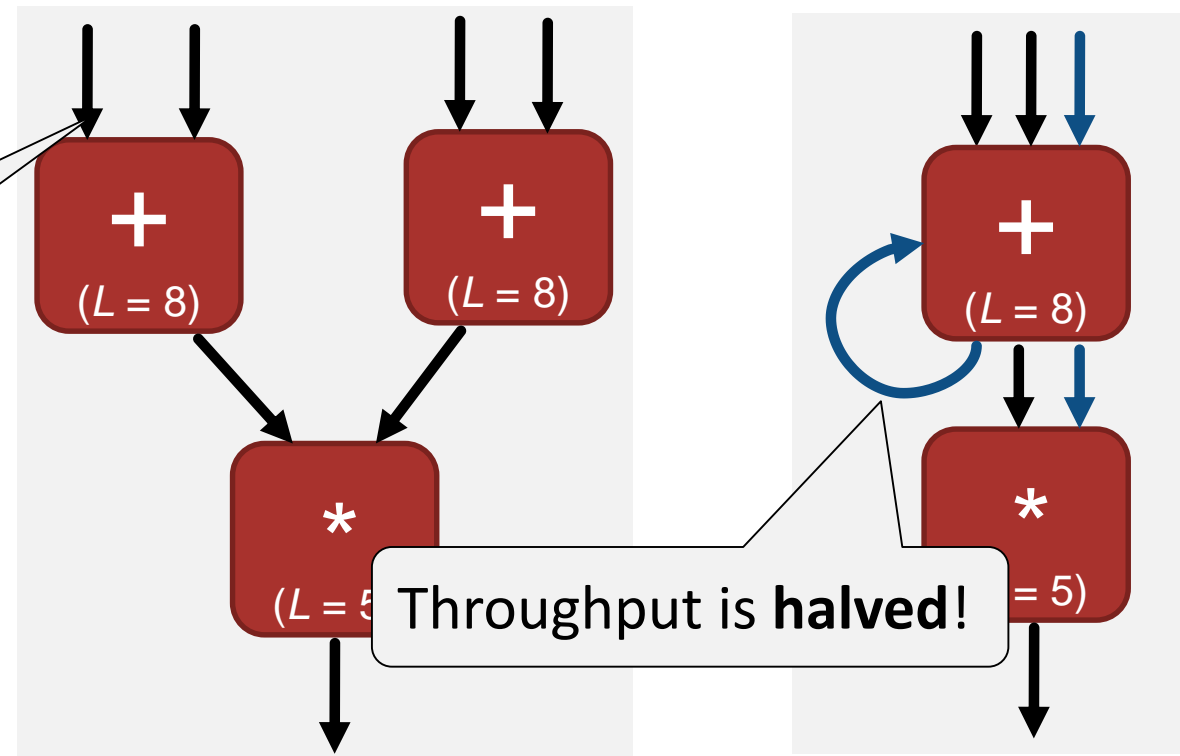
# Initiation interval

In addition to **latency** ($L$), we introduce the property **initiation interval** ("II", here $I$).

**Interpretations**:

1. *No. of cycle* ~~new inputs~~
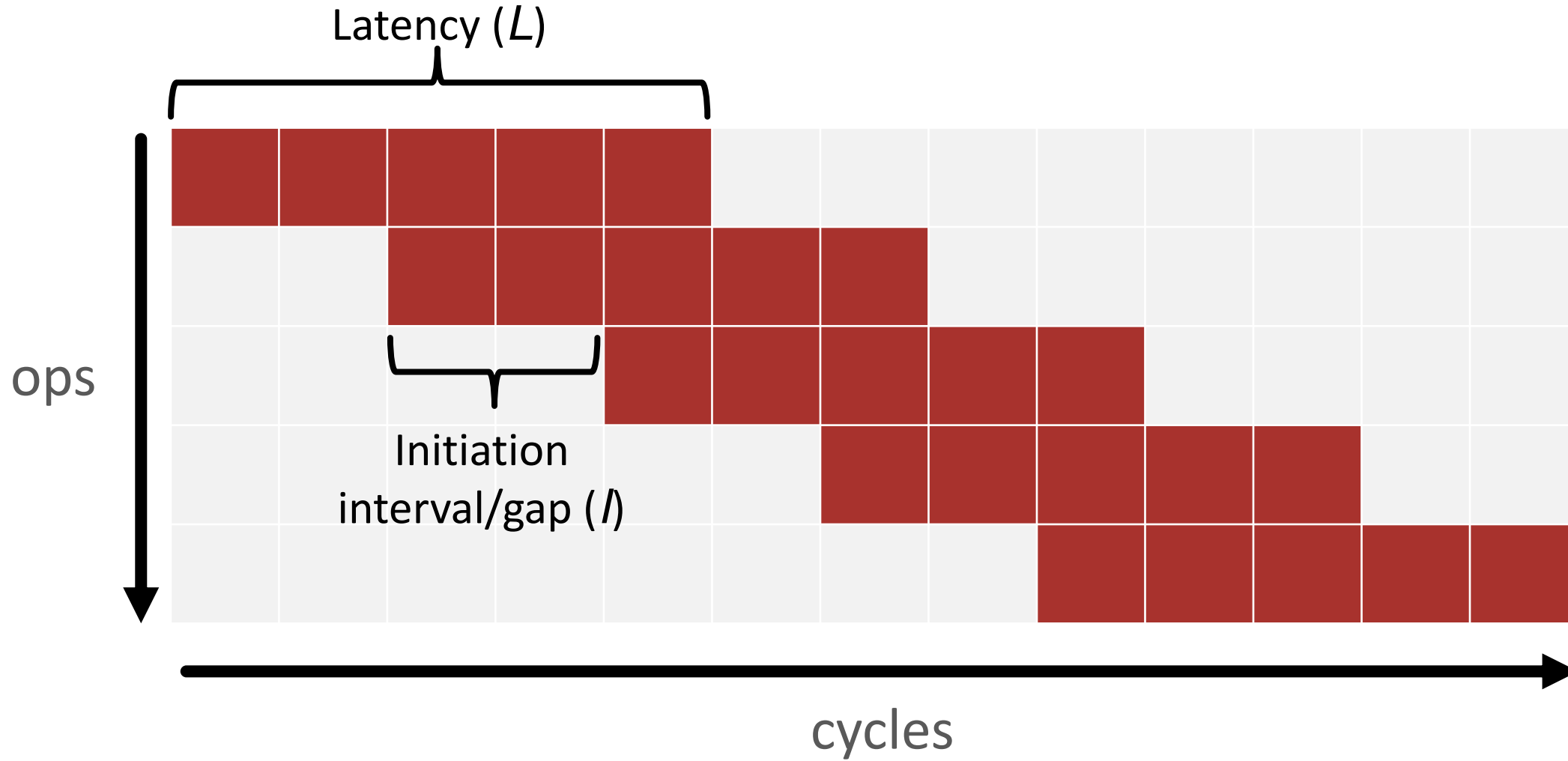
   Can accept all four inputs in parallel

2. *Inverse throughput of the pipeline*

   Throughput is **halved**!

3. *Factor slowdown of your application* ☺



$L$ = 13 cycles
$I$ = 1 cycle
2 adds, 1 mult
**3 op/1 cycle**

$L$ = 14 cycles
$I$ = 2 cycles
1 add, 1 mult
**3 op/2 cycles**

# Pipelines vol. II
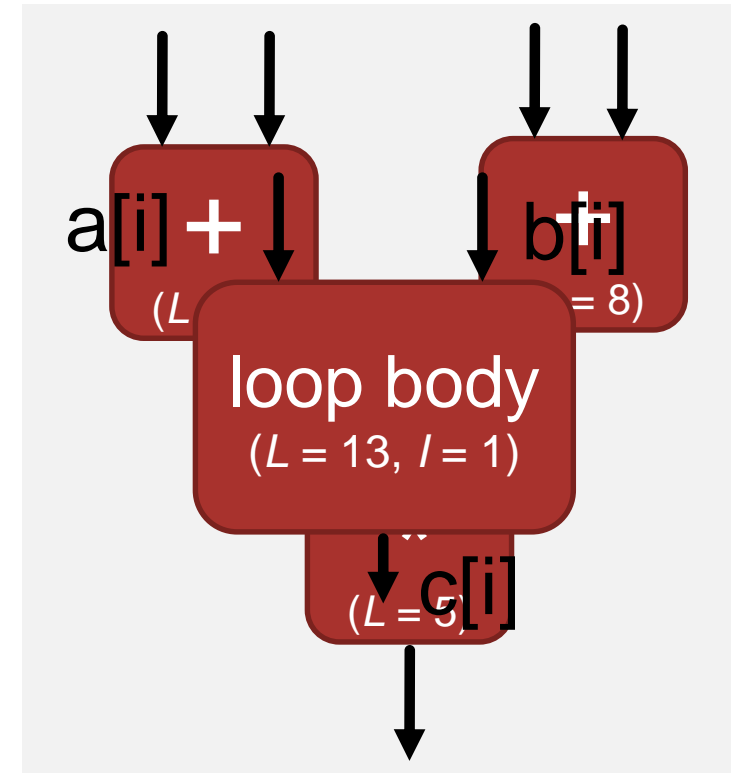
# Adding loops

```
for (int i = 0; i < N; ++i) {
    #pragma HLS PIPELINE II=1
    c[i] = (a[i] + b[i]) *
           (a[i] - b[i]);
}
```

| 1 iteration | 13 + 1 = 14 cycles |
|---|---|
| 10 iterations | 13 + 10 = 23 cycles |
| N iterations | 13 + N cycles |



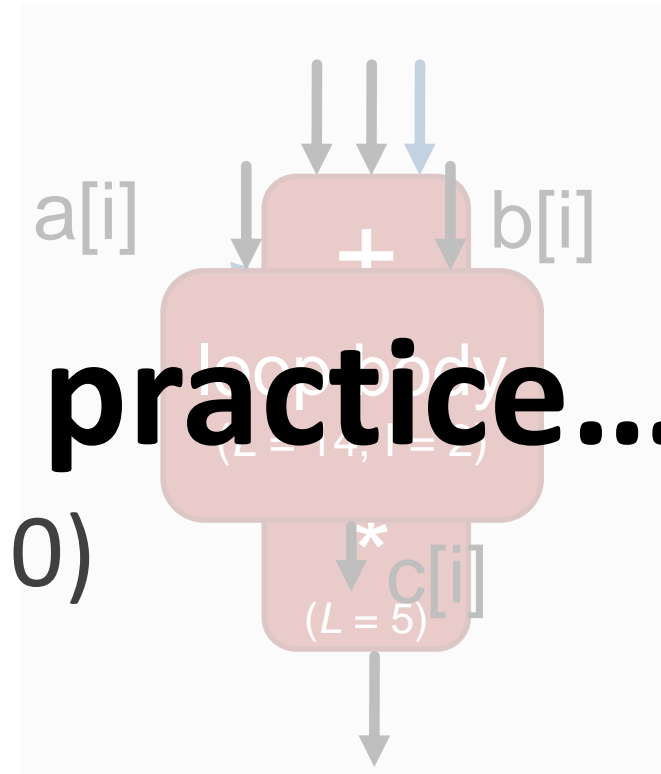Loop iterations affect the runtime **additively**, regardless of body content

# Adding loops

```
for (int i = 0; i < N; ++i) {
    #pragma HLS PIPELINE II=2
    c[i] = (a[i] + b[i]) *
          (a[i] - b[i]);
}
```

Generally:

$$L_{\text{tot}} = L + I \cdot N$$

a[i]    +    b[i]

loop body
(L = 14, II = 2)

*
(L = 5)    c[i]

# Let's see this in practice...
## (example 0)

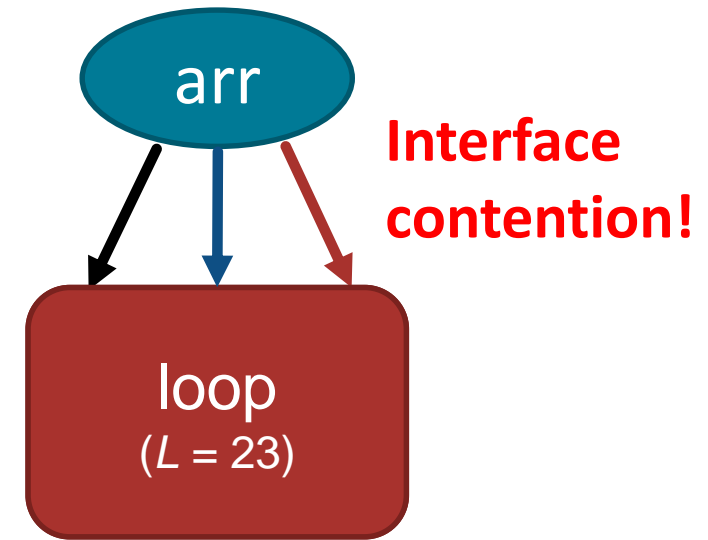Initiation interval paid at *every iteration*

# Pipeline stalls in practice

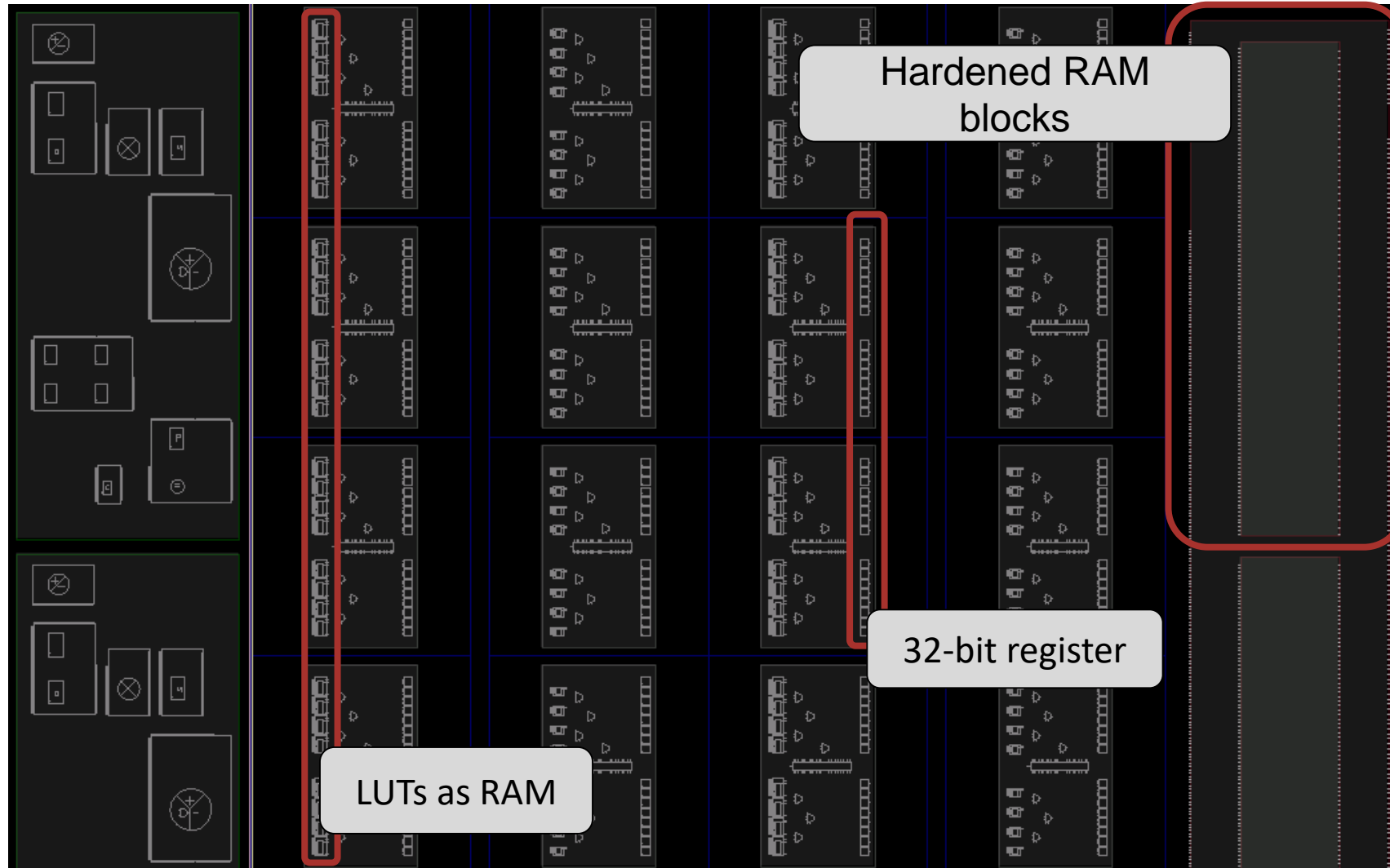Why do we worry so much? Just set *I* = 1...

Increased II is commonly found in the wild:

1. Intra-iteration (now):
   - *Multiple accesses to the same interface*

2. Inter-iteration (later)
   - *Data dependencies*

3. Low throughput requirements (rare)
   - *e.g. input only received every 16 cycles*

```
for (int i = 1; i < N - 1; ++i) {
    #pragma HLS PIPELINE II=1 II=3
    res[i] = 0.3333 *
        (arr[i-1] + arr[i] + arr[i+1]);
}
```
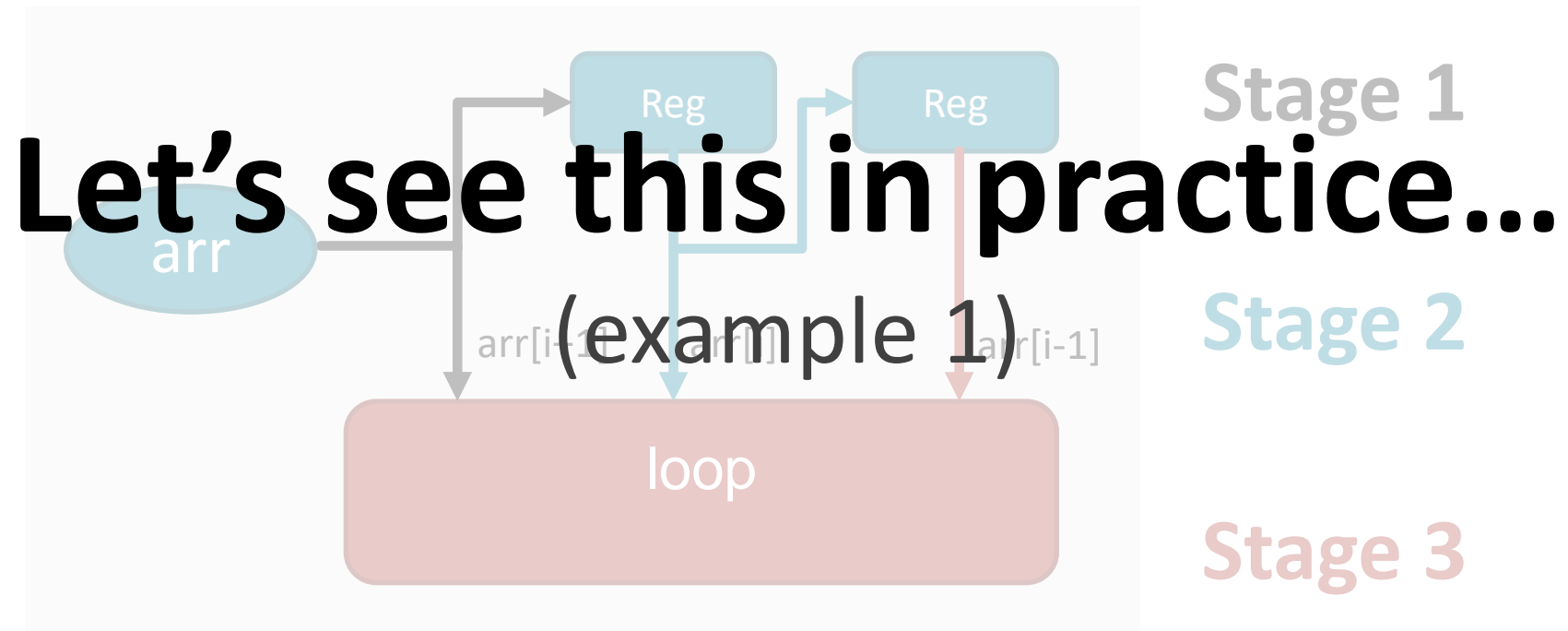
arr

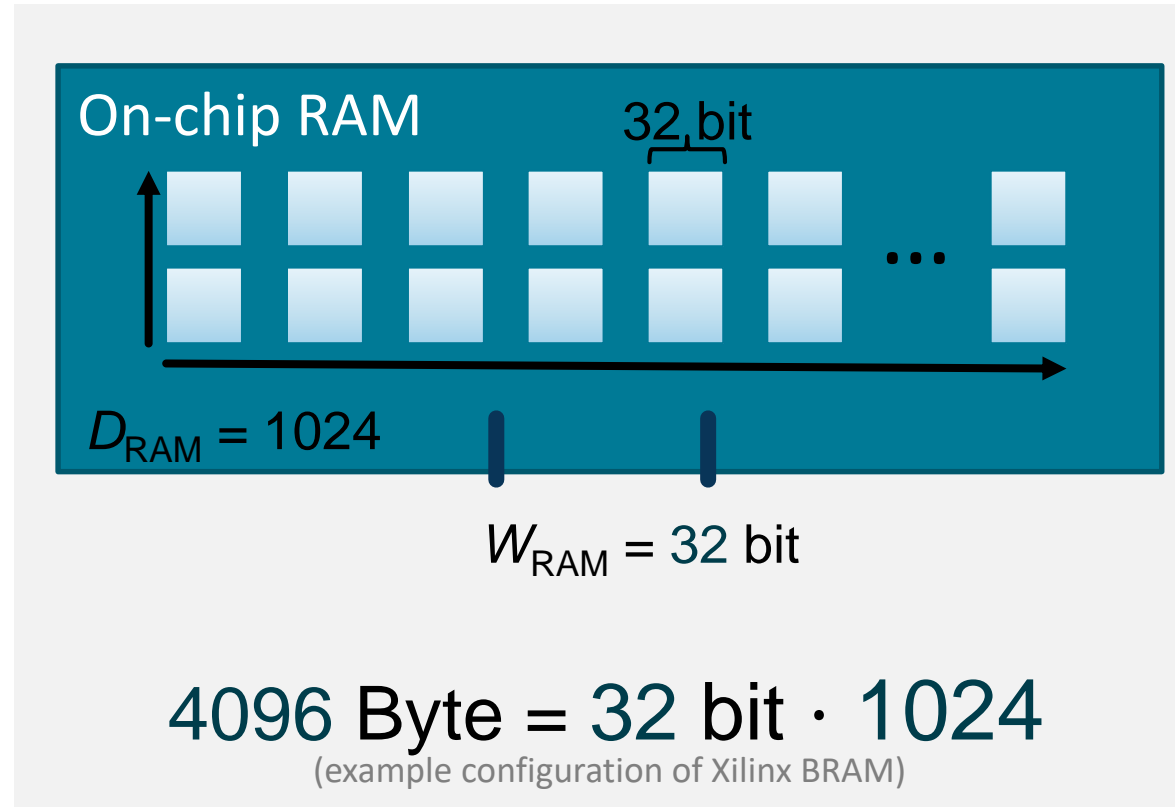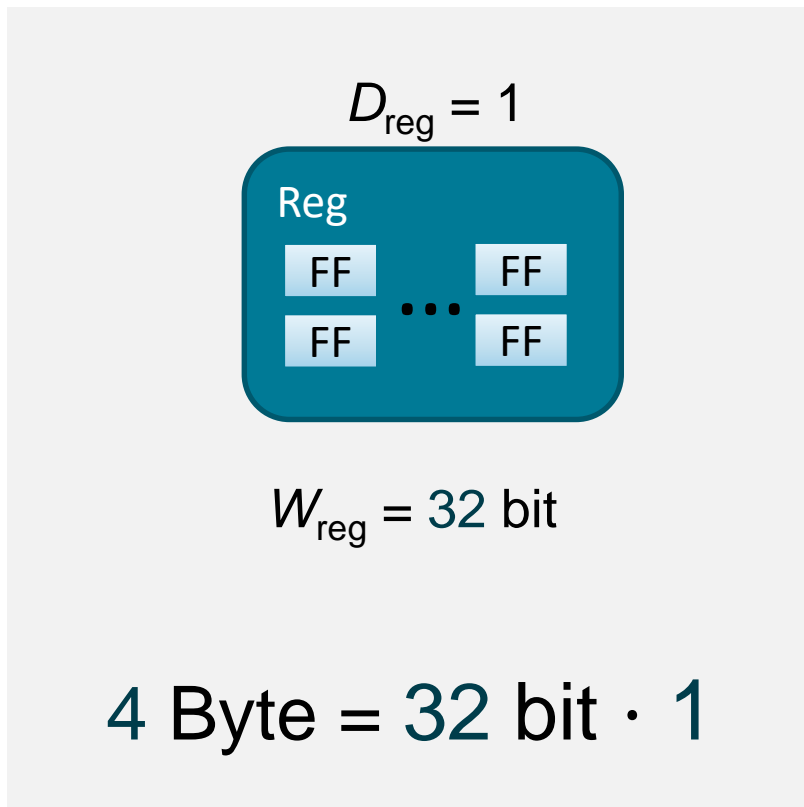**Interface contention!**

loop
(*L* = 23)

# Fast memory everywhere



Hardened RAM blocks

32-bit register

LUTs as RAM

**Inserting registers**

```
for (int i = 1; i < N - 1; ++i) {
  res[i] = 0.3333 *
    (arr[i-1] + arr[i] + arr[i+1]);
}
```

# Let's see this in practice…
## (example 1)

Stage 1

Stage 2

Stage 3

Reg

Reg

arr

loop

# ≥ Two classes of storage

$D_{reg} = 1$

Reg

FF    FF

...

FF    FF

$W_{reg} = 32$ bit

$$4 \text{ Byte} = 32 \text{ bit} \cdot 1$$

On-chip RAM    32 bit

...

$D_{RAM} = 1024$

$W_{RAM} = 32$ bit

$$4096 \text{ Byte} = 32 \text{ bit} \cdot 1024$$

(example configuration of Xilinx BRAM)

Useful to think of memory in terms of **depth** ($D$) and **width** ($W$)

# Buffer depth

**Memory arrangement**

*1D stencil program:*

$$W = 2, \; D = 1$$

*2D row-major:*

$$W = 2, \; D = N$$

# Modified example

# Thinking in width and depth

```
float reg;
```

W = 32 bit
D = 1

```
Vec<float, 4> vec_reg;
```

W = 128 bit
D = 1

```
float buffer[N];
Stream<float> fifo(N);
```
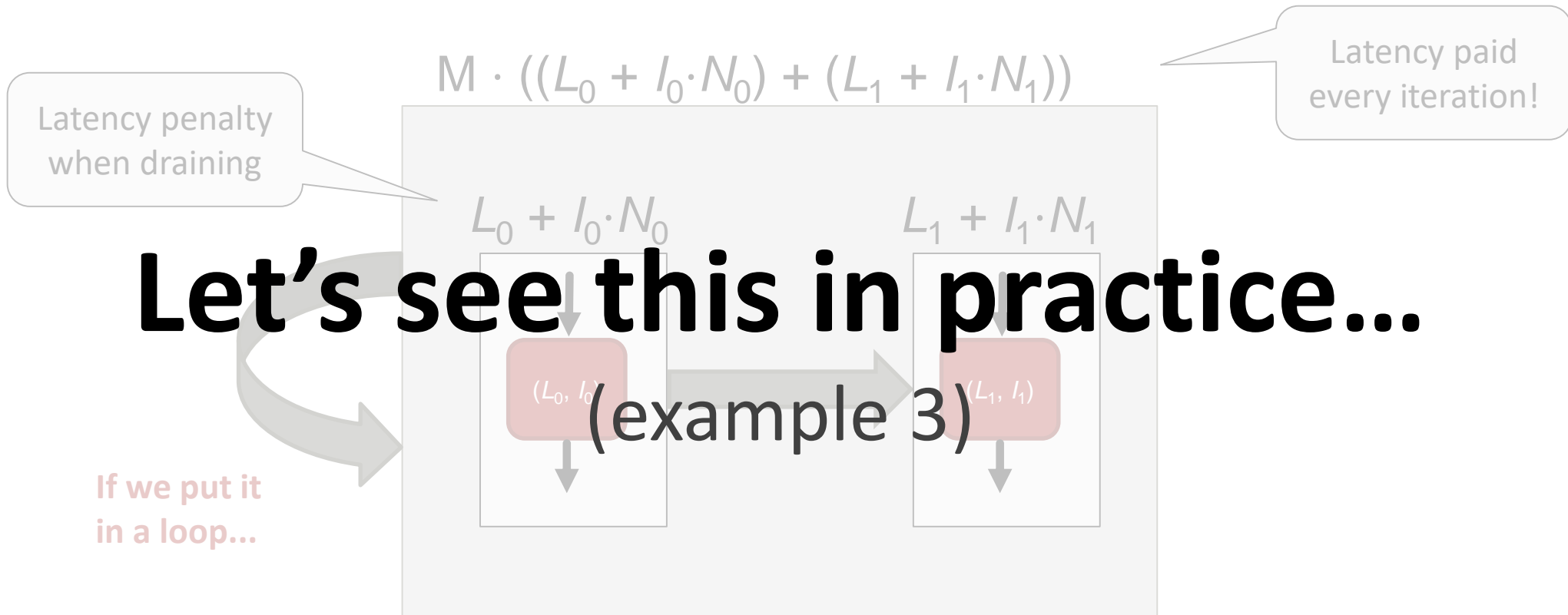
W = 32 bit
D = N

```
Vec<float, 4> vec_buffer[N];
Stream<Vec<float, 4>> fifo(N);
```

W = 128 bit
D = N

# Finite state machine

$L = N + L_0$

arr

**initialize above**

above

FSM Stage 0

*State transition*

$L = N + L_1$

arr

**initialize center**

center

FSM Stage 1

*State transition*

$L = NM + L_2$

arr

**compute**

result

FSM Stage 2

**Sequential** execution of (pipeline-)**parallel** stages

# Flattening

$$M \cdot ((L_0 + I_0 \cdot N_0) + (L_1 + I_1 \cdot N_1))$$

$L_0 + I_0 \cdot N_0$

$L_1 + I_1 \cdot N_1$

$(L_0, I_0)$

$(L_1, I_1)$

$L_2 + I \cdot MN$

$(L, I)$

FSM moved into the pipeline

We will see an example of this shortly

# Multiple concurrent pipelines

$$(L_0 + (L_1 + (L_2 \cdot N_2) \cdot N_1) \cdot N_0)$$

# Streams
(FIFOs, queues...)

depth *D* buffer

Can be 1, or even 0

Input

*Synchronize*

*W*

*Synchronize*

Output

*W*

Same properties as a buffer (because it is one)

# Processing elements



**Properties *L* and *I* remain**

Connected by streams

$(L_0, I_0)$

$(L_3, I_3)$

$(L_1, I_1)$

Global pipeline

*Expressed in HLS as:*
- OpenCL kernels (Intel)
- Dataflow functions (Xilinx)

$(L_2, I_2)$

$(L_3, I_3)$

# Properties of the global pipeline

$L_{\text{tot},0,\text{out}} = L_{\text{tot},1,\text{in}}$

*What goes in must come out*:
Every stream write needs a corresponding read

$L_0$
$I_{0,\text{in}}$
$I_{0,\text{out}}$

$I_{0,1} = \max(I_{0,\text{out}}, I_{1,\text{in}})$

$I_{1,\text{in}}$
$I_{1,\text{out}}$

*Slow dominates*:
Can only stream with the highest initiation interval (producer/consumer) at each interface

# Let's see this in practice...
## (example 4)

$L_{\text{tot}} = I\,N + L_0 + L_1 + L_3 \approx I\,N$

$I, L_0$
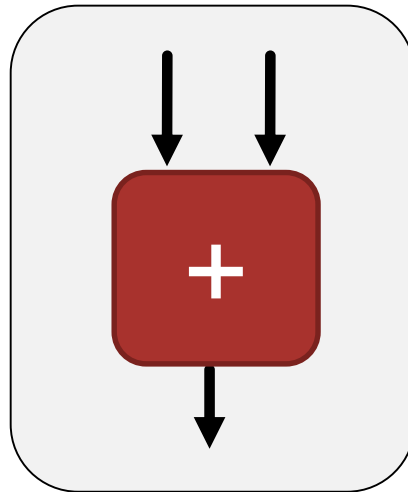
$I, L_1$

$I, L_3$

$I, L_2$

*Depth is "free"*:
In a perfect pipeline for large $N$, the influence of pipeline latency is negligible w.r.t. the total runtime
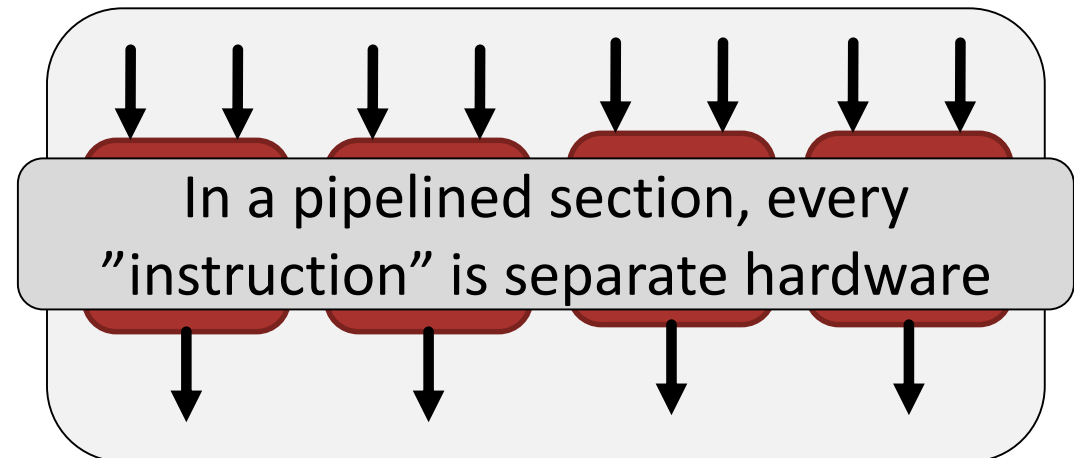
# Unrolling

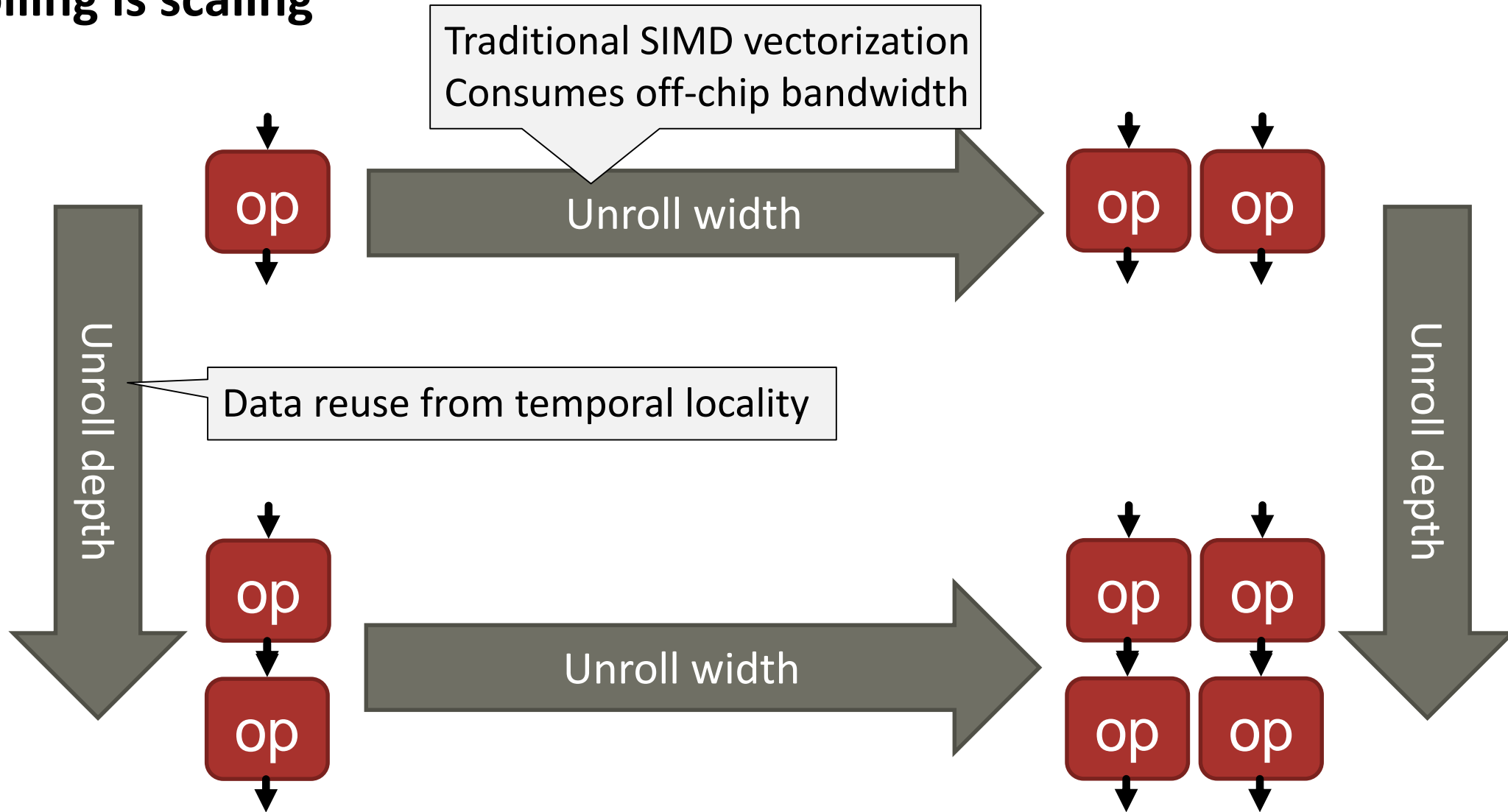Much stronger meaning on FPGA than for an instruction-based architecture.

```
for (int w = 0; w < 4; ++w) {
  #pragma HLS PIPELINE II=1
  res[w] = a[w] + b[w];
}
```

```
#pragma HLS PIPELINE II=1
for (int w = 0; w < 4; ++w) {
  #pragma HLS UNROLL
  res[0] = a[0] + b[0];
  res[1] = a[1] + b[1];
  res[2] = a[2] + b[2];
  res[3] = a[3] + b[3];
}
```



+



In a pipelined section, every "instruction" is separate hardware

# Unrolling is scaling



Traditional SIMD vectorization
Consumes off-chip bandwidth

Unroll width

Unroll depth

Data reuse from temporal locality

op

op op

op

op

op op

op op

Unroll width

Unroll depth

# Matrix-matrix multiplication

```
for (int n = 0; n < N; ++n)
   for (int p = 0; p < P; ++p)
      for (int m = 0; m < M; ++m)
         C[n,p] += A[n,m] * B[m,p];
```
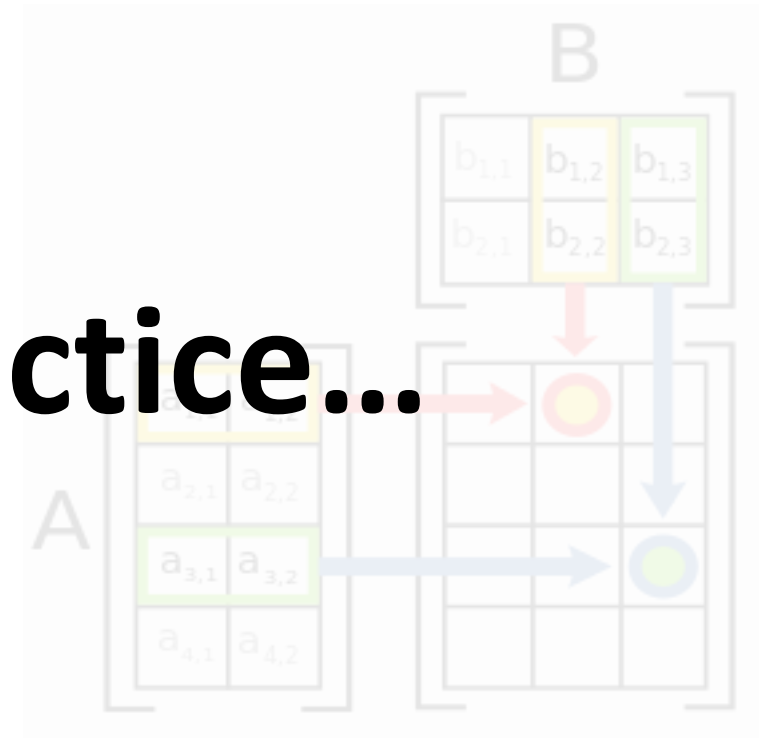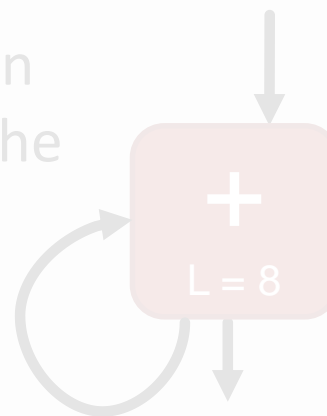
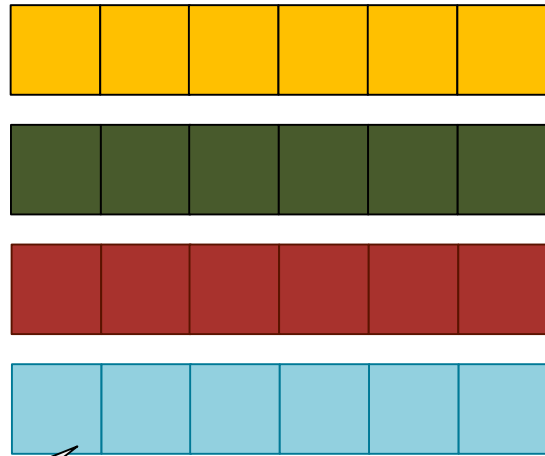# Let's see this in practice…
## (example 5)

Our intuition of temporal locality does not work here!

Every iteration depends on the previous:

**+**

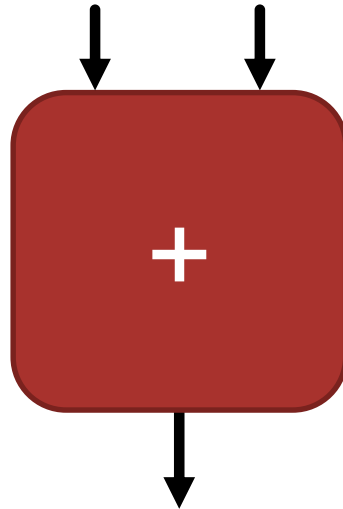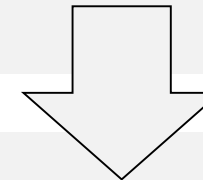L = 8

*l* = 8

# Solving loop-carried data dependencies

```
for (int i = 0; i < N; ++i) {
    float acc = 0;
    for (int j = 0; j < M; ++j) {
        acc += ...
    }
    out[i*M+j] = acc;
```

**Transpose** the iteration space.

Comes at the cost of buffer space.

```
float acc[N];
for (int j = 0; j < M; ++j) {
    for (int i = 0; i < N; ++i) {
        acc[i] += ...
    }
}
for (int i = 0; i < N; ++i) {
    out[i*M+i] = acc[i];
}
```

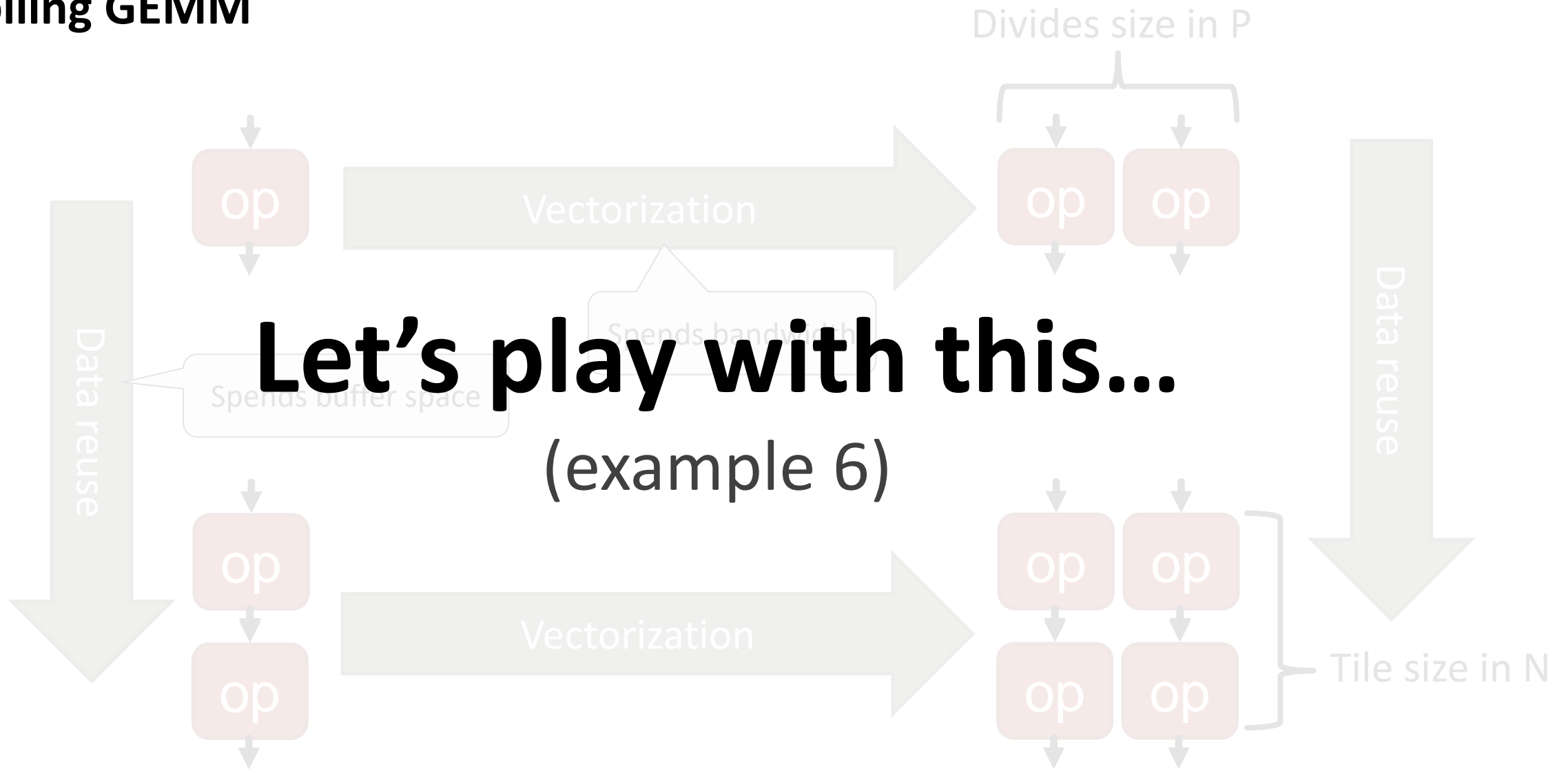Only needs to be larger than the latency

# Locality in the program

```
for (int n = 0; n < N; ++n) {
  float acc[P];

  for (int k = 0; k < K; ++k) {
    const auto a = A[n*K + k];

    for (int m = 0; m < M; ++m) {
      #pragma HLS PIPELINE II=1
      const float prev = (k == 0) ? 0 : acc[m];
      acc[m] = prev + a * B[k*M + m];
    }
  }
  // ...
}
```
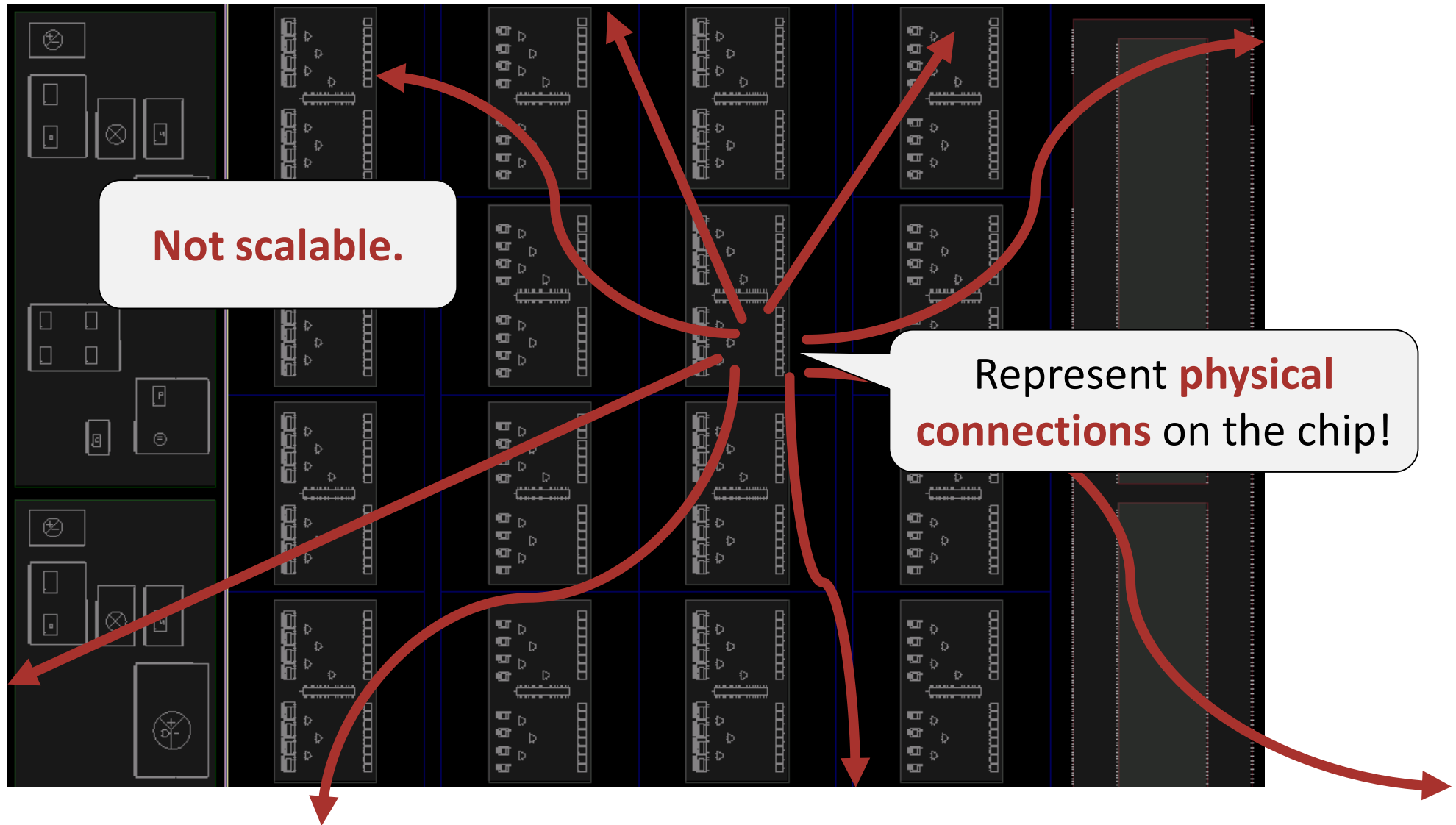
**Temporal locality:** Reused $P$ times. Load more of these and tile $N$!
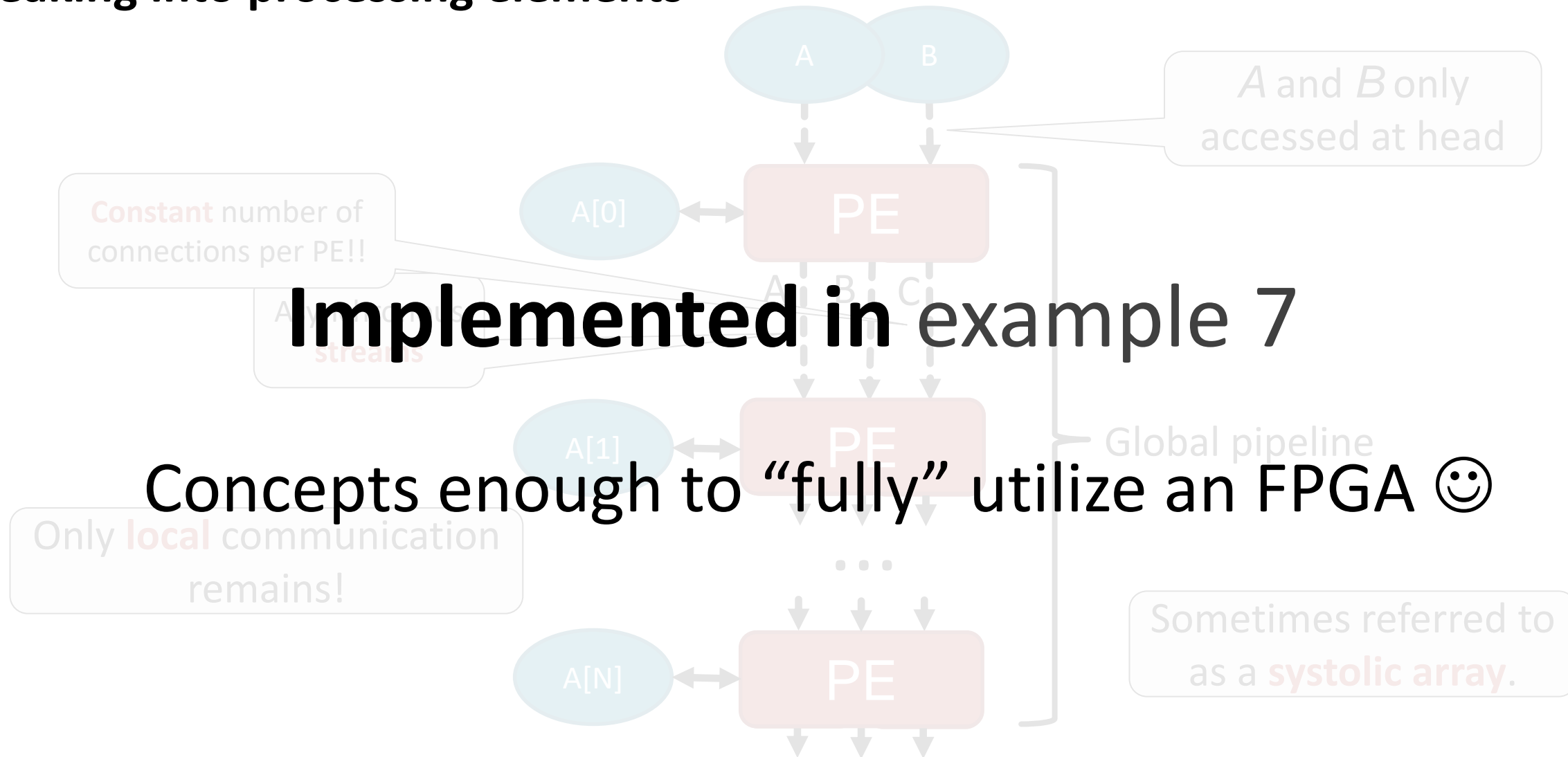
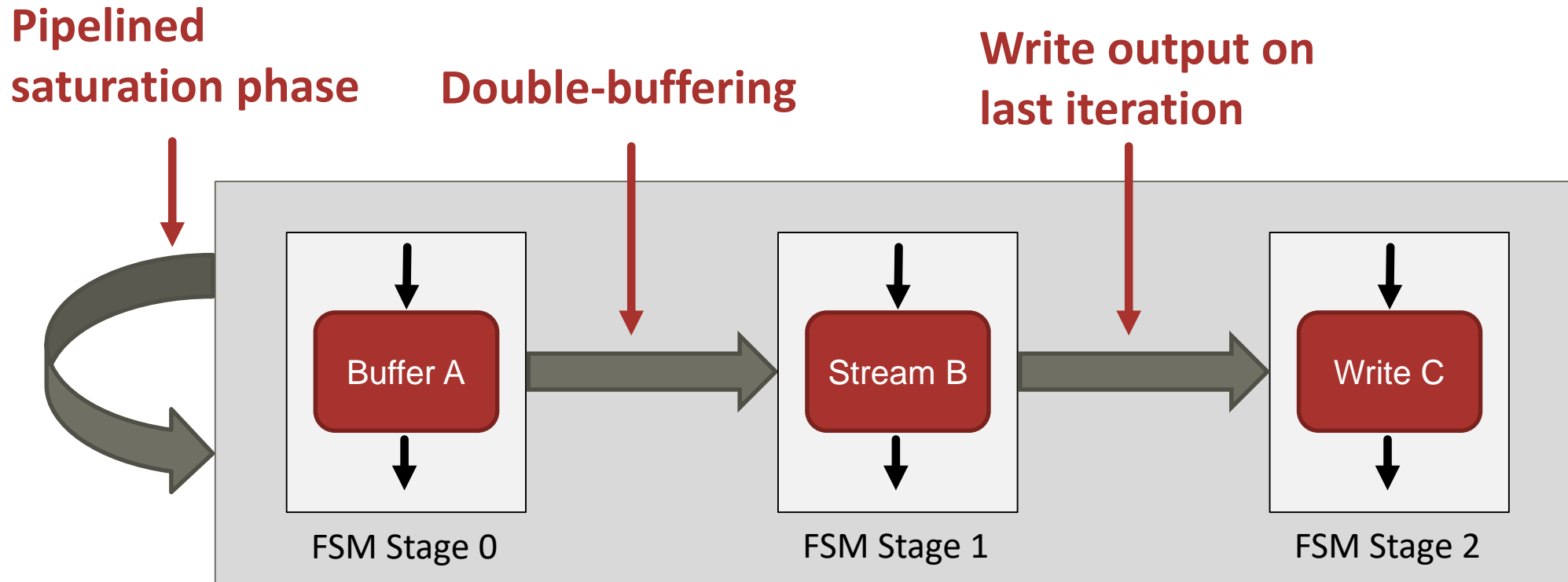**Spatial locality:** Vectorizable.

**Unrolling GEMM**

Divides size in P

op

Vectorization

op op

Data reuse

Spends bandwidth

Data reuse

# Let's play with this…
## (example 6)

Spends buffer space

op

op

Vectorization

op op

op op

Tile size in N

# Fanout issue



**Not scalable.**

Represent **physical connections** on the chip!

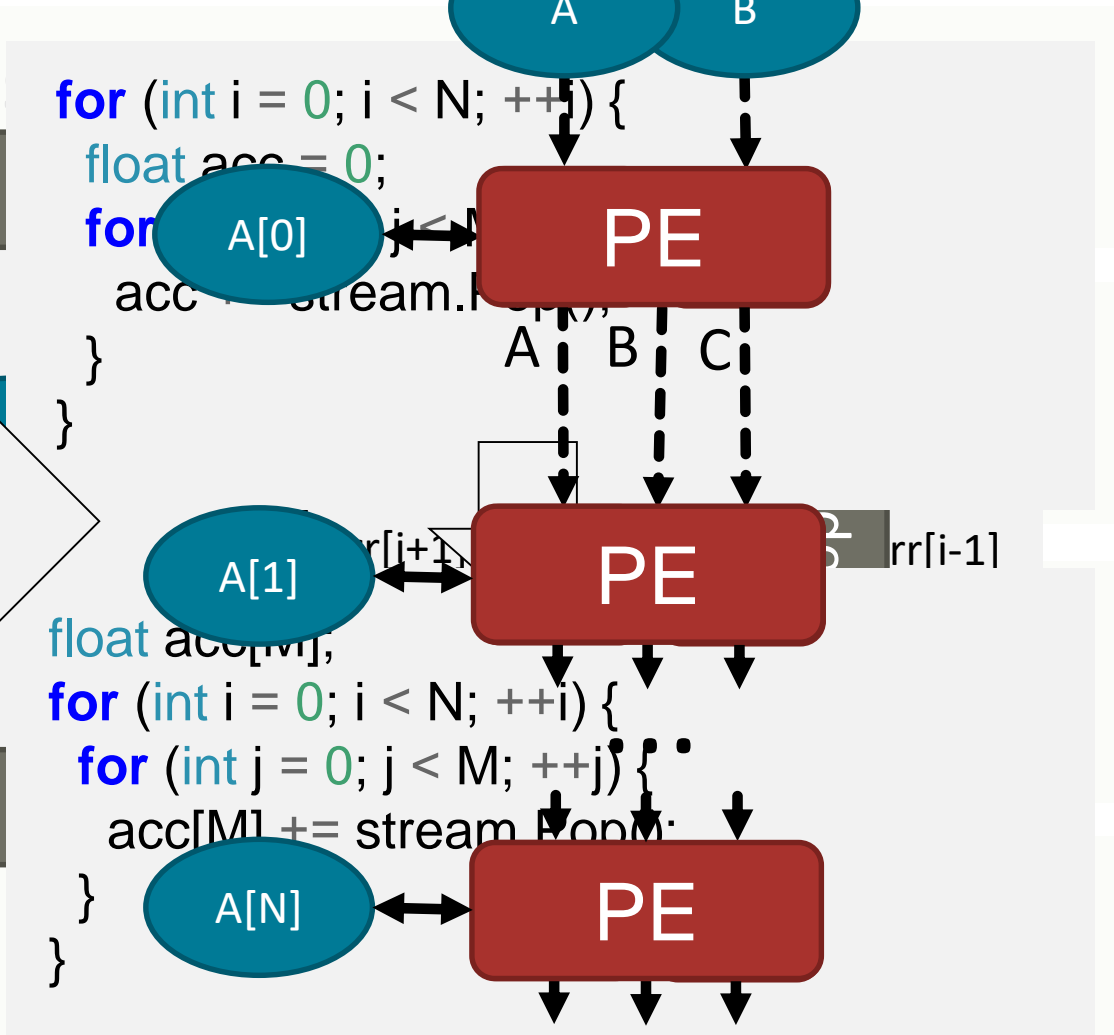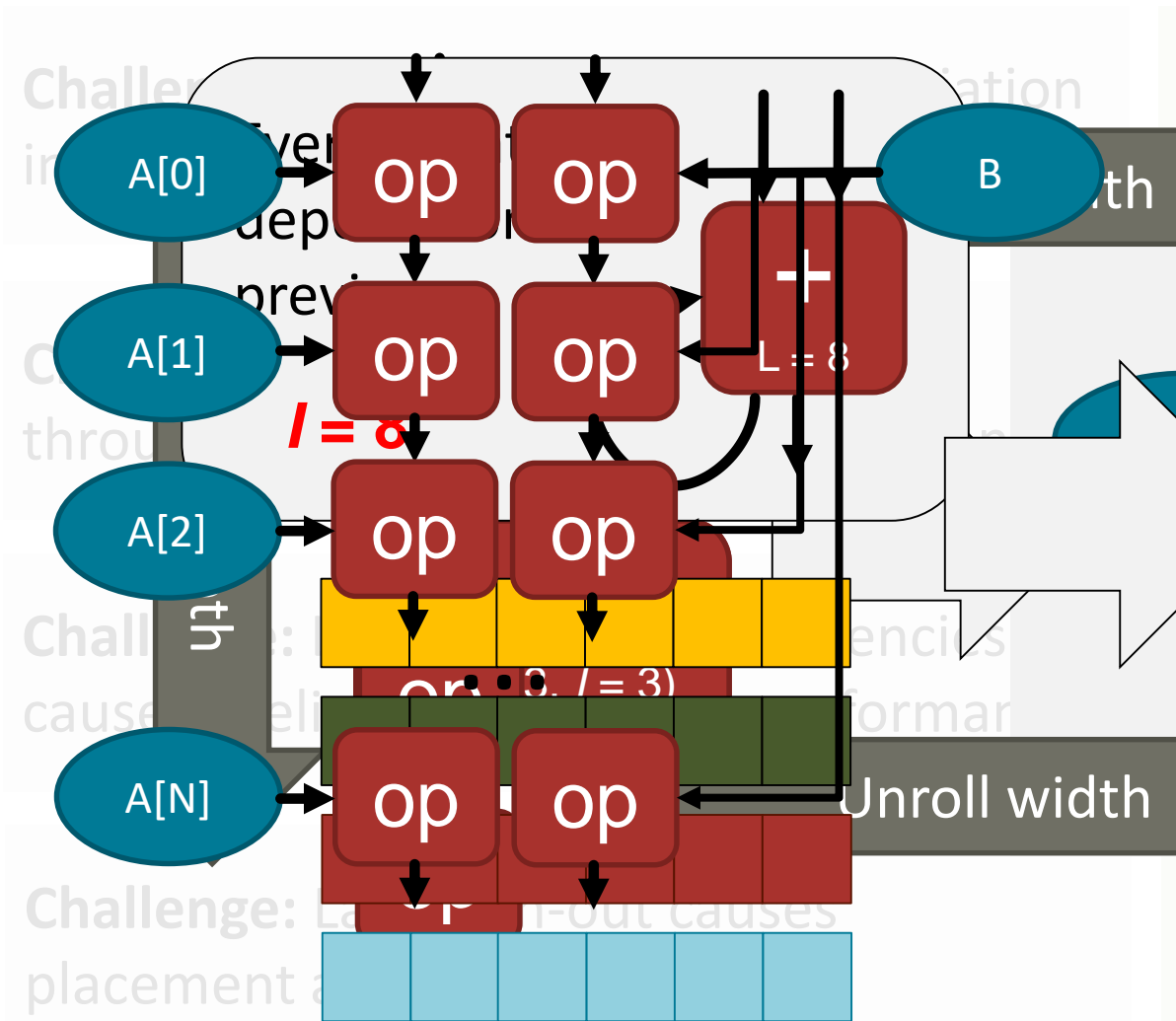**Breaking into processing elements**

# Implemented in example 7

## Concepts enough to "fully" utilize an FPGA ☺

# Remaining optimization potential



Fully optimized implementation at: https://github.com/spcl/gemm_hls

# Summary
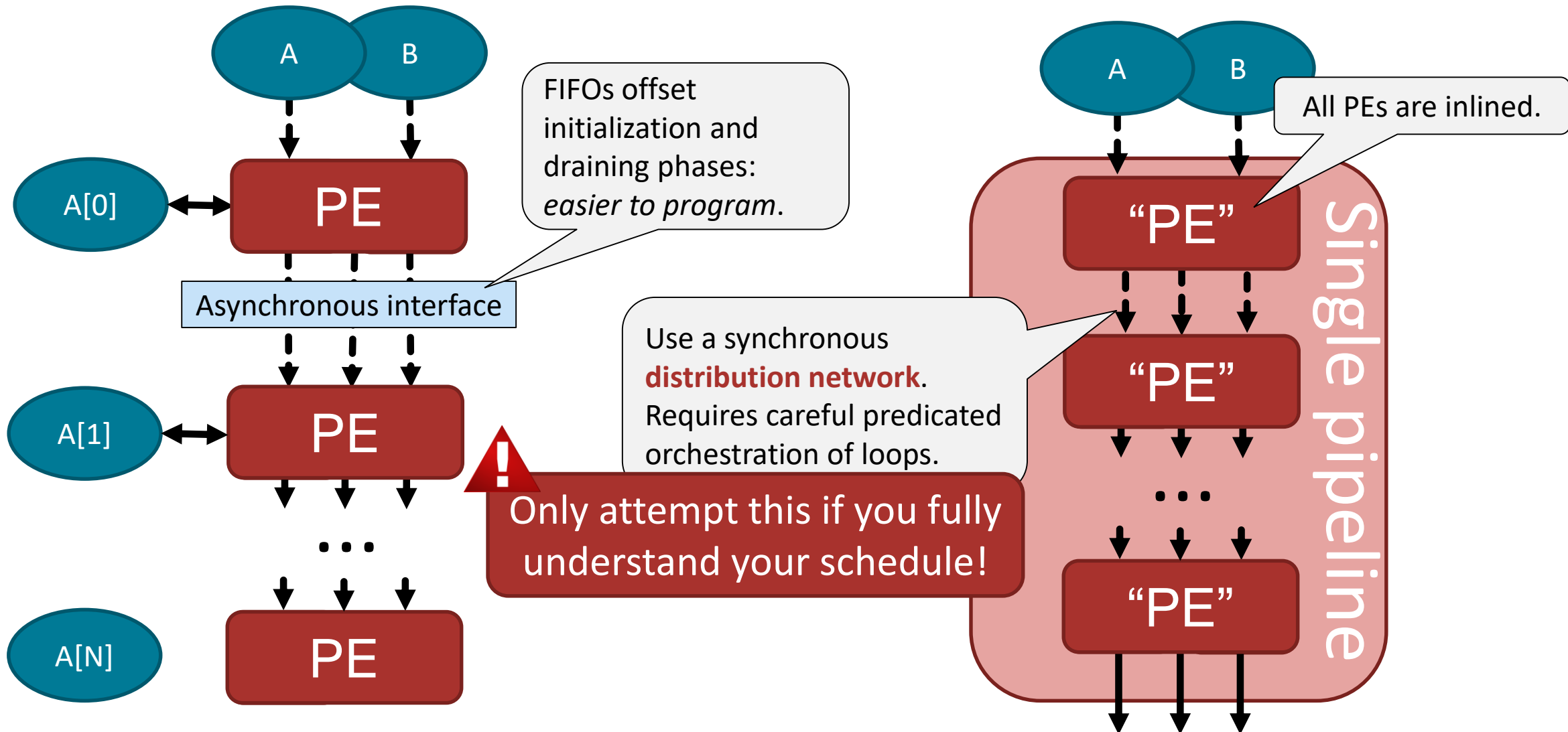
# Thank you for your attention ☺
Reach out at: definelicht@inf.ethz.ch

For a detailed description of HLS
transformations, see:

*Transformations of High-Level Synthesis Codes*
*for High-Performance Computing*
https://arxiv.org/abs/1805.08288

# Two flavors of systolic arrays

# Approaching a new problem

# When HDL should be involved

```
always @(posedge ACLK) begin
    if (~ARESETN | system_reset) begin
        write_state <= WRITE_IDLE;
        write_addr <= 0;
        start_kernel_signal <= 1'b0;
        SC_reset <= 1'b0;
        icap_wr <= 1'b0;
        host2device_wr <= 1'b0;
        host2device <= 32'h0;
        system_reset <= 1'b0;
    end
    else begin
        AWREADY <= 1'b1;
        WREADY <= 1'b0;
        BVALID <= 1'b0;
        system_reset <= 1'b0;
```

```
for (int i = 1; i < N - 1; ++i) {
    for (int j = 0; j < M; ++j) {
        #pragma HLS PIPELINE II=1
```

```
                                              ();
                                              ();
                                              M + j];
```

**Latency critical optimizations**

HDL and HLS can (and do often) happily co-exist!

0.3333;

**Interfacing**

```
        memory_out[i * M + j] = average;
    }
}
```

**Claim**

FPGAs are more energy efficient than GPUs

Compute throughput

At $F_{\text{FPGA}} = F_{\text{GPU}}$, we are looking at a improvement in energy efficiency.

**Verdict**

Only if the performance is competitive

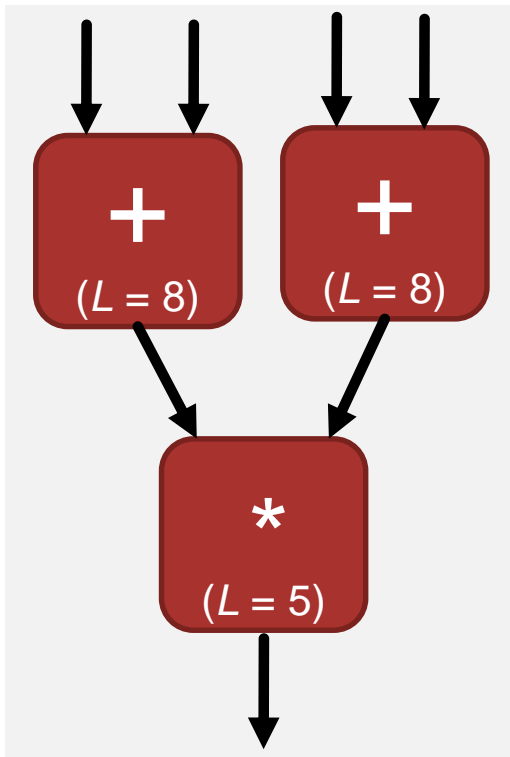For a 10× slowdown, we **lose** a factor 2× in energy efficiency

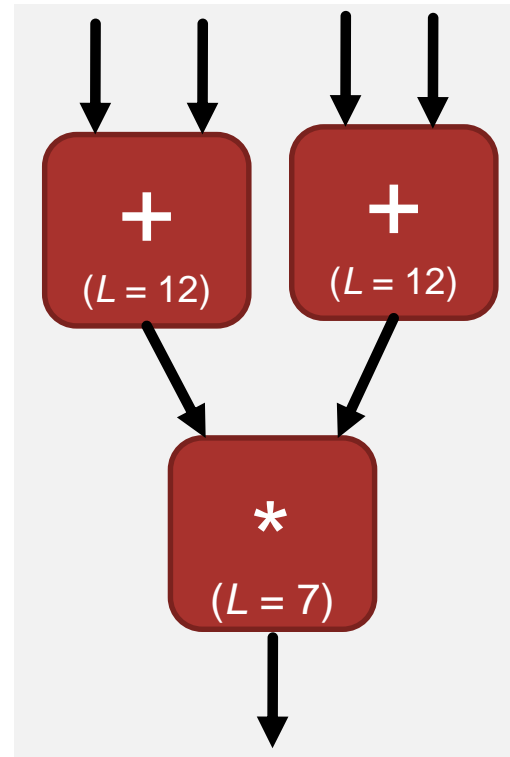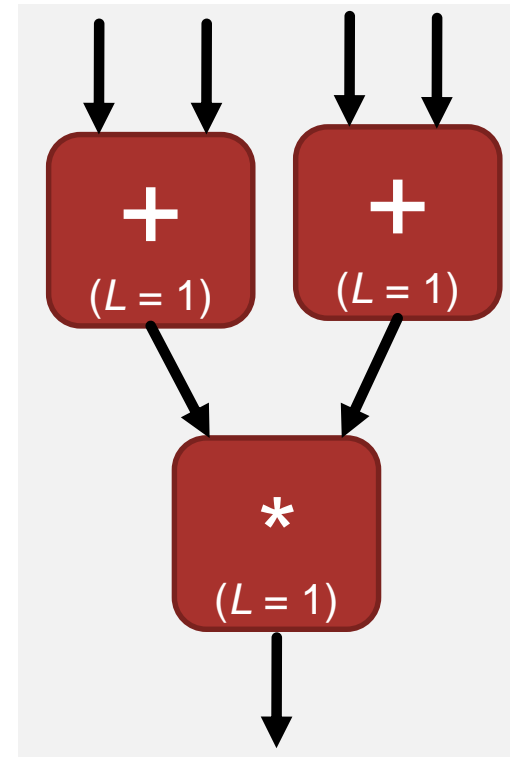Might seem trivial… F matters!

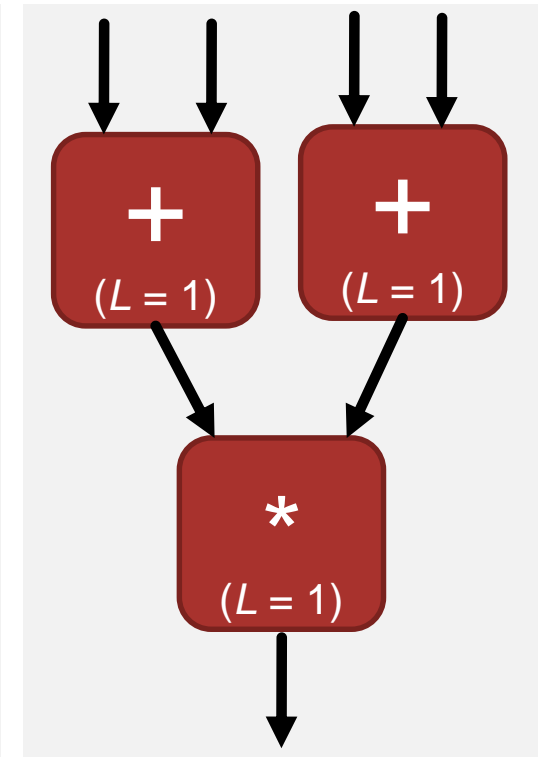# Types on FPGA



float    double    int    fixed_point

Same concepts, different latencies (some problems go away at $L = 1$).