

A Pre-Thanksgiving holiday sampler...



Lecture 16:

Lock-Free Programming Memory Consistency A Bit on Domain-Specific Languages

Parallel Computing
Stanford CS149, Fall 2024

Today: three different topics

- I want to introduce lock-free data structures as an alternative to locks (finishing up the fine-grained synchronization theme from last time)
- The idea of relaxed memory consistency (and why it exists)
- Domain-specific programming languages (and why they are growing in importance in an era of heterogeneous, parallel computing)

Lock-free data structures

(Please see slides from last lecture)

Relaxed memory consistency

Shared memory behavior

- Intuition says loads should return latest value written
 - What is the definition of “latest”?
 - Coherence: only one memory location
 - Consistency: apparent ordering for all locations
 - Order in which memory operations performed by one thread become visible to other threads

- Affects
 - Programmability: how programmers reason about program behavior
 - Allowed behavior of multithreaded programs executing with shared memory
 - Performance: limits HW/SW optimizations that can be used
 - Reordering memory operations to hide latency

Today: who should care

- **Anyone who:**
 - **Wants to implement a synchronization library**
 - **Will ever work a job in kernel (or driver) development**
 - **Seeks to implement lock-free data structures**

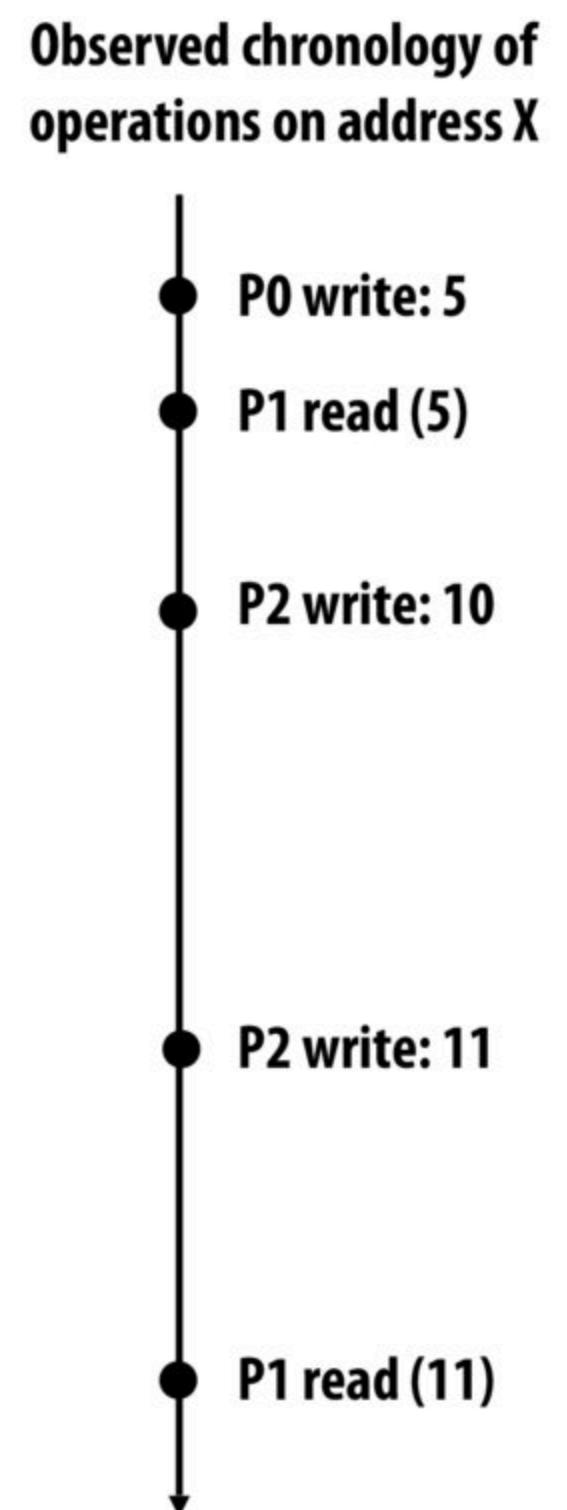
Memory coherence vs. memory consistency

- **Memory coherence** defines requirements for the observed behavior of reads and writes to the same memory location

- All processors must agree on the order of reads/writes to X
- In other words: it is possible to put all operations involving X on a timeline such that the observations of all processors are consistent with that timeline

- **Memory consistency** defines the behavior of reads and writes to different locations (as observed by other processors)

- Coherence only guarantees that writes to address X will eventually propagate to other processors
- Consistency deals with when writes to X propagate to other processors, relative to reads and writes to other addresses



Coherence vs. consistency

(said again, perhaps more intuitively this time)

- **The goal of cache coherence is to ensure that the memory system in a parallel computer behaves as if the caches were not there**
 - Just like how the memory system in a uni-processor system behaves as if the cache was not there
- **A system without caches would have no need for cache coherence**
- **Memory consistency defines the allowed behavior of loads and stores to different addresses in a parallel system**
 - The allowed behavior of memory should be specified whether or not caches are present (and that's what a memory consistency model does)

Memory operation ordering

- A program defines a sequence of loads and stores
(this is the “program order” of the loads and stores)
- Four types of memory operation orderings
 - $W_X \rightarrow R_Y$: write to X must commit before subsequent read from Y *
 - $R_X \rightarrow R_Y$: read from X must commit before subsequent read from Y
 - $R_X \rightarrow W_Y$: read to X must commit before subsequent write to Y
 - $W_X \rightarrow W_Y$: write to X must commit before subsequent write to Y

* To clarify: “write must commit before subsequent read” means:

When a write comes before a read in program order, the write must commit (its results are visible) by the time the read occurs.

Multiprocessor execution

Initially $A = B = 0$

Proc 0

(1) $A = 1$

(2) print B

Proc 1

(3) $B = 1$

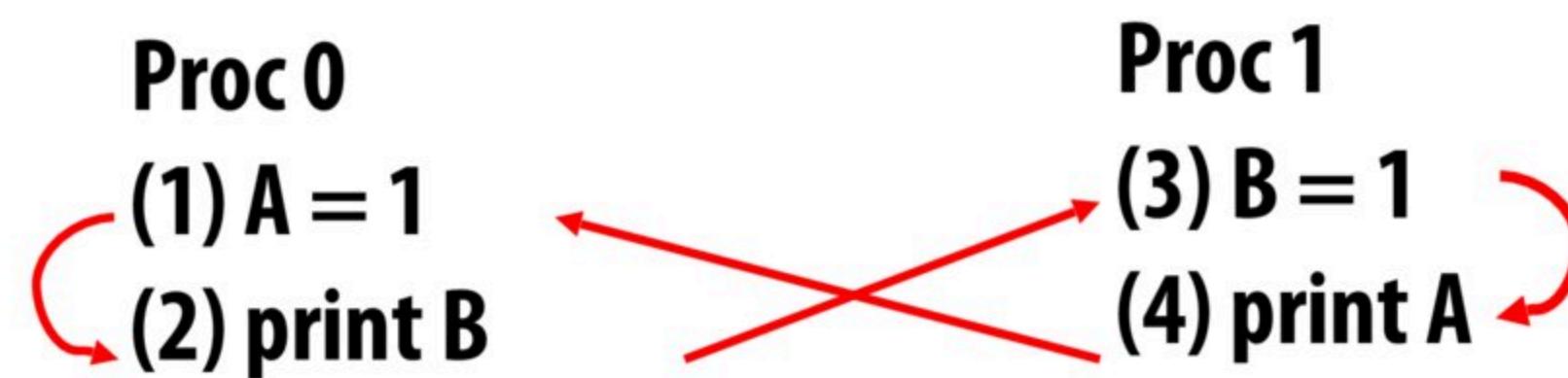
(4) print A

What can be printed?

- “01”?
- “10”?
- “11”?
- “00”?

Orderings That Should Not Happen

Initially $A = B = 0$



- The program should not print “00” or “10”
- A “happens-before” graph shows the order in which events must execute to get a desired outcome
- If there’s a cycle in the graph, an outcome is impossible—an event must happen before itself!

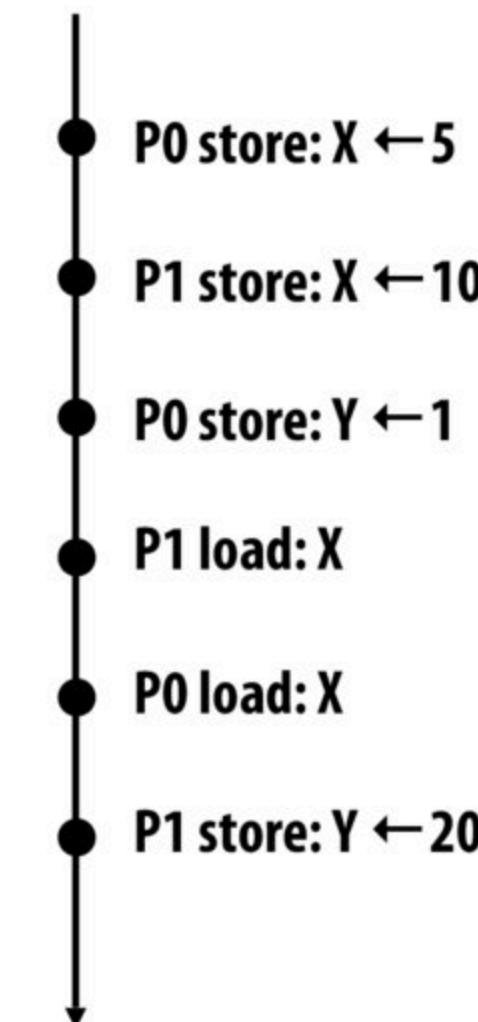
What should programmers expect

■ Sequential Consistency

- Lamport 1976 (Turing Award 2013)
- All operations executed in some sequential order
 - As if they were manipulating a single shared memory
 - Each thread's operations happen in program order

■ A sequentially consistent memory system maintains all four memory operation orderings ($W_X \rightarrow R_Y$, $R_X \rightarrow R_Y$, $R_X \rightarrow W_Y$, $W_X \rightarrow W_Y$)

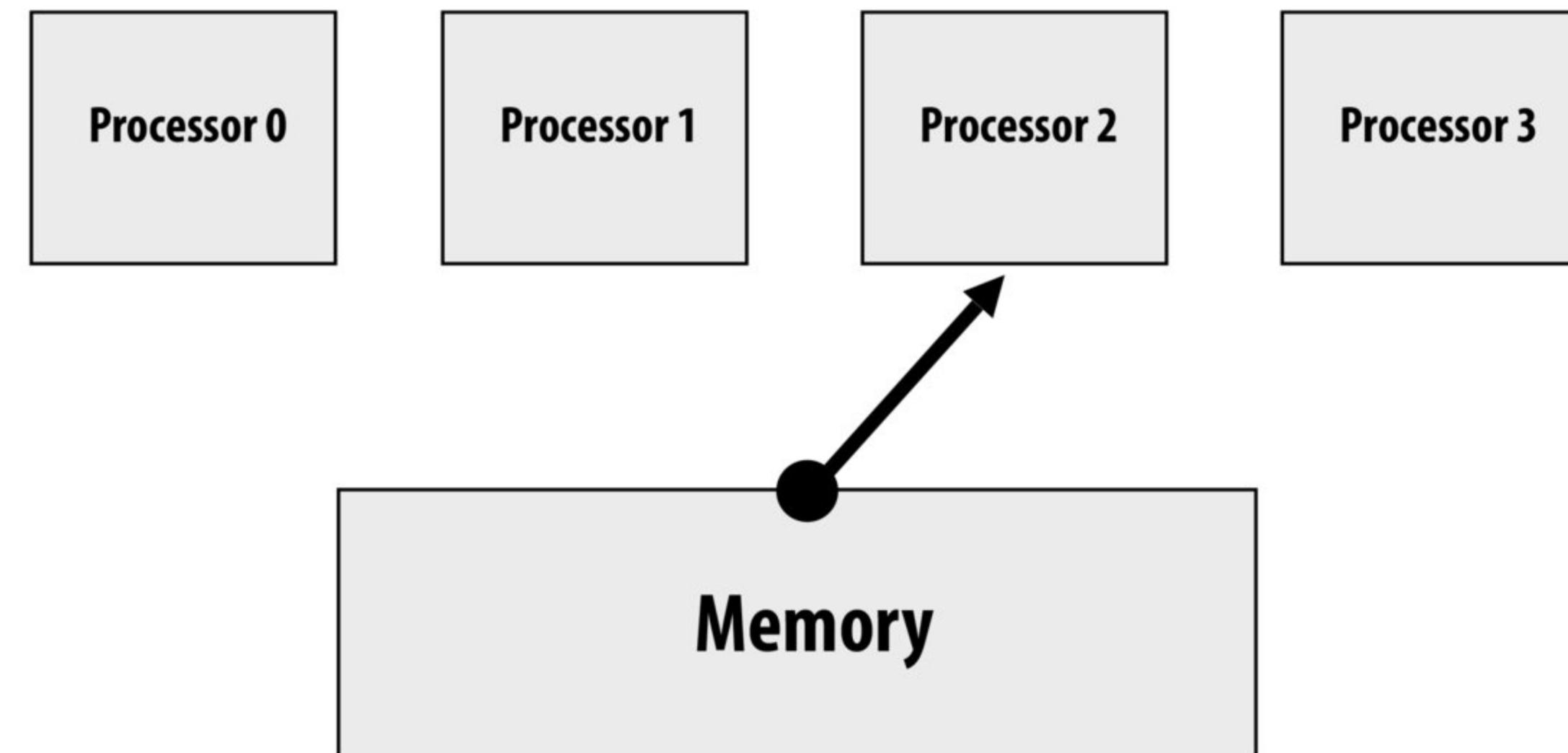
*There is a chronology of all memory operations
that is consistent with observed values*



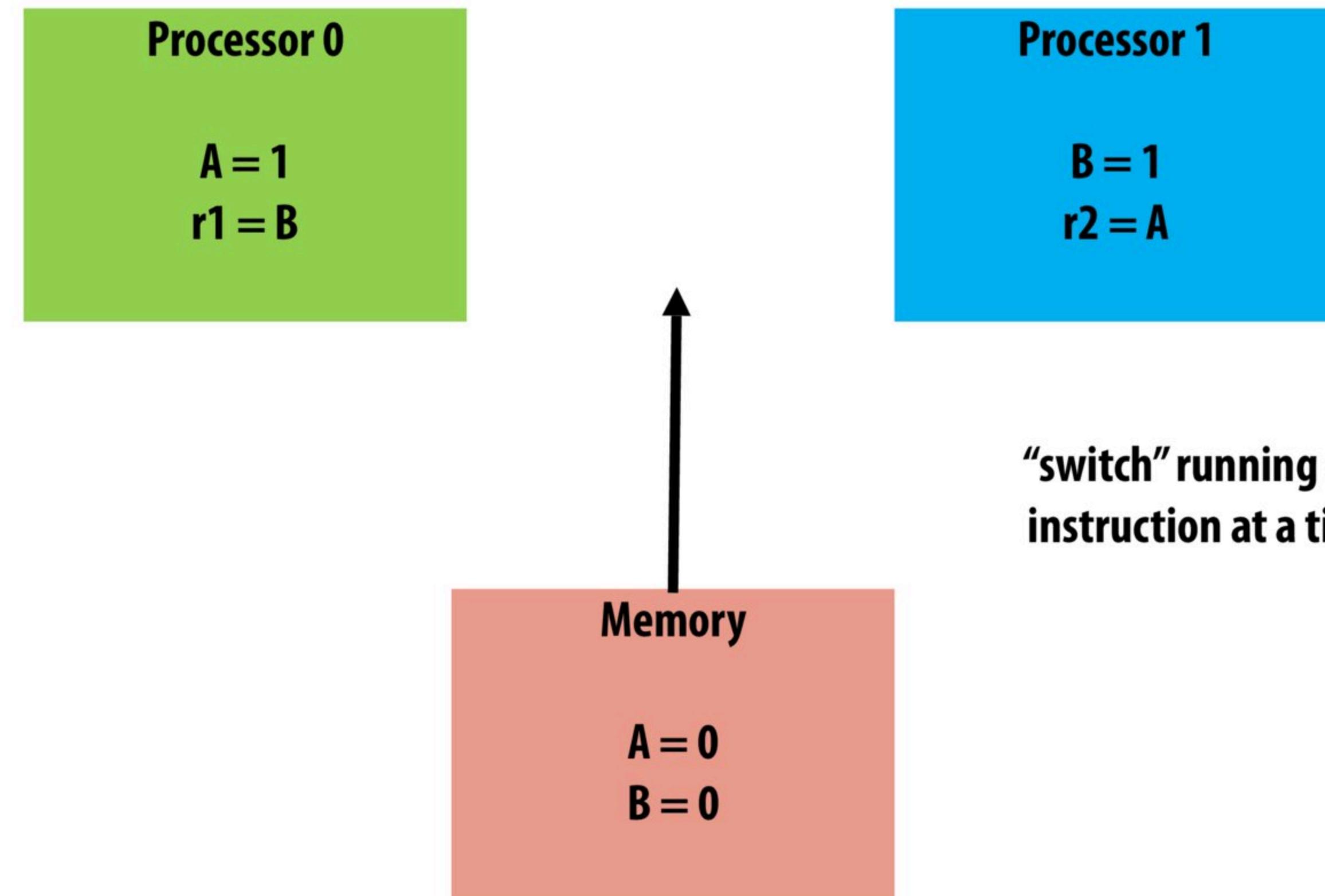
Note:
now timeline lists
operations to addresses
X and Y

Sequential consistency (switch metaphor)

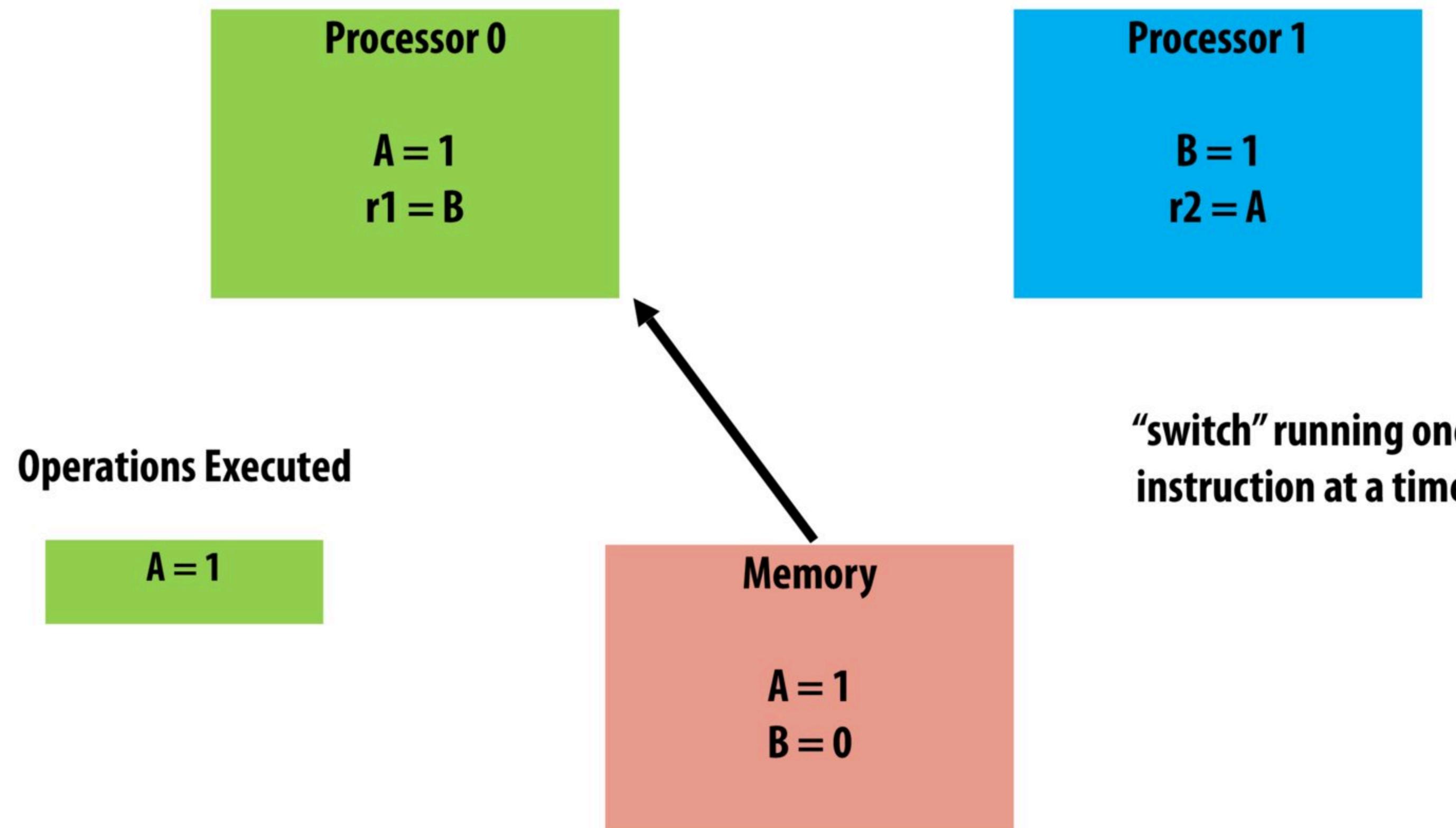
- All processors issue loads and stores in program order
- Memory chooses a processor at random, performs a memory operation to completion, then chooses another processor, ...



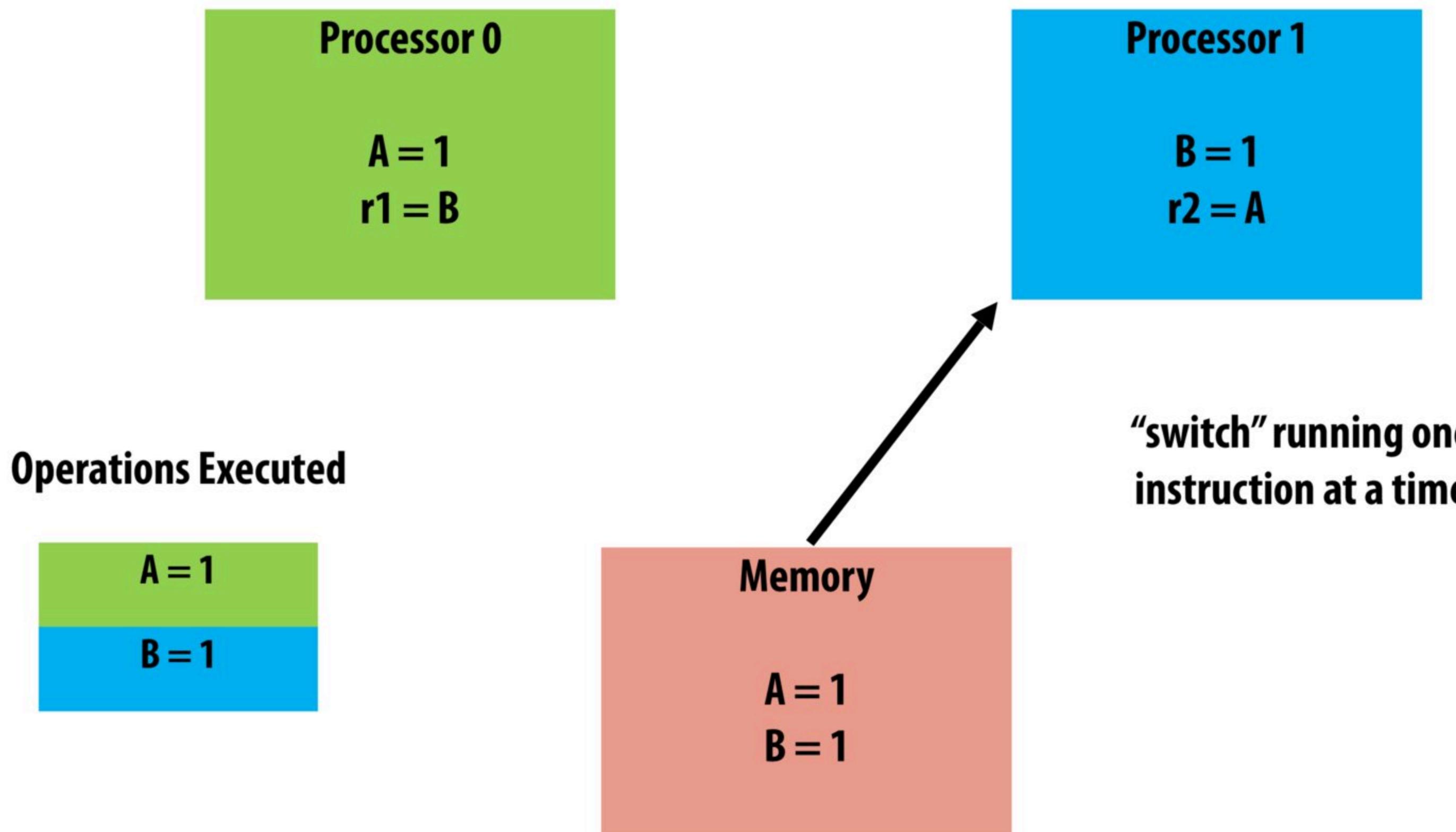
Sequential consistency example



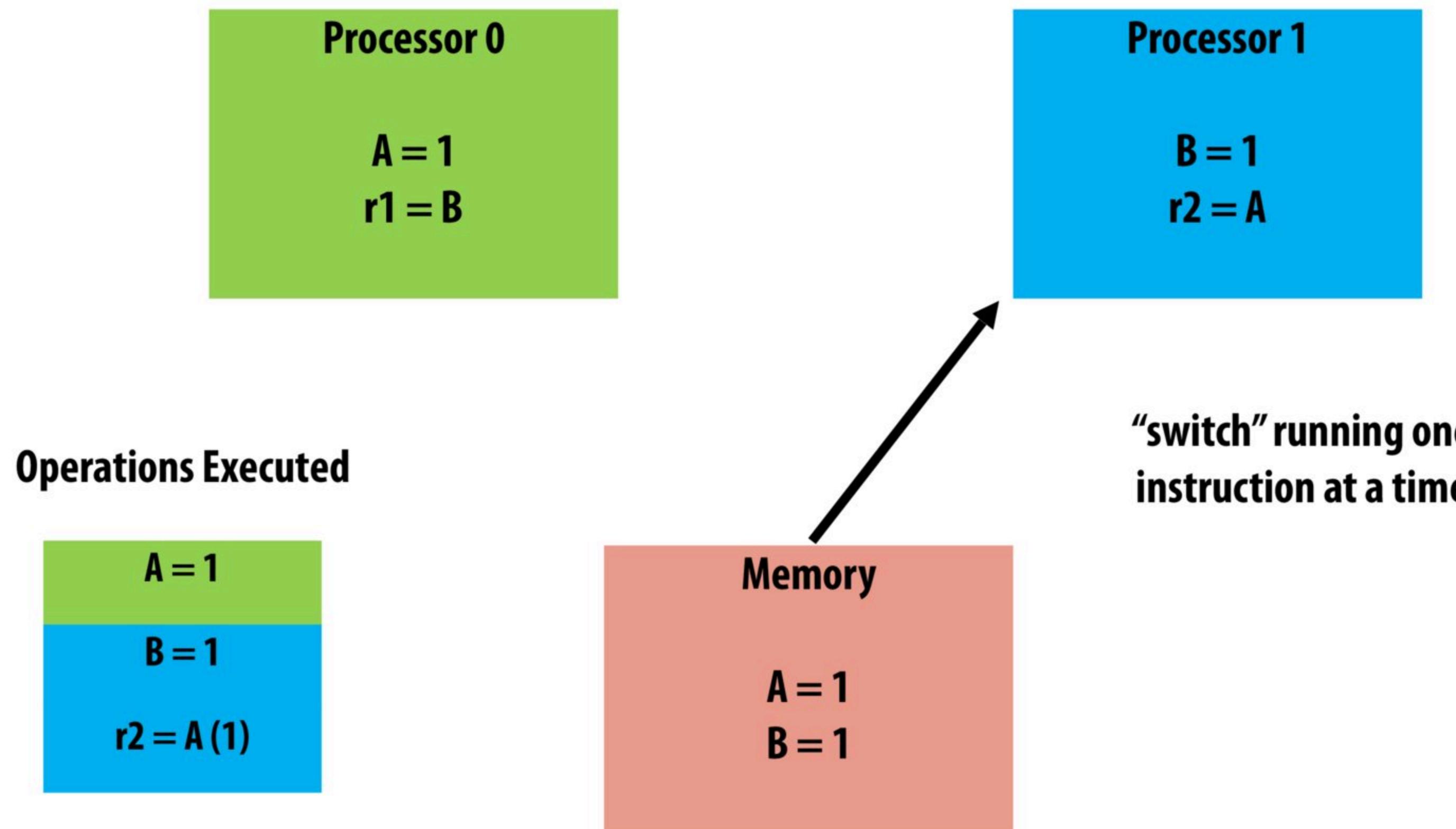
Sequential consistency example



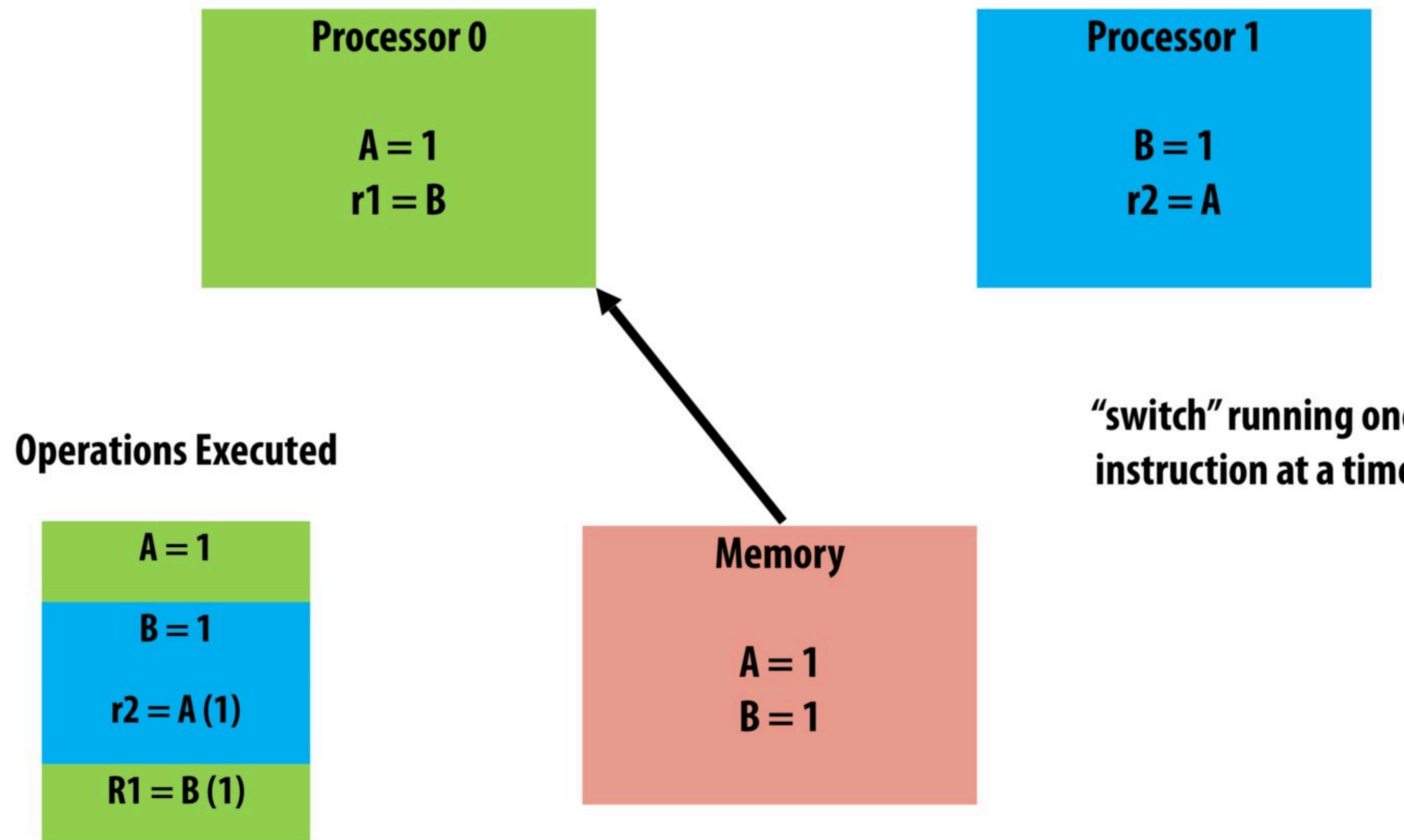
Sequential consistency example



Sequential consistency example



Sequential consistency example



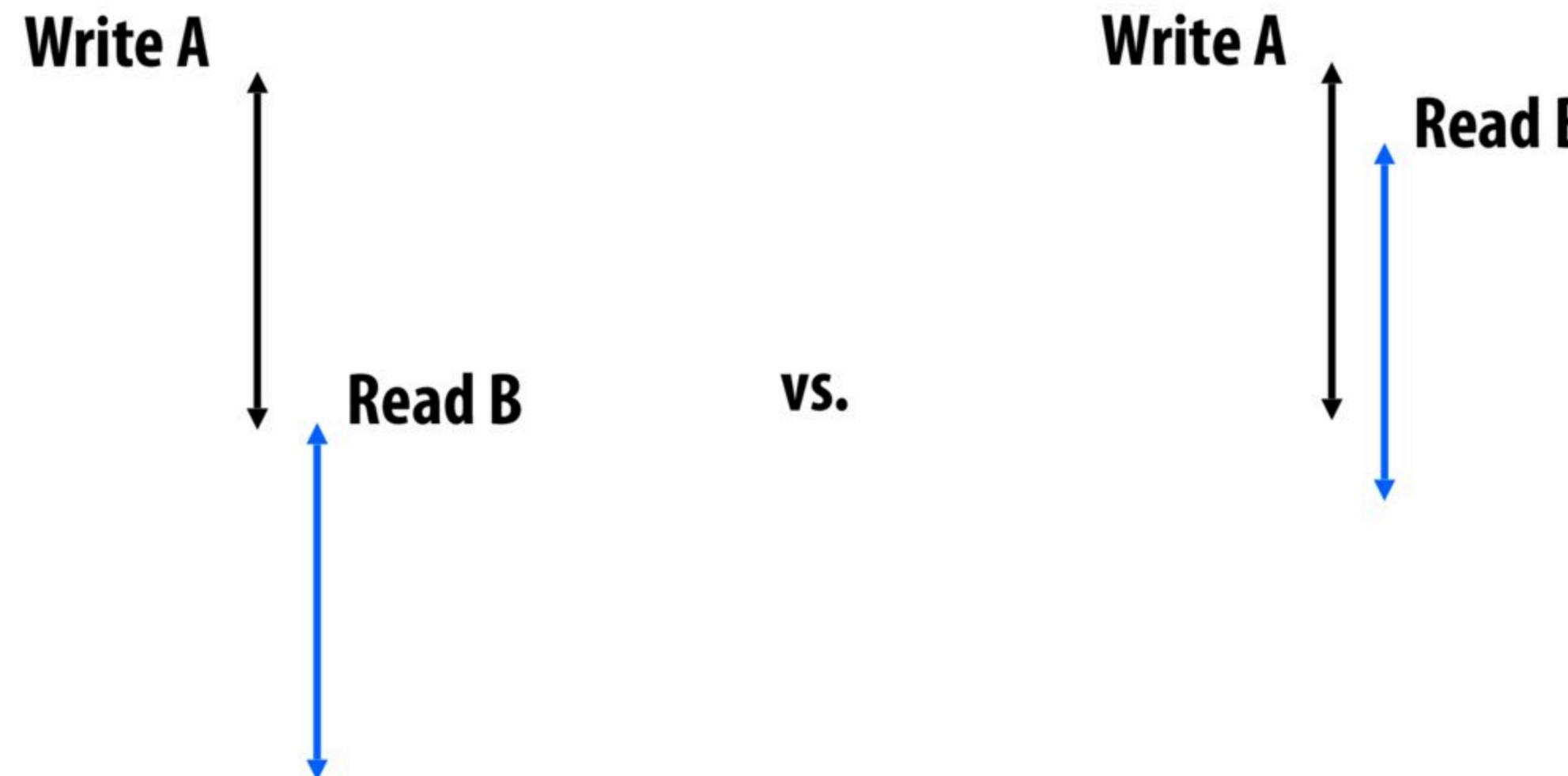
Relaxing memory operation ordering

- A sequentially consistent memory system maintains all four memory operation orderings ($W_x \rightarrow R_y, R_x \rightarrow R_y, R_x \rightarrow W_y, W_x \rightarrow W_y$)
- Relaxed memory consistency models allow certain orderings to be violated

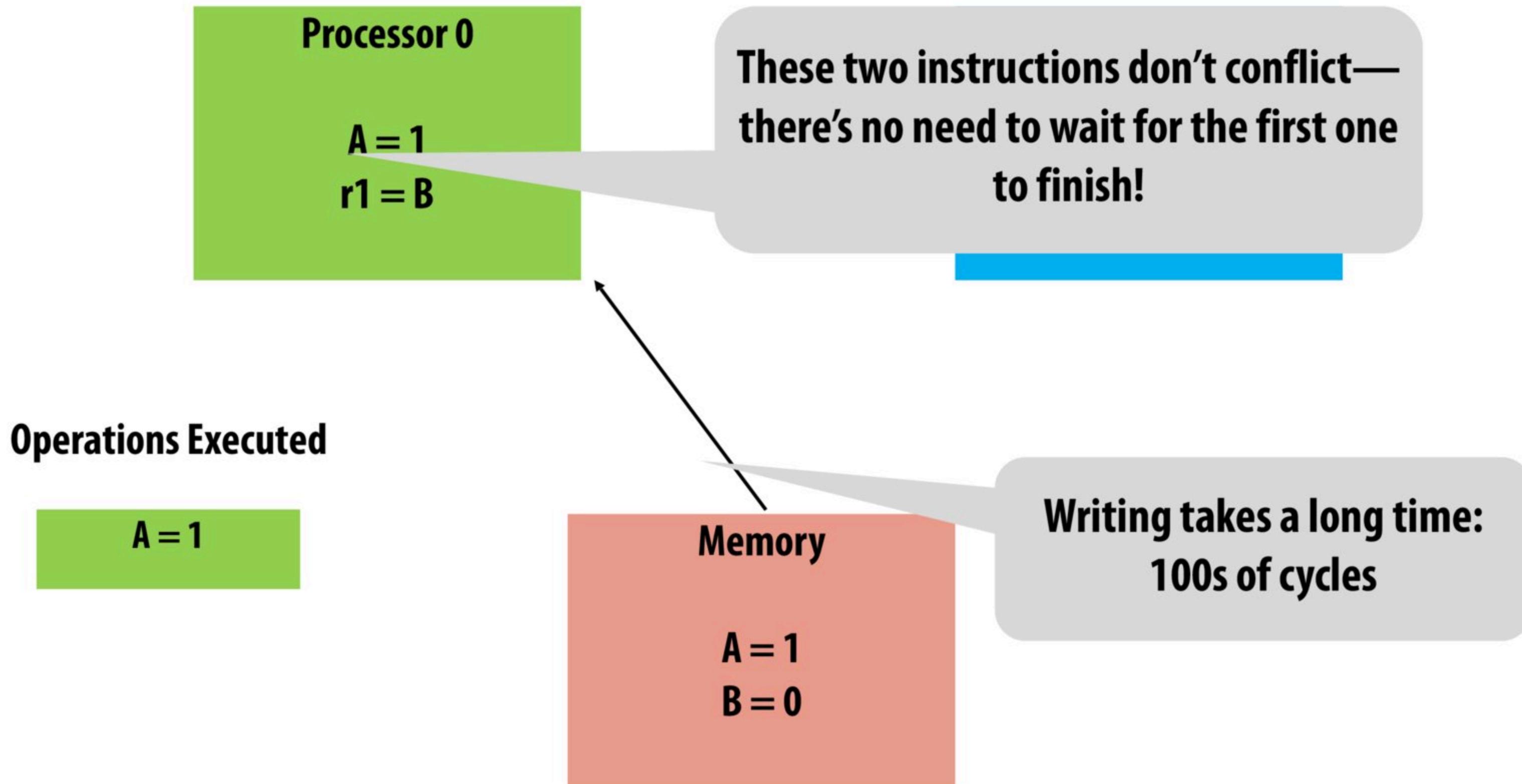
Motivation for relaxed consistency: hiding latency

■ Why are we interested in relaxing ordering requirements?

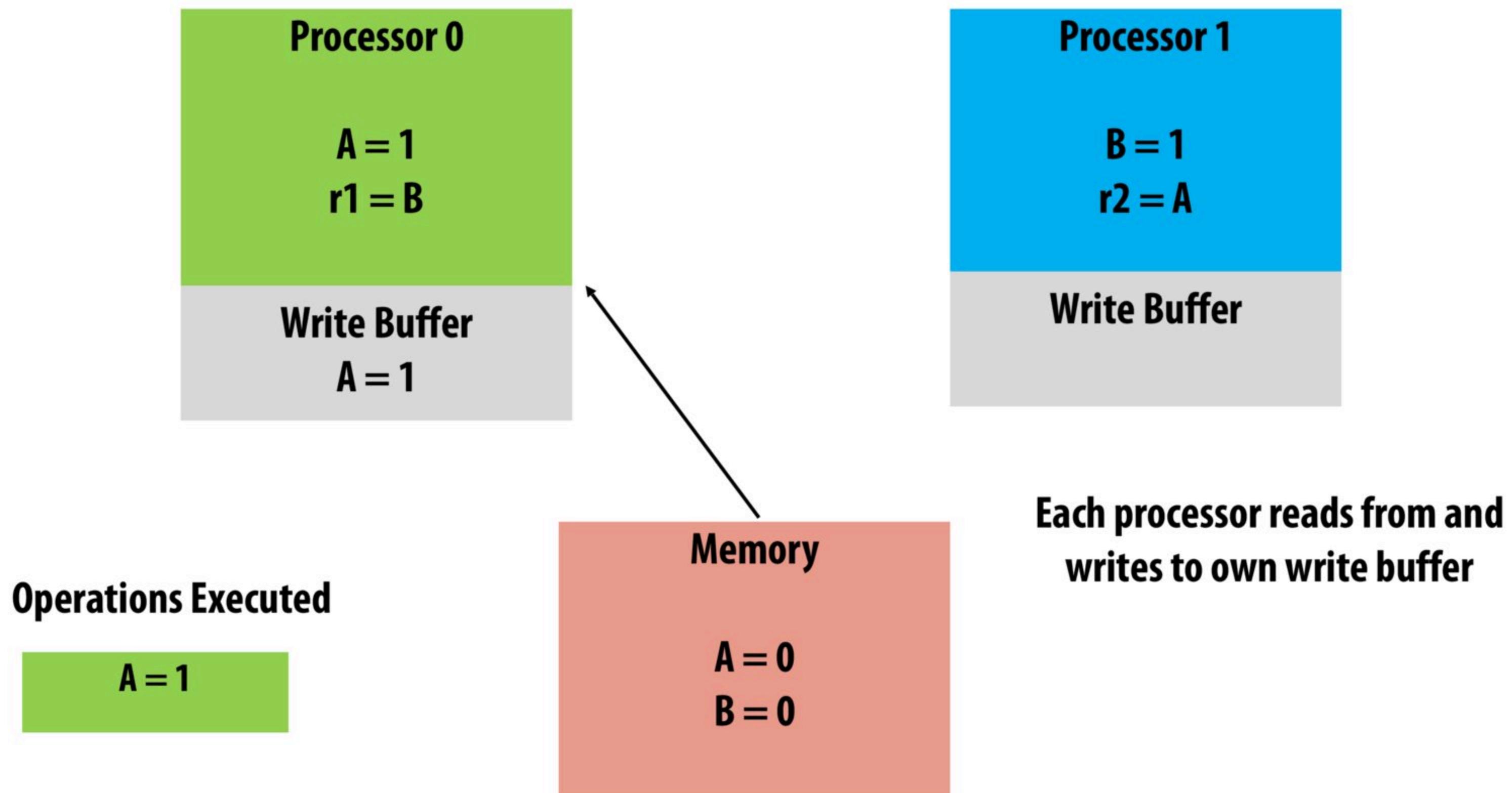
- To gain performance
- Specifically, hiding memory latency: overlap memory access operations with other operations when they are independent
- Remember, memory access in a cache coherent system may entail much more work then simply reading bits from memory (finding data, sending invalidations, etc.)



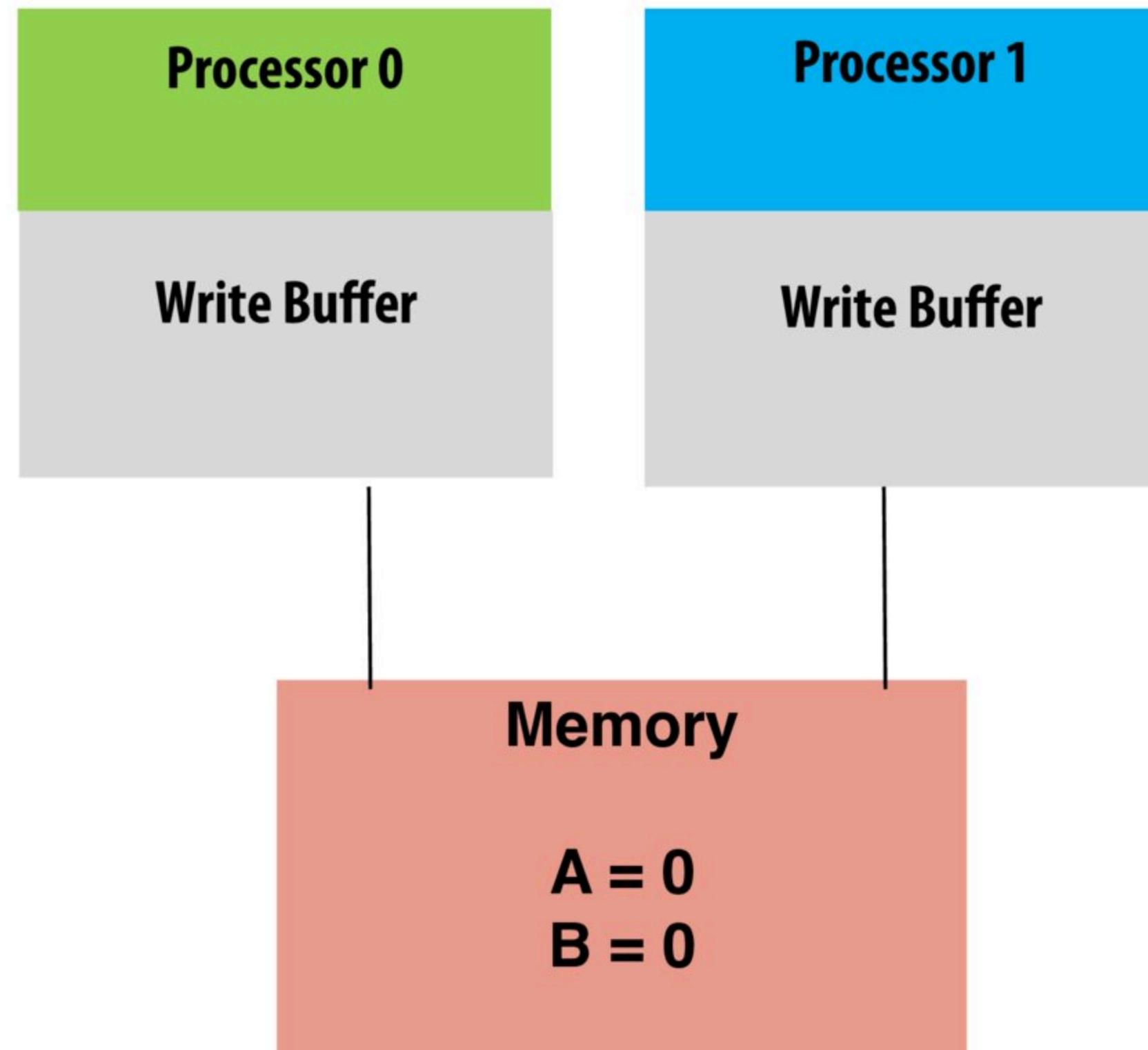
Problem with SC



Optimization: write buffer



Write buffers change memory behavior



Initially $A = B = 0$

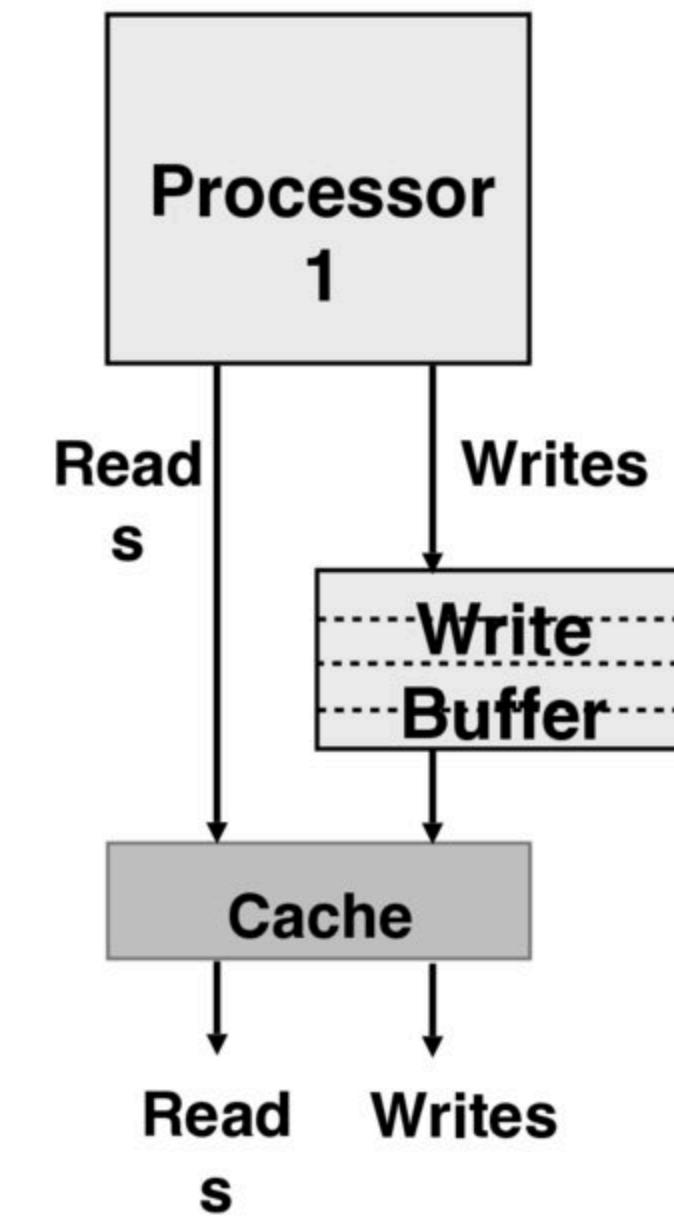
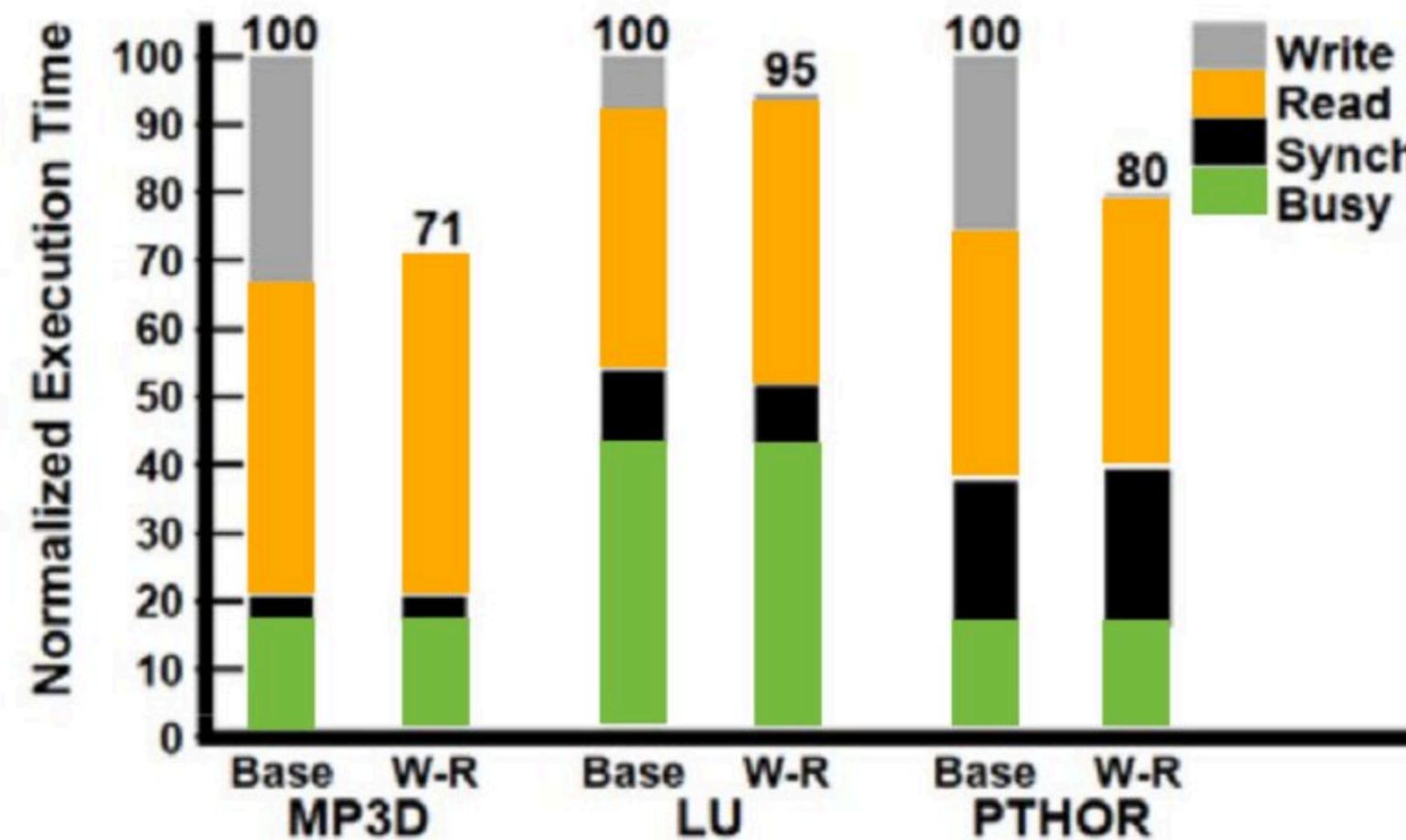
Proc 0	Proc 1
(1) $A = 1$	(3) $B = 1$
(2) $r_1 = B$	(4) $r_2 = A$

Can $r_1 = r_2 = 0$?

SC: No

Write buffers:

Write buffer performance



Base: Sequentially consistent execution. Processor issues one memory operation at a time, stalls until completion

W-R: relaxed W→R ordering constraint (write latency almost fully hidden)

Write buffers: who cares?

- **Performance improvement**
- **Every modern processor uses them**
 - **Intel x86, ARM, SPARC**
- **Need a weaker memory model**
 - **TSO: Total Store Order**
 - **Slightly harder to reason about than SC**
 - **x86 uses an incompletely specified form of TSO**

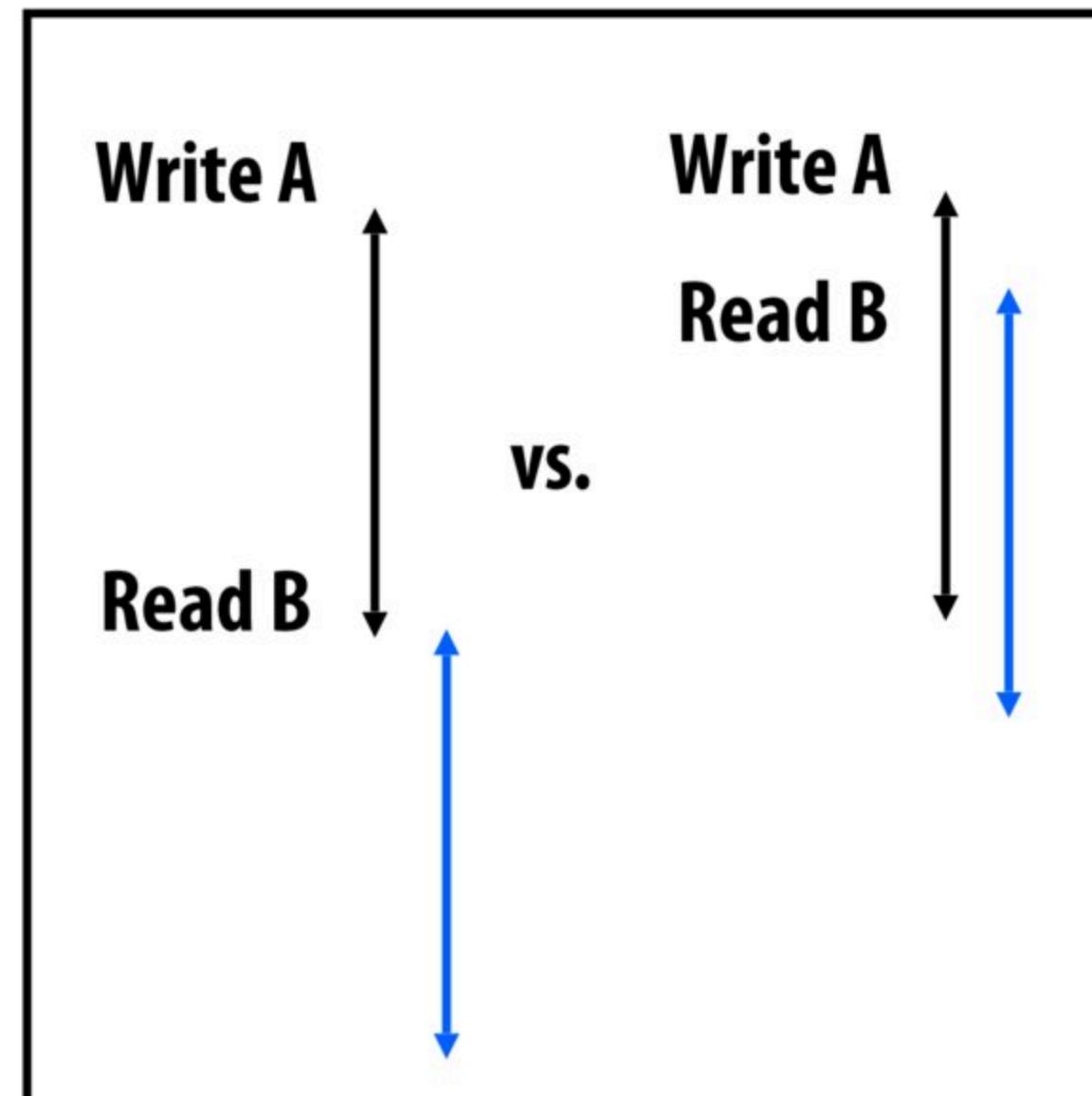
Allowing reads to move ahead of writes

■ Four types of memory operation orderings

- ~~$W_x \rightarrow R_y$: write must complete before subsequent read~~
- $R_x \rightarrow R_y$: read must complete before subsequent read
- $R_x \rightarrow W_y$: read must complete before subsequent write
- $W_x \rightarrow W_y$: write must complete before subsequent write

■ Allow processor to hide latency of writes

- Total Store Ordering (TSO)
- Processor Consistency (PC)



Allowing reads to move ahead of writes

- Total store ordering (TSO)
 - Processor P can read B before its write to A is seen by all processors
(processor can move its own reads in front of its own writes)
 - Reads by other processors cannot return new value of A until the write to A is observed by all processors
- Processor consistency (PC)
 - Any processor can read new value of A before the write is observed by all processors
- In TSO and PC, only $W_x \rightarrow R_y$ order is relaxed. The $W_x \rightarrow W_y$ constraint still exists. Writes by the same thread are not reordered (they occur in program order)

Clarification (make sure you get this!)

- The cache coherency problem exists because hardware implements the optimization of duplicating data in multiple processor caches. The copies of the data must be kept coherent.
- Relaxed memory consistency issues arise from the optimization of reordering memory operations. (Consistency is unrelated to whether or not caches exist in the system)

Allowing writes to be reordered

■ Four types of memory operation orderings

- ~~$W_x \rightarrow R_y$: write must complete before subsequent read~~
- $R_x \rightarrow R_y$: read must complete before subsequent read
- $R_x \rightarrow W_y$: read must complete before subsequent write
- ~~$W_x \rightarrow W_y$: write must complete before subsequent write~~

■ Partial Store Ordering (PSO)

- Execution may not match sequential consistency on the following program
(P2 may observe change to flag before change to A)

Thread 1 (on P1)

```
A = 1;  
flag = 1;
```

Thread 2 (on P2)

```
while (flag == 0);  
print A;
```

Why might it be useful to allow more aggressive memory operation reorderings?

- $W_x \rightarrow W_y$: processor might reorder write operations in a write buffer (e.g., one is a cache miss while the other is a hit)
- $R_x \rightarrow W_y, R_x \rightarrow R_y$: processor might reorder independent instructions in an instruction stream (out-of-order execution)
- Keep in mind these are all valid optimizations if a program consists of a single instruction stream

Allowing all reorderings

- Four types of memory operation orderings
 - ~~$W_x \Rightarrow R_y$: write must complete before subsequent read~~
 - ~~$R_x \Rightarrow R_y$: read must complete before subsequent read~~
 - ~~$R_x \Rightarrow W_y$: read must complete before subsequent write~~
 - ~~$W_x \Rightarrow W_y$: write must complete before subsequent write~~
- No guarantees about operations on data!
 - Everything can be reordered
- Motivation is increased performance
 - Overlap multiple reads and writes in the memory system
 - Execute reads as early as possible and writes as late as possible to hide memory latency
- Examples:
 - Weak ordering (WO)
 - Release Consistency (RC)

Synchronization to the Rescue

- Memory reordering seems like a nightmare (it is!)
- Every architecture provides synchronization primitives to make memory ordering stricter
- Fence (memory barrier) instructions prevent reorderings, but are expensive
 - All memory operations complete before any memory operation after it can begin
- Other synchronization primitives (per address):
 - read-modify-write/compare-and-swap, transactional memory, ...

reorderable reads
and writes here

...

MEMORY FENCE

...

reorderable reads
and writes here

...

MEMORY FENCE

Example: expressing synchronization in relaxed models

- Intel x86/x64 ~ total store ordering
 - Provides sync instructions if software requires a specific instruction ordering not guaranteed by the consistency model
 - mm_lfence (“load fence”: wait for all loads to complete)
 - mm_sfence (“store fence”: wait for all stores to complete)
 - mm_mfence (“mem fence”: wait for all memory operations to complete)

■ ARM processors: very relaxed consistency model

A cool post on the role of memory fences in x86:

<http://bartoszmilewski.com/2008/11/05/who-ordered-memory-fences-on-an-x86/>

ARM has some great examples in their programmer's reference:

http://infocenter.arm.com/help/topic/com.arm.doc.genc007826/Barrier_Litmus_Tests_and_Cookbook_A08.pdf

A great list of academic papers:

<http://www.cl.cam.ac.uk/~pes20/weakmemory/>

Problem: data races

- Every example so far has involved a data race
 - Two accesses to the same memory location
 - At least one is a write
 - Unordered by synchronization operations

Conflicting data accesses

- Two memory accesses by different processors conflict if...

- They access the same memory location
 - At least one is a write

- Unsynchronized program

- Conflicting accesses not ordered by synchronization (e.g., a fence, operation with release/acquire semantics, barrier, etc.)
 - Unsynchronized programs contain data races: the output of the program depends on relative speed of processors (non-deterministic program results)

Synchronized Programs

- Synchronized programs yield SC results on non-SC systems
 - Synchronized programs are data-race-free
- If there are no data races, reordering behavior doesn't matter
 - Accesses are ordered by synchronization, and synchronization forces sequential consistency
- In practice, most programs you encounter will be synchronized (via locks, barriers, etc. implemented in synchronization libraries)
 - Rather than via ad-hoc reads/writes to shared variables like in the example programs

Summary: relaxed consistency

- Motivation: obtain higher performance by allowing reordering of memory operations (reordering is not allowed by sequential consistency)
- One cost is software complexity: programmer or compiler must correctly insert synchronization to ensure certain specific operation orderings when needed
 - But in practice complexities encapsulated in libraries that provide intuitive primitives like lock/unlock, barrier (or lower-level primitives like fence)
 - Optimize for the common case: most memory accesses are not conflicting, so don't design a system that pays the cost as if they are
- Relaxed consistency models differ in which memory ordering constraints they ignore

Languages need memory models too

Thread 1

```
X = 0  
for i=0 to 100:  
    X = 1  
    print X
```

compiler

Thread 1

```
X = 1  
for i=0 to 100:  
    print X
```

Languages need memory models too

**Single threaded case: optimization of moving write to X out of the loop
is not visible to programmer**

Thread 1

```
X = 0
for i=0 to 100:
    X = 1
    print X
```

1111111111...

Thread 1

```
X = 1
for i=0 to 100:
    print X
```

1111111111...

Languages need memory models too

**Multi-threaded case: optimization of moving write to X outside the loop
is visible to programmer**

Thread 1	Thread 2
X = 0	X = 0
for i=0 to 100:	
X = 1	
print X	

1111111111...

1111011111...

Thread 1	Thread 2
X = 1	X = 0
for i=0 to 100:	
print X	

1111111111...

1111000000...

Language must provide a contract to programmers about how their memory operations will be reordered by the compiler e.g. no reordering of shared memory operations

Language level memory models

- Modern (C11, C++11) and not-so-modern (Java 5) languages guarantee sequential consistency for data-race-free programs (“SC for DRF”)
 - Compilers will insert the necessary synchronization to cope with the hardware memory model
- No guarantees if your program contains data races!
 - The intuition is that most programmers would consider a racy program to be buggy
- Use a synchronization library!

Summary: memory consistency models

- **Define the allowed reorderings of memory operations by hardware and compilers**
- **A contract between hardware or compiler and application software**
- **Motivation is more performant/more efficient hardware**
- **Details of memory model can be hidden in synchronization library**
 - **Requires data race free (DRF) programs**

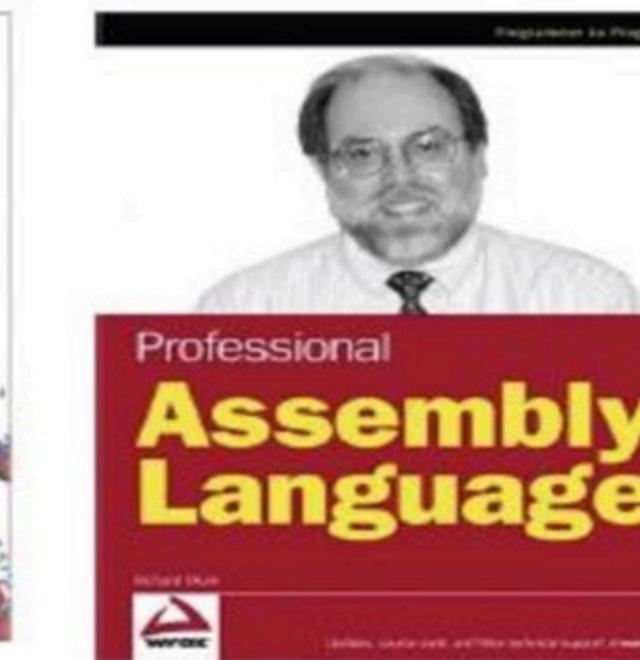
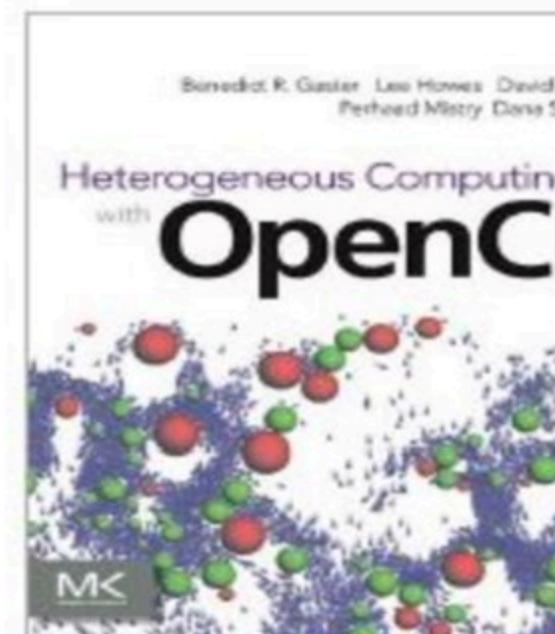
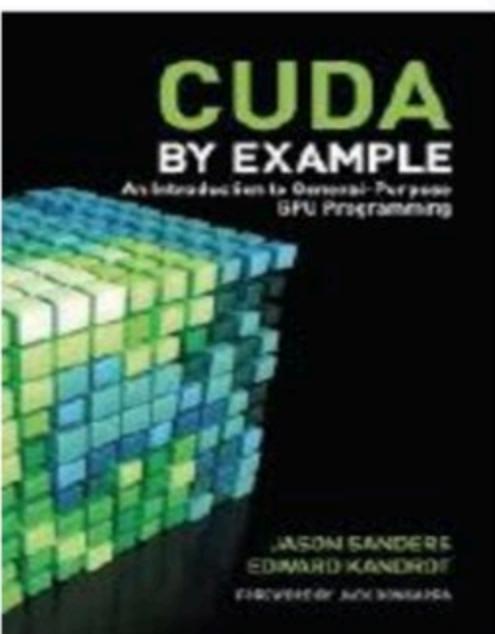
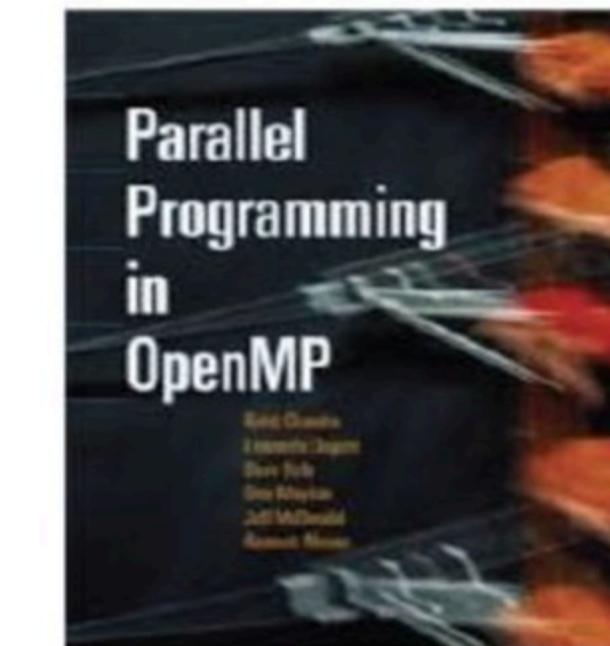
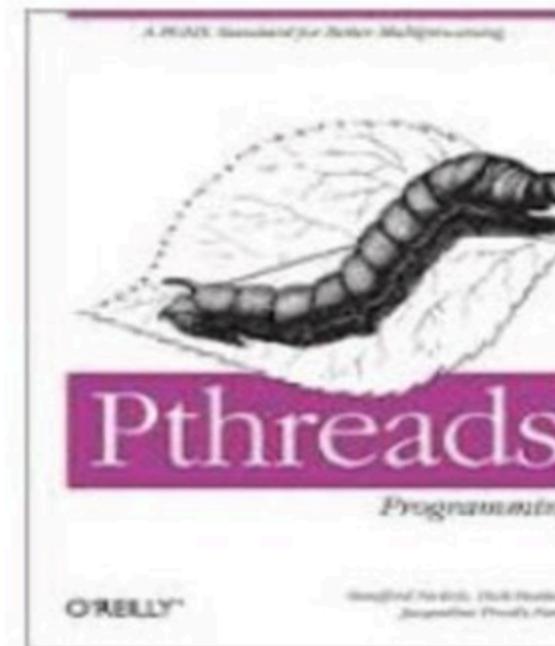
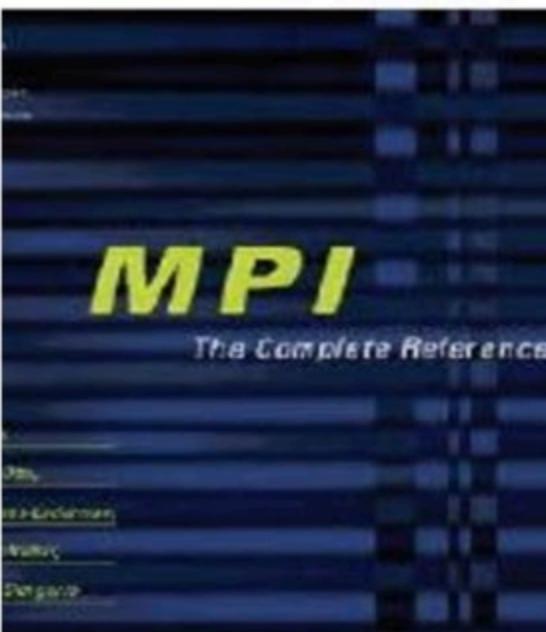
Relaxed memory consistency

Today

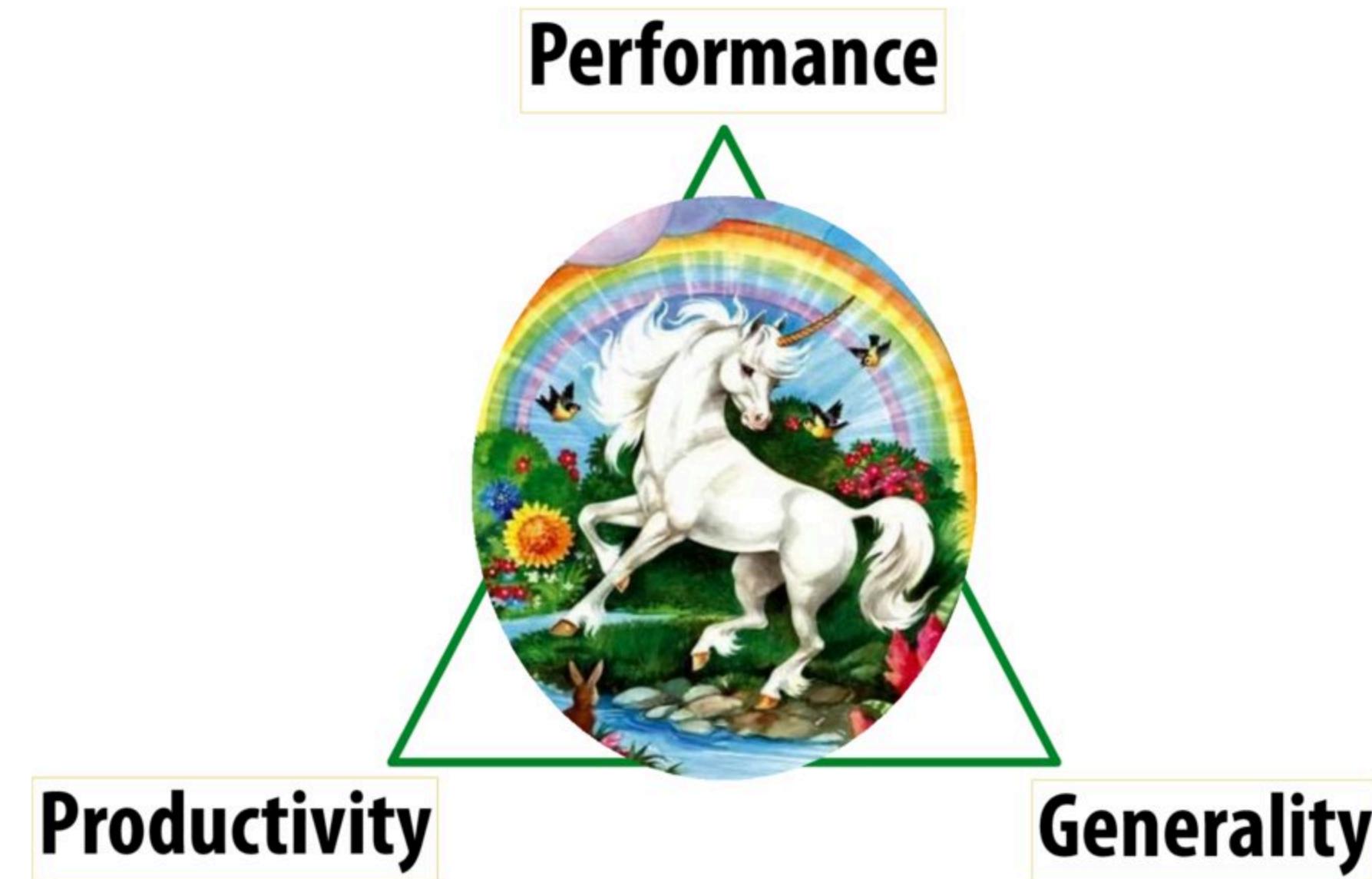
- Deeper dive into the idea of choosing the right abstractions for the job
- What is a domain specific programming language (DSL)?
- Two concrete examples in the slides:
 - Image processing in Halide (discussed in class, if time)
 - Physical simulation in Lizst (in extra slides for optional extra reading)
- Key concept: what are the advantages of performance-oriented application development using DSLs

CS149 educated programmers = hard to find

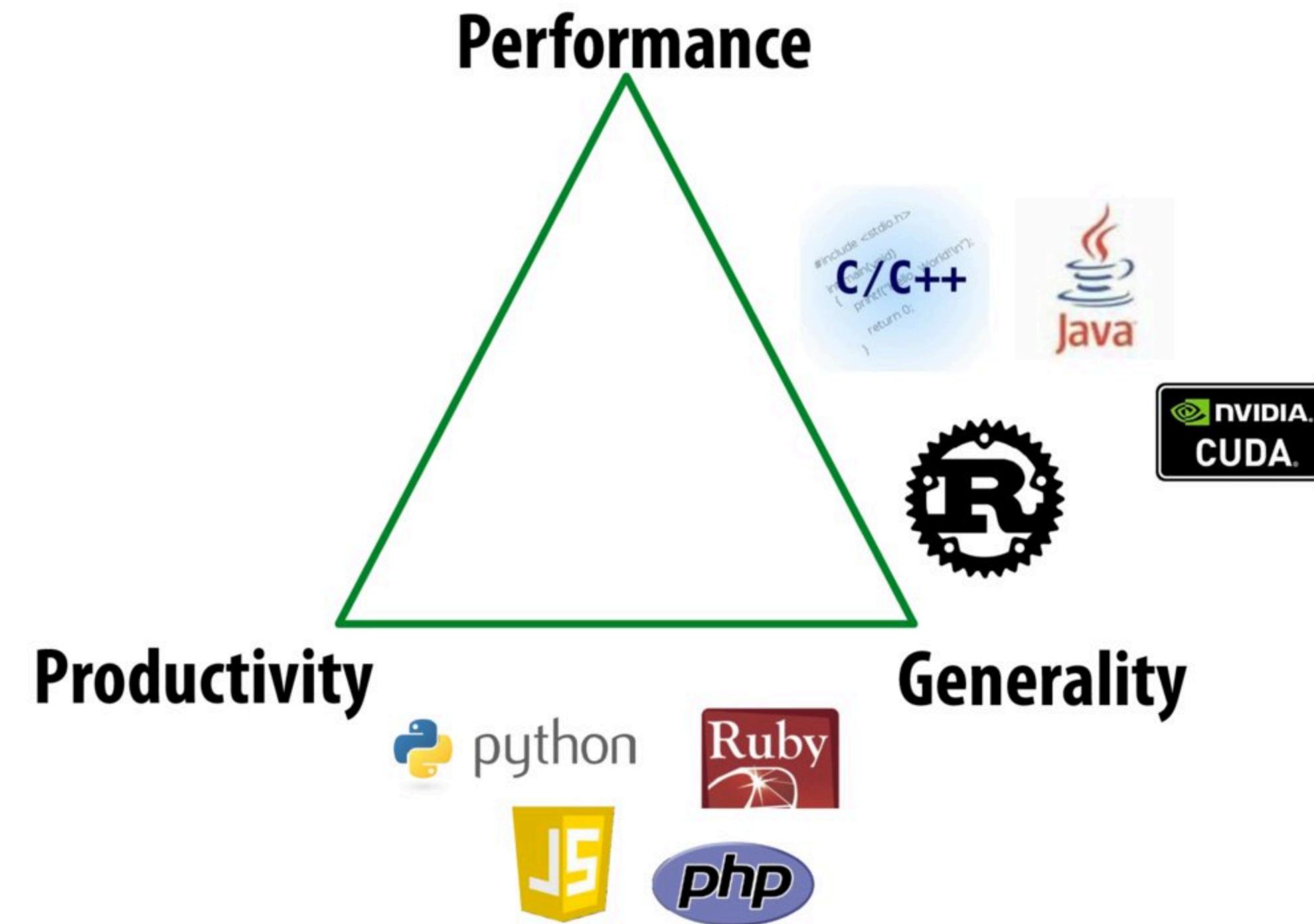
Performance optimization in languages like C++, ISPC, CUDA = low productivity (Proof by assignments 1, 2, 3, 4, etc...)



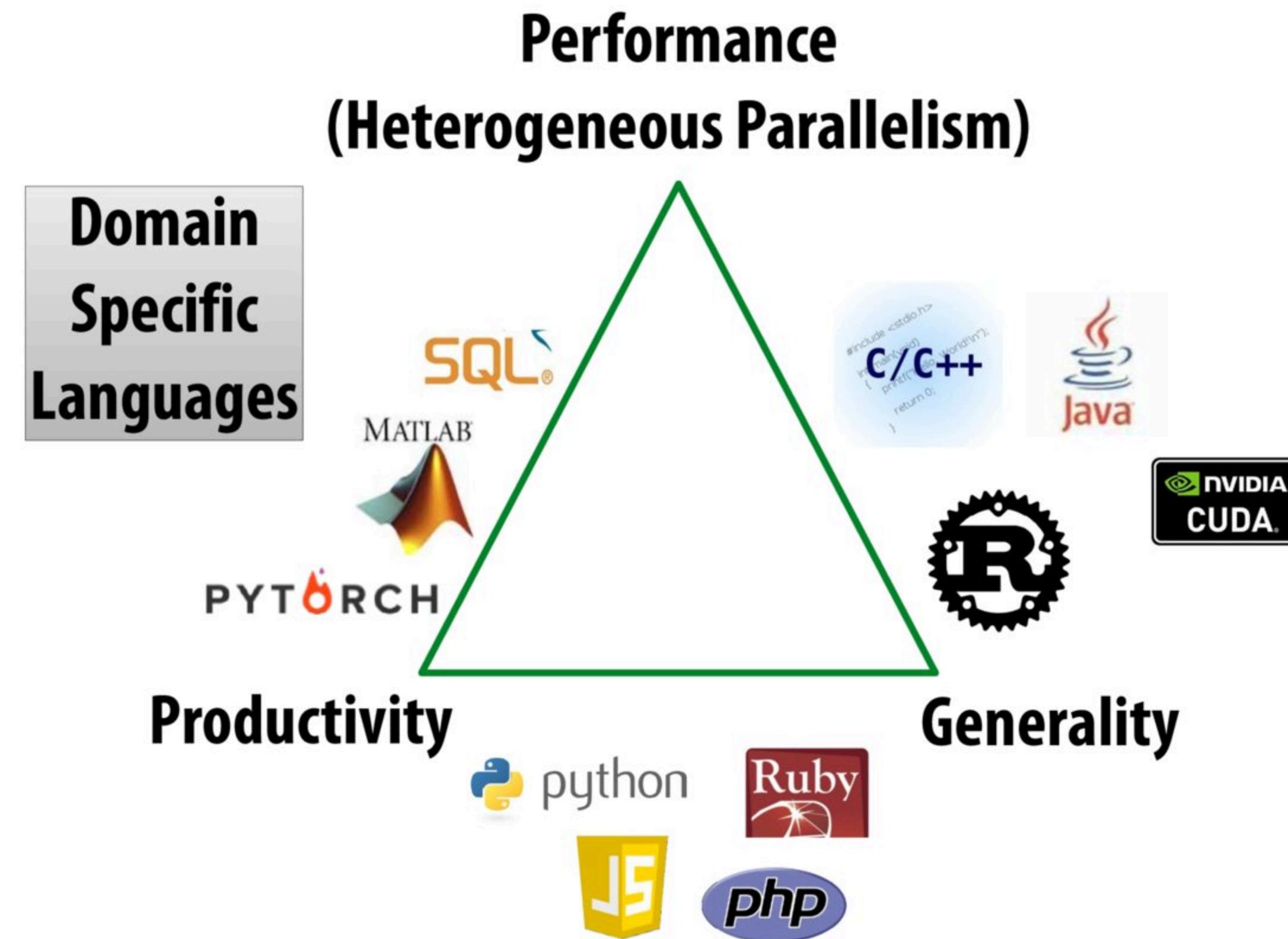
The ideal parallel programming language



Popular languages (not exhaustive ;-))



Way forward ⇒ domain-specific languages



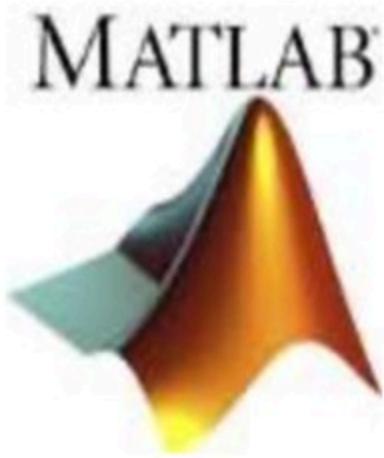
DSL hypothesis

**It is possible to write one program...
and**

run it efficiently on a range of heterogeneous parallel systems

Domain specific languages

- Domain Specific Languages (DSLs)
 - Programming language with restricted expressiveness for a particular domain
 - High-level, usually declarative, and deterministic



TensorFlow



TeX



Domain-specific programming systems

- Main idea: raise level of abstraction for expressing programs
 - Goal: write one program, and run it efficiently on different machines
- Introduce high-level programming primitives specific to an application domain
 - Productive: intuitive to use, portable across machines, primitives correspond to behaviors frequently used to solve problems in targeted domain
 - Performant: system uses domain knowledge to provide efficient, optimized implementation(s)
 - Given a machine: system knows what algorithms to use, parallelization strategies to employ for this domain
 - Optimization goes beyond efficient mapping of software to hardware! The hardware platform itself can be optimized to the abstractions as well
- Cost: loss of generality/completeness

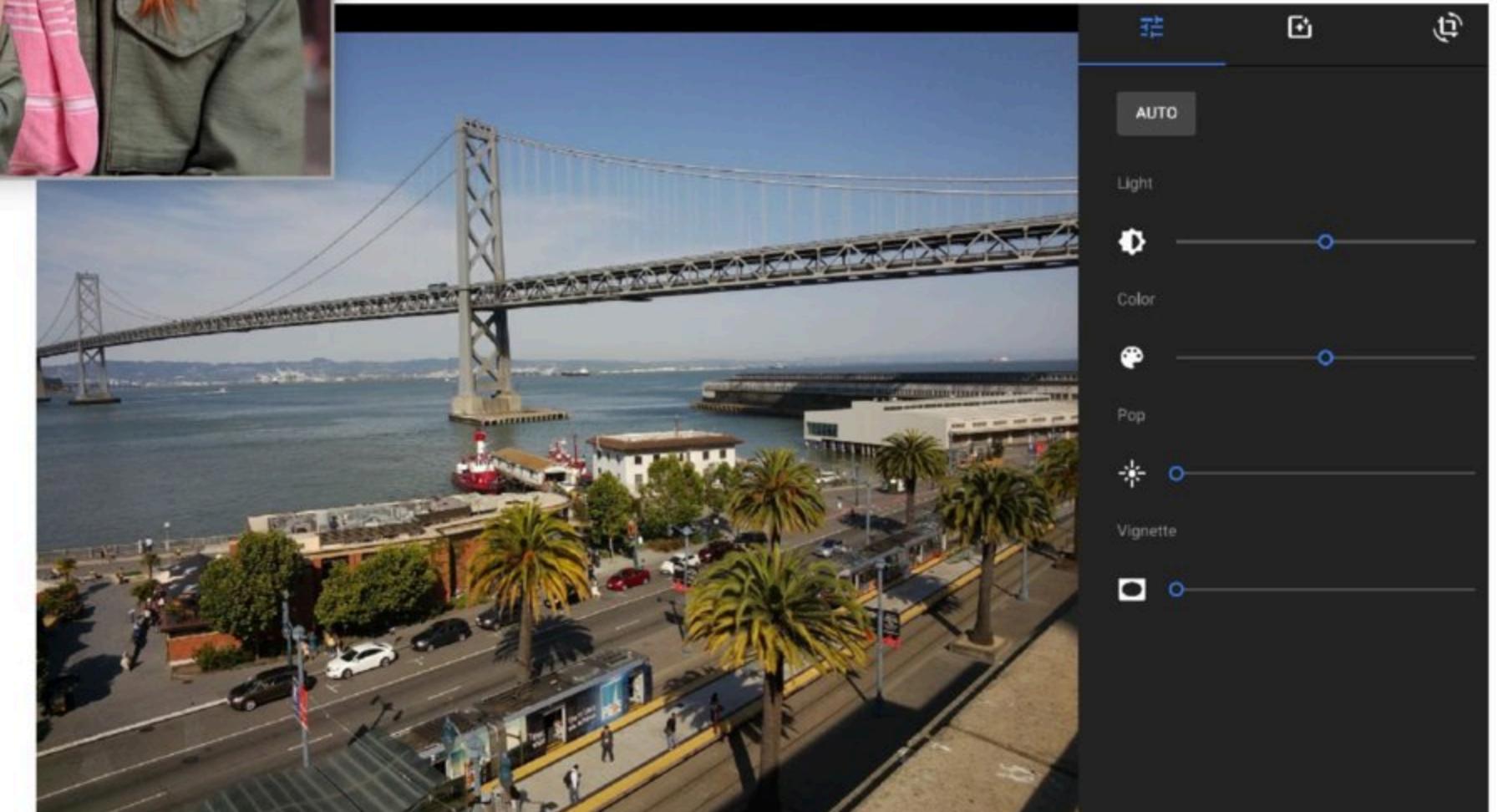
A DSL example:

Halide: a domain-specific language for image processing

Jonathan Ragan-Kelley, Andrew Adams et al.
[SIGGRAPH 2012, PLDI 13]

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.



A quick tutorial on high-performance image processing

What does this code do?



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);
                    }
                }
            }
        }
    }
}
```

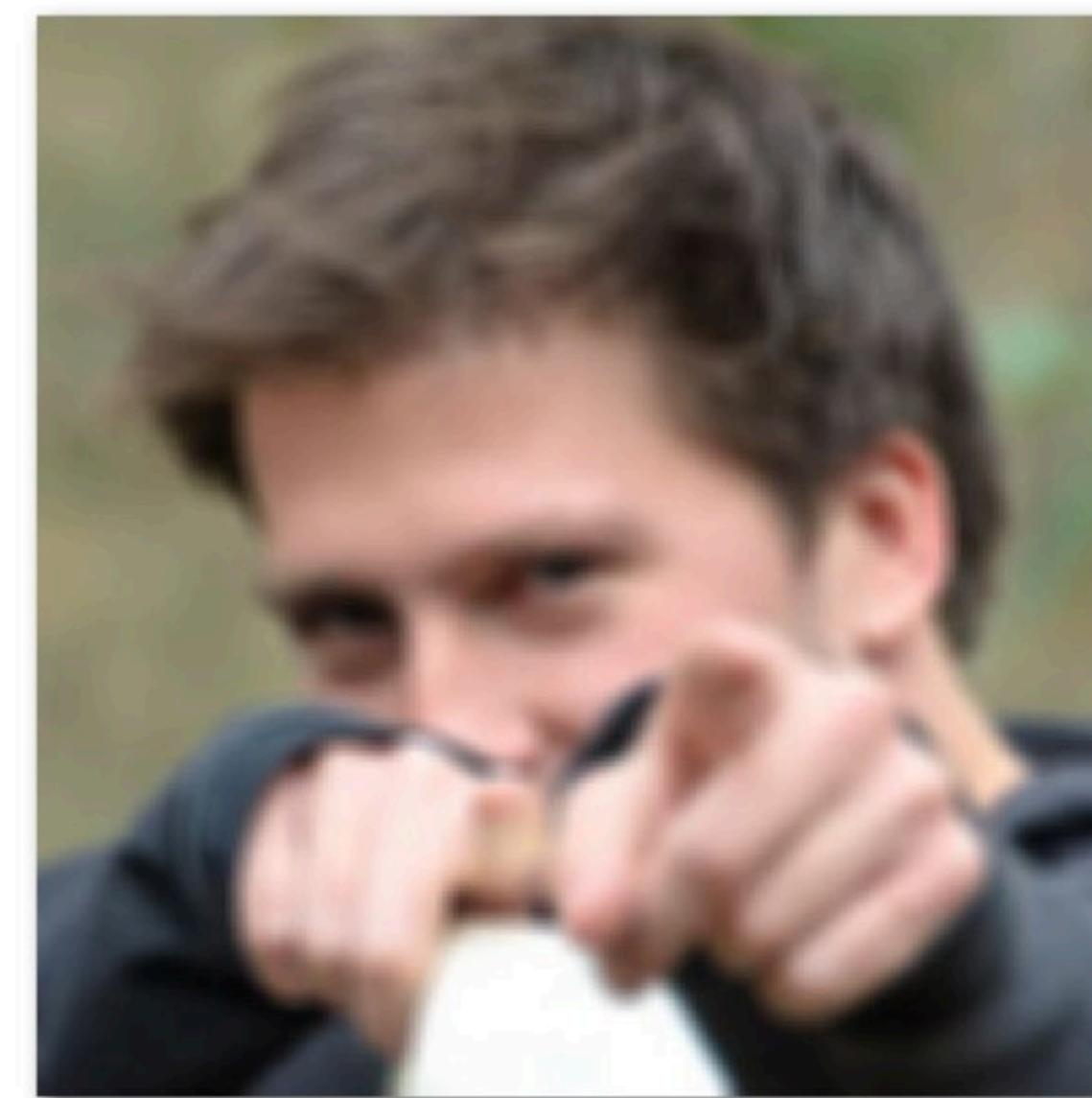
What does this C code do?

```
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;
    }
}
```

The code on the previous slide performed a 3x3 box blur



(Zoomed view)

3x3 image 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;
    }
}
```

Total work per image = $9 \times \text{WIDTH} \times \text{HEIGHT}$

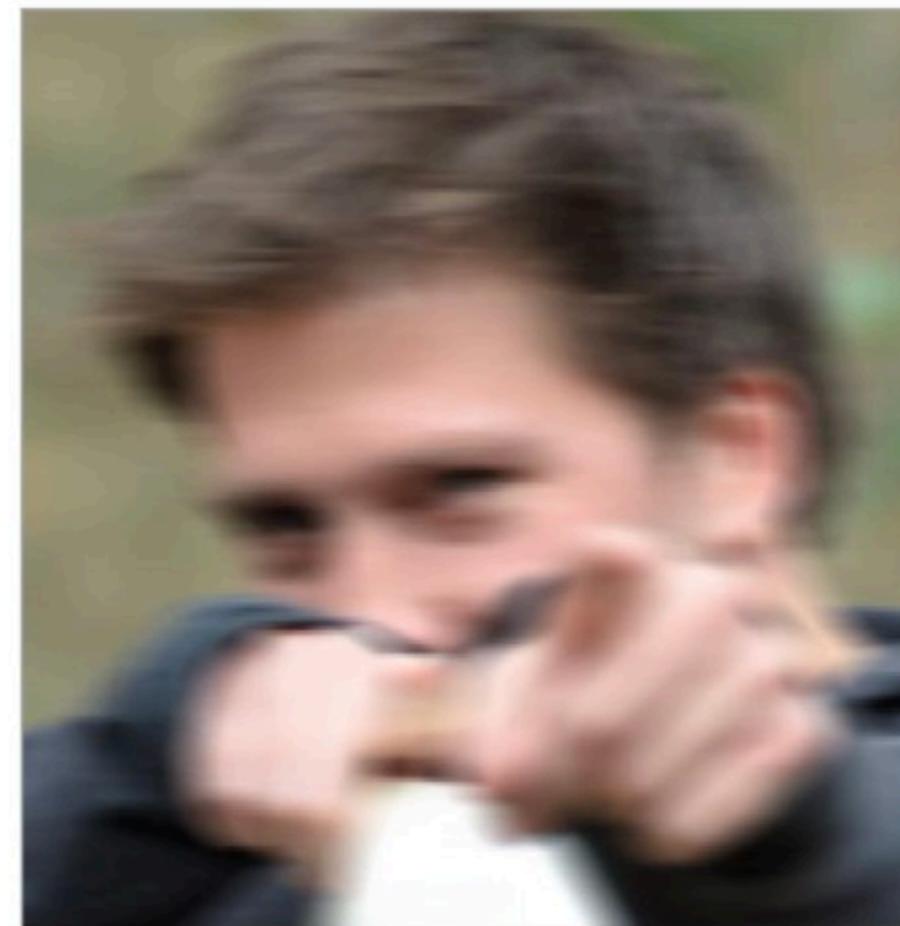
For NxN filter: $N^2 \times \text{WIDTH} \times \text{HEIGHT}$

Two-pass blur

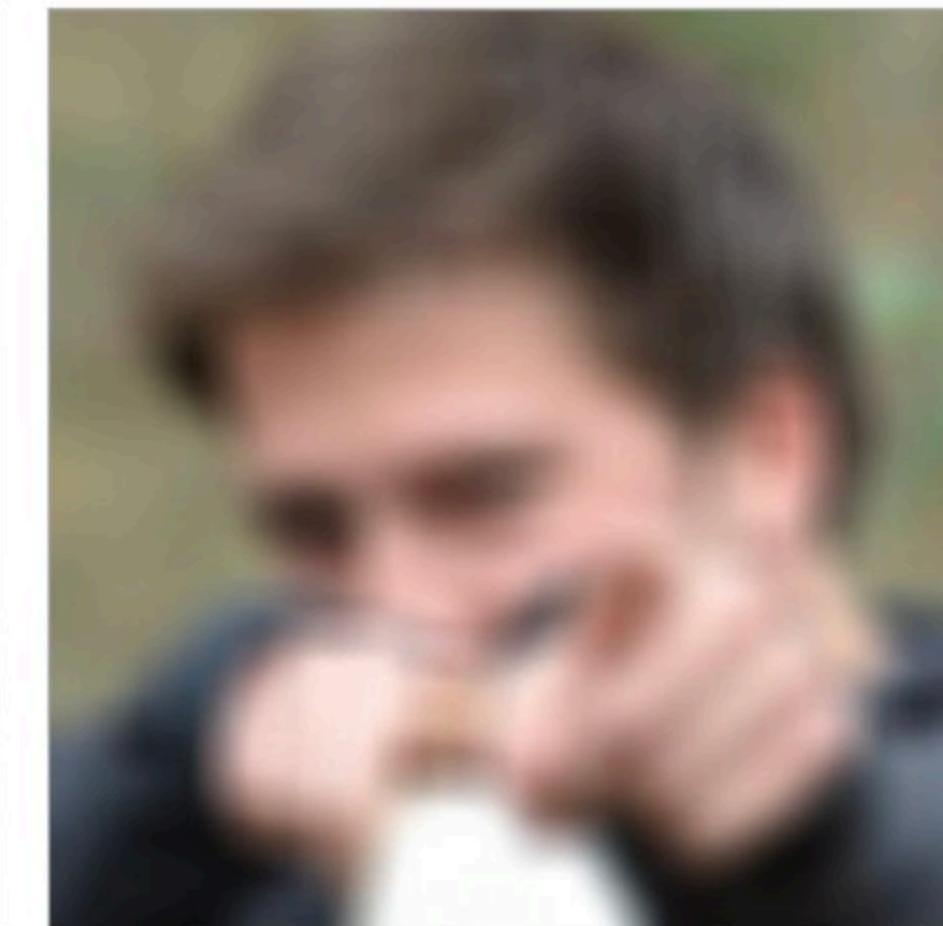
A 2D separable filter (such as a box filter) can be evaluated via two 1D filtering operations



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

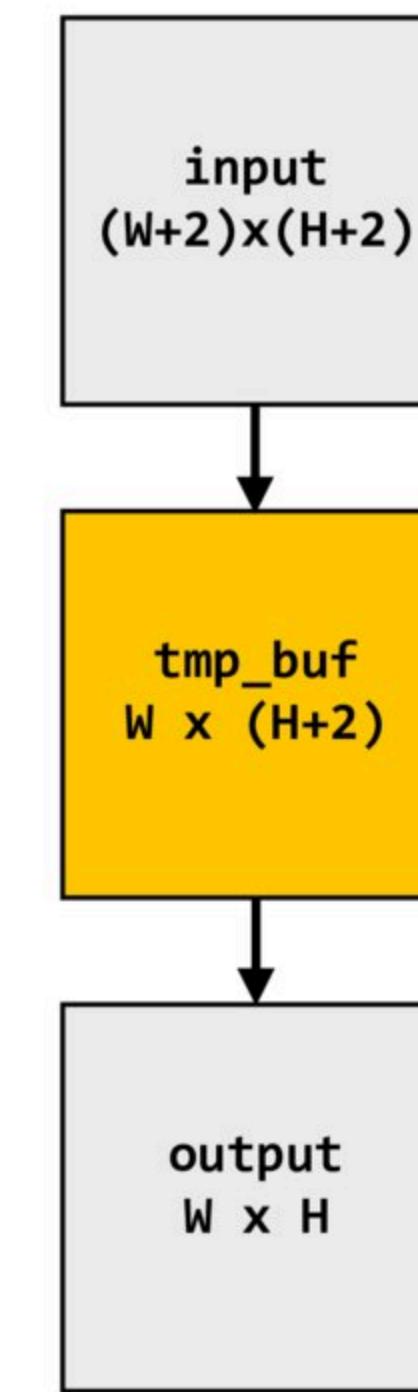
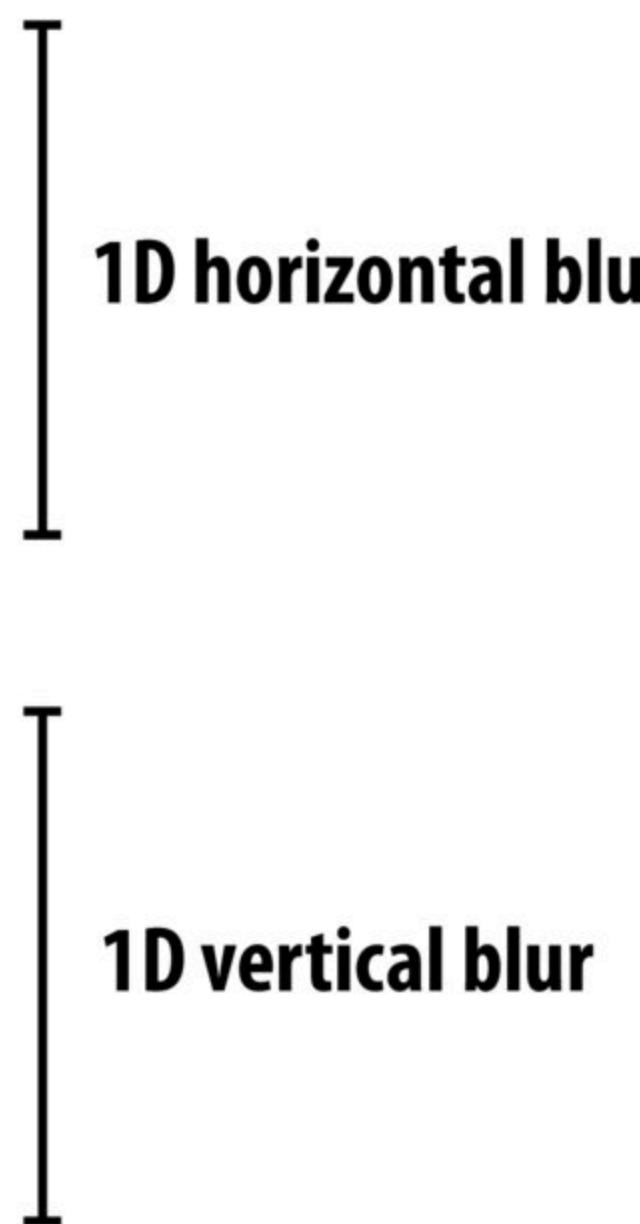
```
int WIDTH = 1024;  
int HEIGHT = 1024;  
float input[(WIDTH+2) * (HEIGHT+2)];  
float tmp_buf[WIDTH * (HEIGHT+2)];  
float output[WIDTH * HEIGHT];  
  
float weights[] = {1.f/3, 1.f/3, 1.f/3};  
  
for (int j=0; j<(HEIGHT+2); j++)  
    for (int i=0; i<WIDTH; i++) {  
        float tmp = 0.f;  
        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;  
    }  
}
```

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

For NxN 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: thinking about locality

```
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;
    }
}
```

Intrinsic bandwidth requirements of blur algorithm:
Application must read each element of input image
and must write each element of output image.

Data from `input` reused three times. (immediately reused in next two i-loop iterations after first load, never loaded again.)

- Perfect cache behavior: never load required data more than once
- Perfect use of cache lines (don't load unnecessary data into cache)

Two pass: loads/stores to `tmp_buf` are overhead (this memory traffic is an artifact of the two-pass implementation: it is not intrinsic to computation being performed)

Data from `tmp_buf` reused three times (but three rows of image data are accessed in between)

- Never load required data more than once... if cache has capacity for three rows of image
- Perfect use of cache lines (don't load unnecessary data into cache)

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)

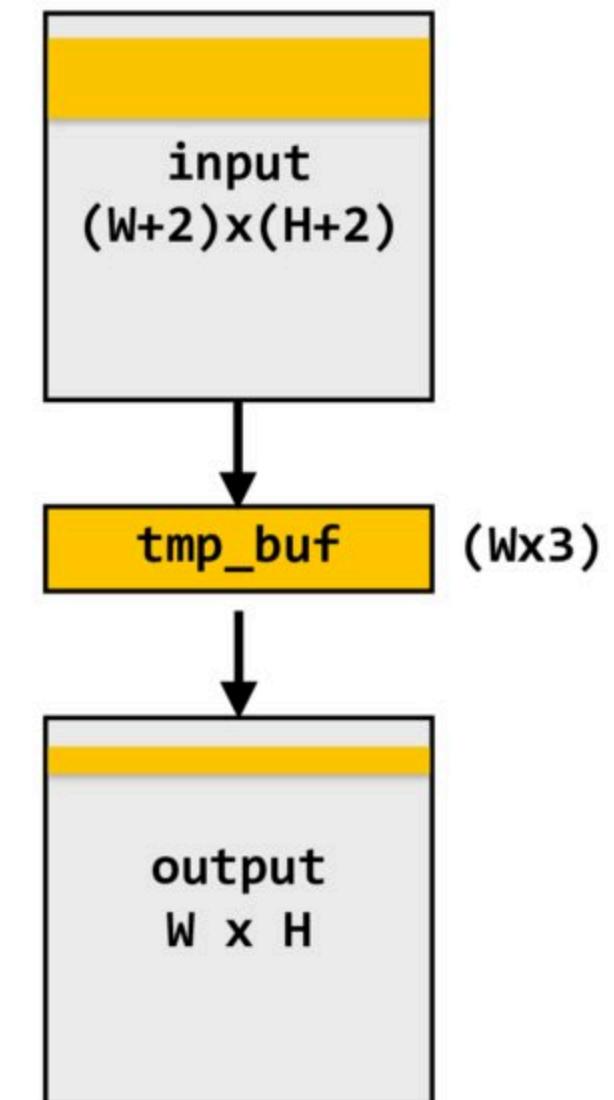
Combine them together to get one row of output

Total work per row of output:

- step 1: $3 \times 3 \times \text{WIDTH}$ work
- step 2: $3 \times \text{WIDTH}$ work

Total work per image = $12 \times \text{WIDTH} \times \text{HEIGHT}$????

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

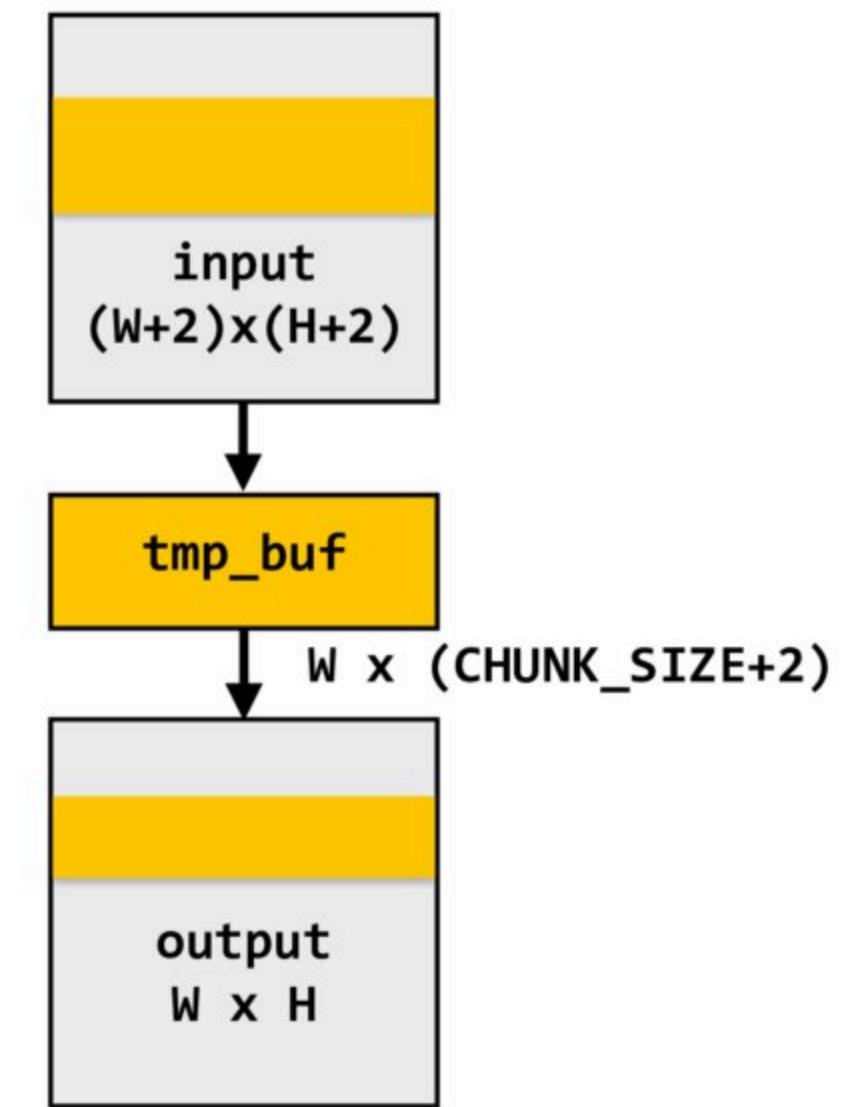
Total work per chuck of output: (assume **CHUNK_SIZE** = 16)

- Step 1: $18 \times 3 \times \text{WIDTH}$ work
- Step 2: $16 \times 3 \times \text{WIDTH}$ 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 $6 \times \text{WIDTH} \times \text{HEIGHT}$ as **CHUNK_SIZE** is increased!



Still not done

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

Optimized C++ code: 3x3 image blur



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

Halide language

[Ragan-Kelley / Adams 2012]

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

```
Var x, y;
Func blurx, blury, 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));
blury(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(blury(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 (0:1024,0:1024)
Halide::Buffer<uint8_t> result = out.realize(1024, 1024);
```

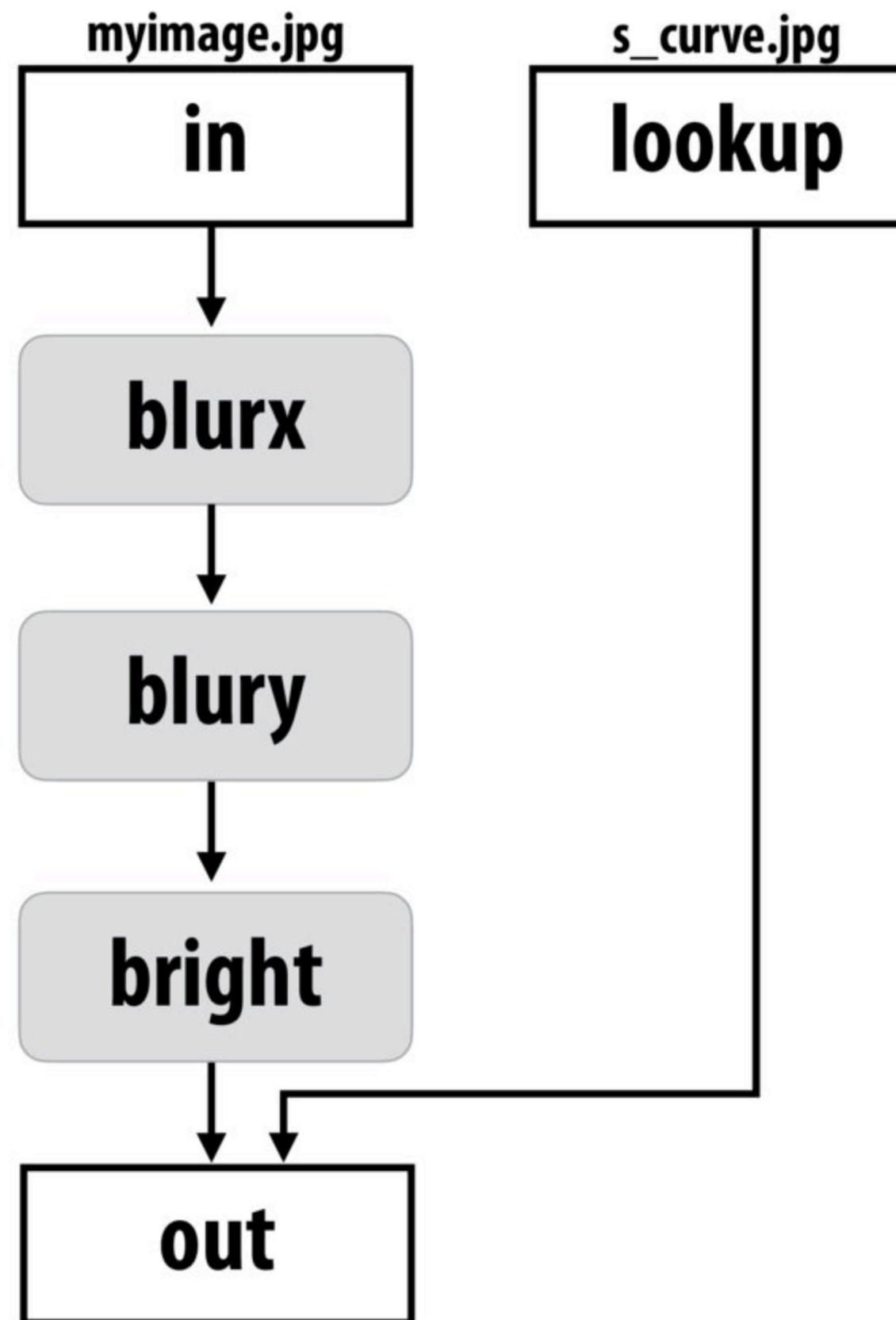
“Functions” map integer coordinates to values
(e.g., colors of corresponding pixels)

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

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.

Image processing application as a DAG of tensor operations



Key aspects of representation

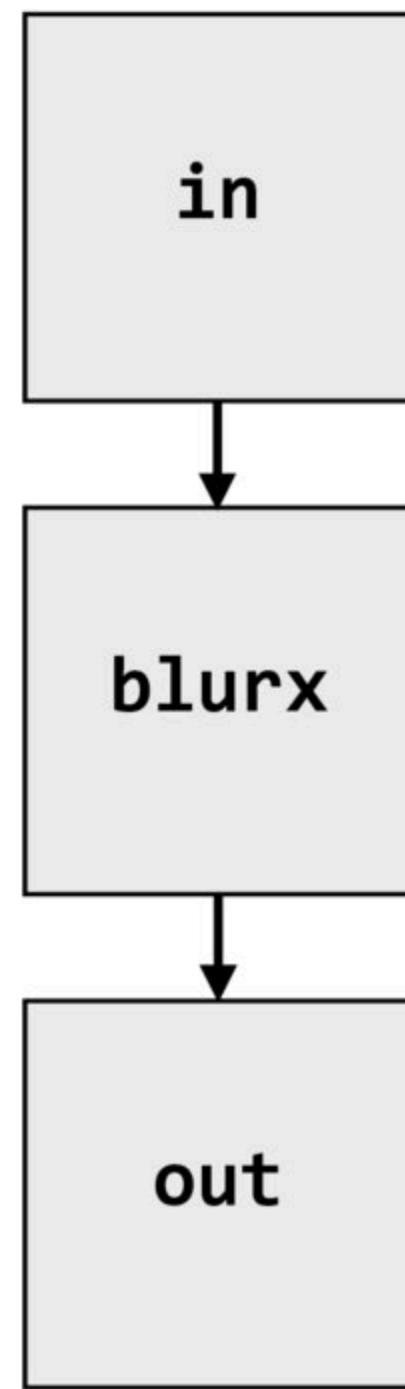
■ Intuitive expression:

- Adopts local “point wise” view of expressing algorithms
- Halide language is declarative. It does not define order of iteration, or what values in domain are stored!
 - It only defines what is 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 1024x1024
Halide::Buffer<uint8_t> result = out.realize(1024, 1024);
```



Real-world image processing pipelines feature complex sequences 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.

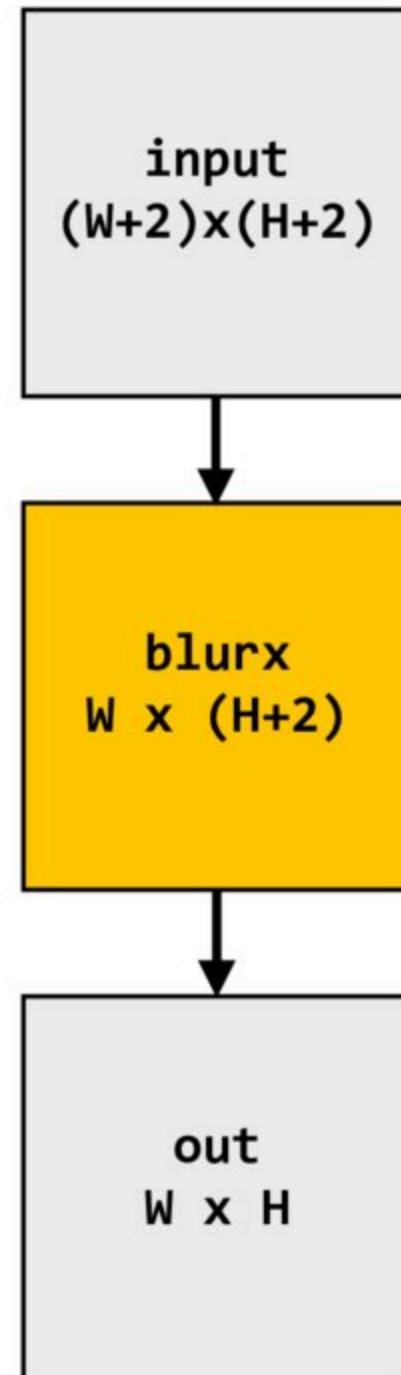
One (serial) implementation of Halide

```
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, conversation of quantities, 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.

(Aka... doing a CS149 assignment)

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

Specifying loop iteration order and parallelism

```
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;
```

Given this schedule for the function “out”...

```
out.tile(x, y, xi, yi, 256, 32).vectorize(xi,8).parallel(y);
```

Halide compiler will generate this parallel, vectorized loop nest for computing elements of out...

```
for y=0 to HEIGHT  
  for x=0 to WIDTH  
    blurx(x,y) = ...  
  
  for y=0 to num_tiles_y:          // parallelize this loop with threads  
    for x=0 to num_tiles_x:  
      for yi=0 to 32:  
        for xi=0 to 256 by 8:    // vectorize this loop with SIMD instr  
          idx_x = x*256+xi;  
          idx_y = y*32+yi  
          out(idx_x, idx_y) = ... (simd arithmetic here)
```

Primitives for how to interleave producer/consumer processing (perform fusion optimizations)

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

blurx.compute_root(); Do not compute blurx within out's loop nest.
 Compute all of blurx, then all of out

```
allocate buffer for all of blurx(x,y)  
for y=0 to HEIGHT:  
  for x=0 to WIDTH:  
    blurx(x,y) = ...
```

[all of blurx is computed here]

```
for y=0 to num_tiles_y:  
  for x=0 to num_tiles_x:  
    for yi=0 to 32:  
      for xi=0 to 256:  
        idx_x = x*256+xi;  
        idx_y = y*32+yi  
        out(idx_x, idx_y) = ...
```

[values of blurx consumed here]

Primitives for how to interleave producer/consumer processing

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

```
blurx.compute_at(out, xi);
```

Compute necessary elements of blurx
within out's xi loop nest

```
for y=0 to num_tiles_y:  
  for x=0 to num_tiles_x:  
    for yi=0 to 32:  
      for xi=0 to 256:  
        idx_x = x*256+xi;  
        idx_y = y*32+yi
```

Note: Halide compiler performs analysis that the
output of each iteration of the xi loop required 3
elements of blurx

allocate 3-element buffer for tmp_blurx

```
// compute 3 elements of blurx needed for out(idx_x, idx_y) here  
for (blur_x=0 to 3)  
  tmp_blurx(blur_x) = ...  
  
out(idx_x, idx_y) = ...
```

Primitives for how to interleave producer/consumer processing

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

```
blurx.compute_at(out, x);
```

Compute necessary elements of blurx within out's x loop nest (all necessary elements for one tile of out)

```
for y=0 to num_tiles_y:  
  for x=0 to num_tiles_x:  
  
    allocate 258x34 buffer for tile blurx  
    for yi=0 to 32+2:  
      for xi=0 to 256+2:  
        tmp_blurx(xi,yi) = // compute blurx from in  
        [ tile of blurx is computed here ]  
  
    for yi=0 to 32:  
      for xi=0 to 256:  
        idx_x = x*256+xi;  
        idx_y = y*32+yi  
        out(idx_x, idx_y) = ...  
        [ tile of blurx is consumed here ]
```

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 258x34 buffer for tile blurx
        for yi=0 to 32+2:
            for xi=0 to 256+2 BY 8:
                tmp_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 of how to schedule the application efficiently on machine (but it's slow and tedious), so Halide gives programmer a language to express 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 details of the schedule in terms of mechanisms available on the target machine (pthreads, AVX intrinsics, etc.)

Constraints on language

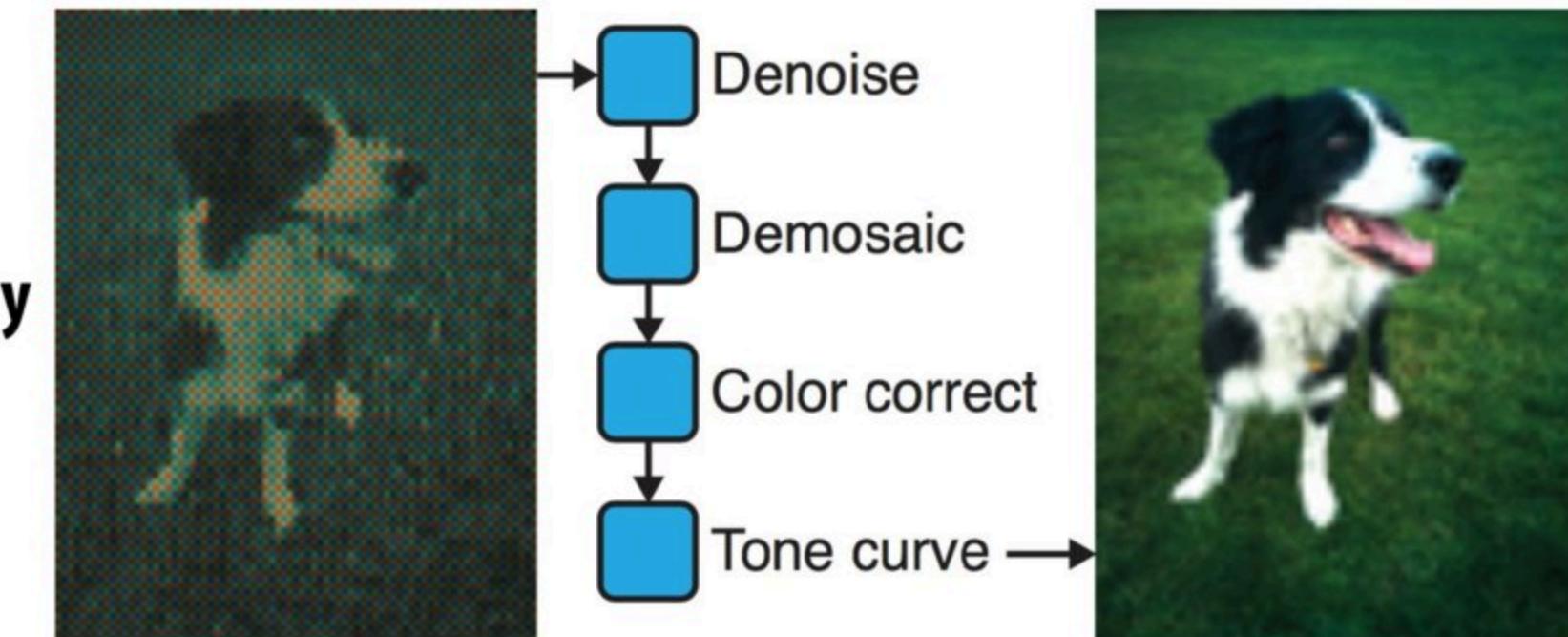
(to enable compiler to provide desired services)

- 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**

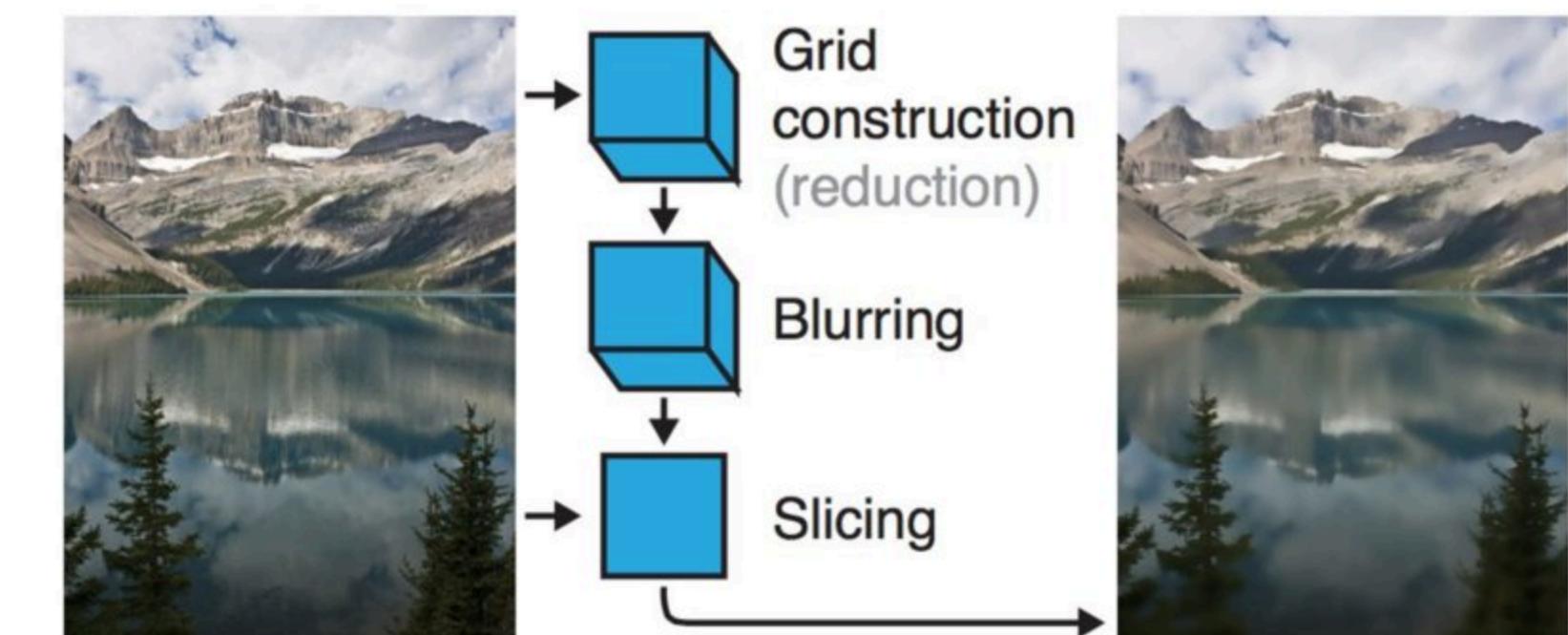
Initial academic Halide results

[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**
 - **Halide provides primitives for a programmer (that has strong knowledge of code optimization) to rapidly express what optimizations the system should apply**
 - **Halide compiler carries out the nitty-gritty of mapping that strategy to a machine**

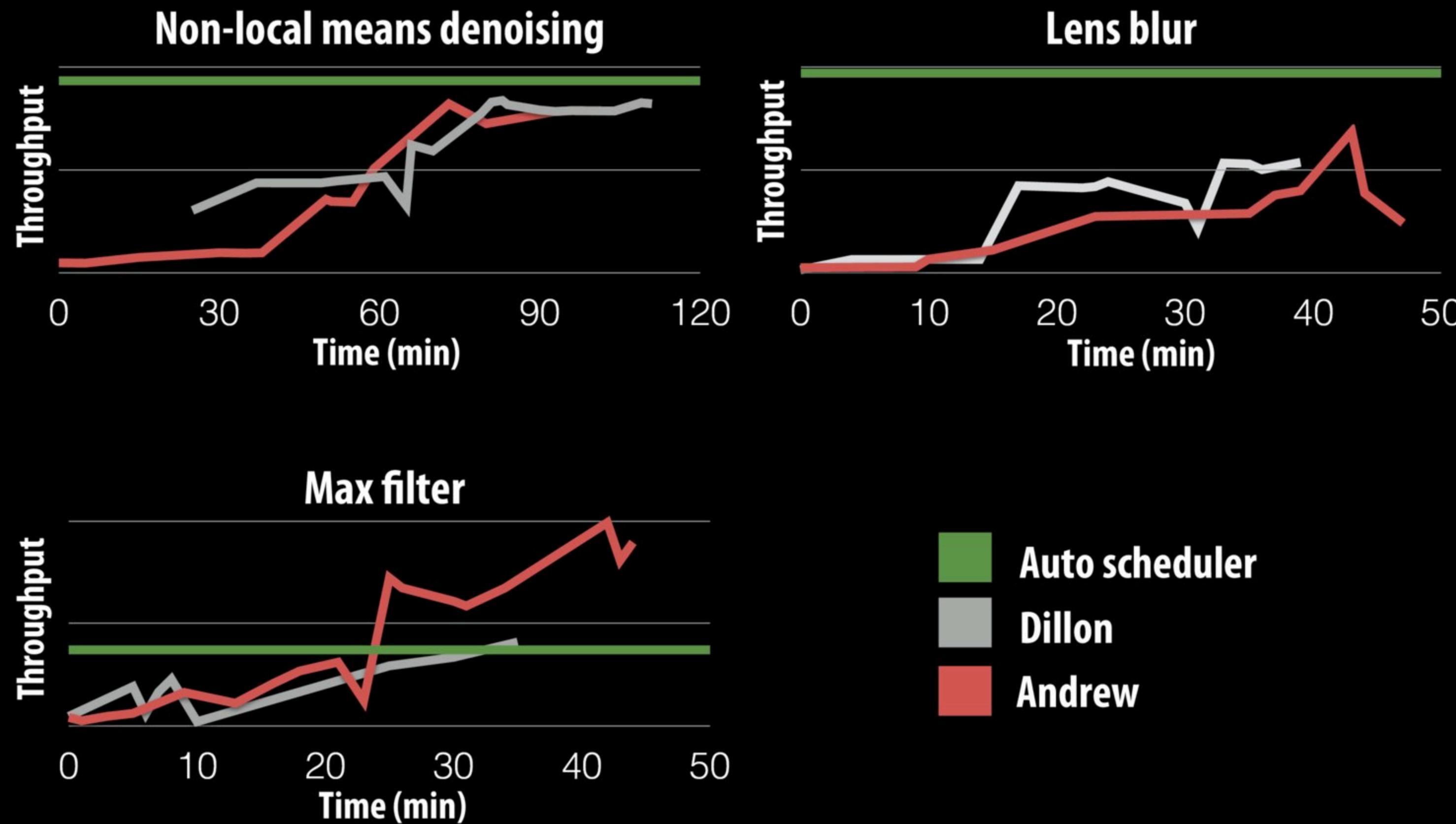
Automatically generating Halide schedules

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

See "*Learning to Optimize Halide with Tree Search and Random Programs*", Adams et al. SIGGRAPH 2019

Autoscheduler saves time for experts

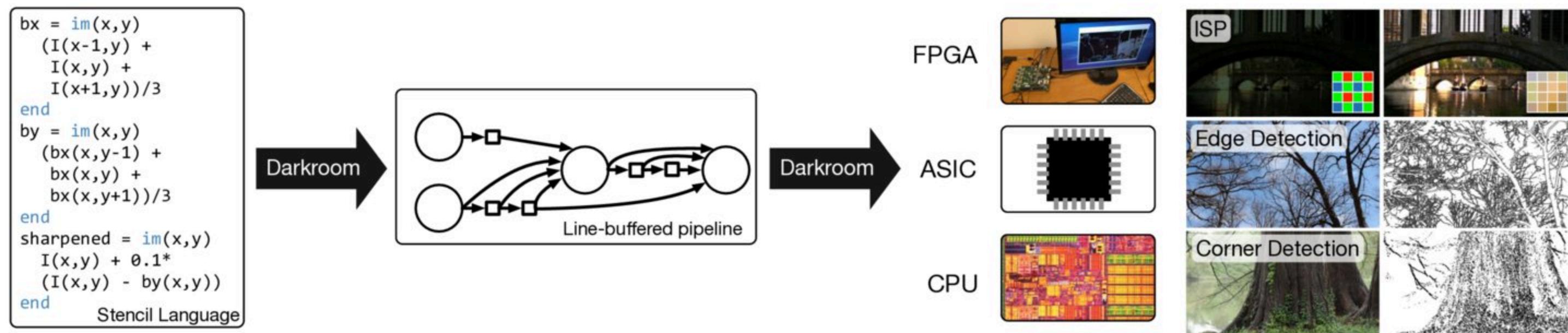
Early results from [Mullapudi 2016]



Darkroom/Rigel/Aetherling

[Hegarty 2014, Hegarty 2016, Durst 2020]

Goal: directly synthesize ASIC or FPGA implementation of image processing pipelines from a high-level algorithm description (a constrained “Halide-like” language)



Goal: very-high efficiency image processing

Many other recent domain-specific programming systems



Less domain specific than examples given today,
but still designed specifically for:
data-parallel computations on big data for
distributed systems (“Map-Reduce”)



DSL for graph-based machine learning computations
Also see Ligra
(DSLs for describing operations on graphs)



Model-view-controller paradigm for
web-applications



Language for real-time 3D graphics



DSL for defining deep neural
networks and training/inference
computations on those networks



Numerical computing

Ongoing efforts in many domains...

Languages for physical simulation: Simit [MIT], Ebb [Stanford]

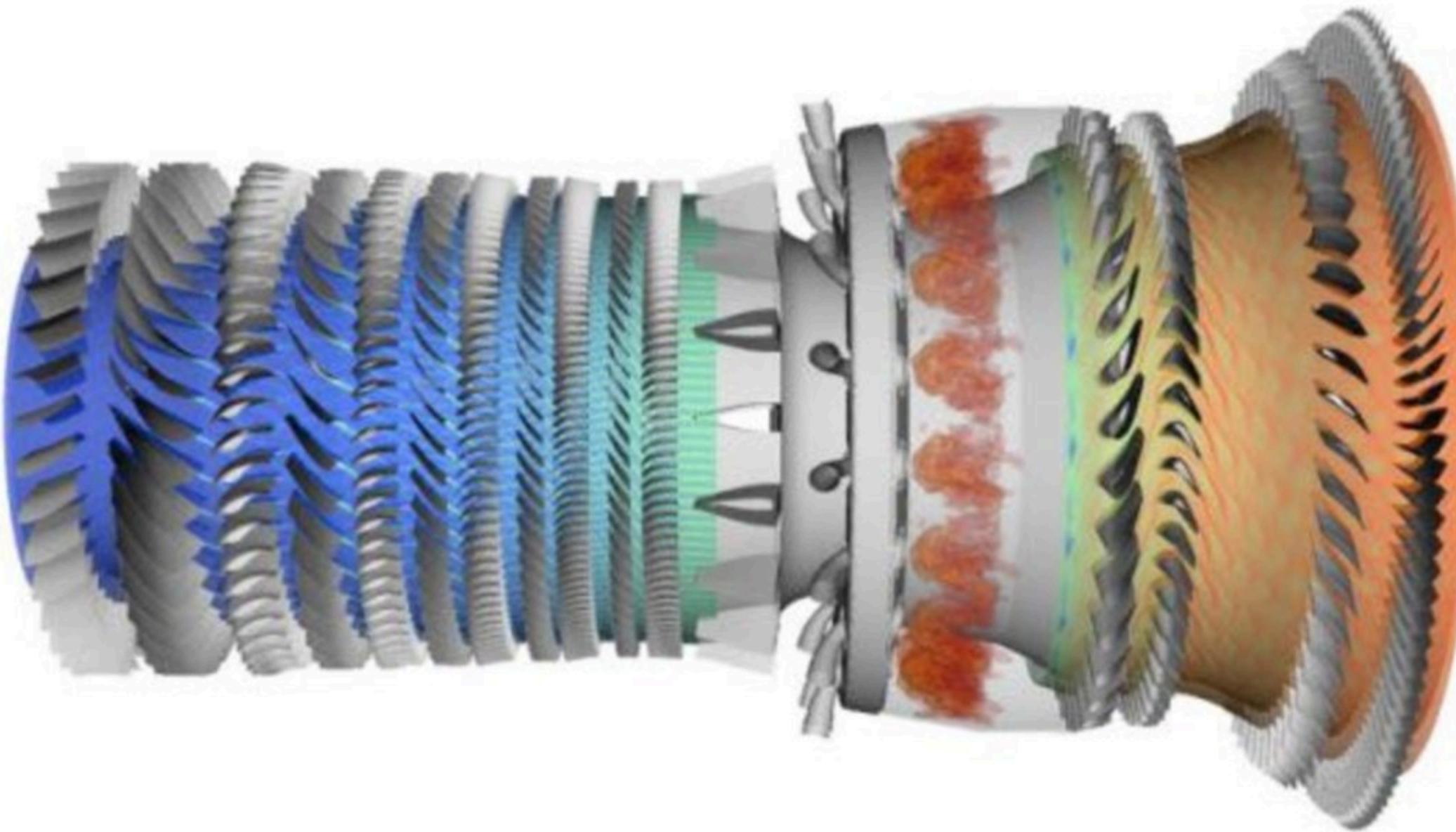
Opt: a language for non-linear least squares optimization [Stanford]

Summary

- **Modern machines: parallel and heterogeneous**
 - Only way to increase compute capability in energy-constrained world
- **Most software uses small fraction of peak capability of machine**
 - Very challenging to tune programs to these machines
 - Tuning efforts are not portable across machines
- **Domain-specific programming environments trade-off generality to achieve productivity, performance, and portability**
 - Case study today: Halide
 - Leverage explicit dependencies, domain restrictions, domain knowledge for system to synthesize efficient implementations

Another DSL example: Lizst: a language for solving PDE's on meshes

[DeVito et al. Supercomputing 11, SciDac'11]



Slide credit for this section of lecture:

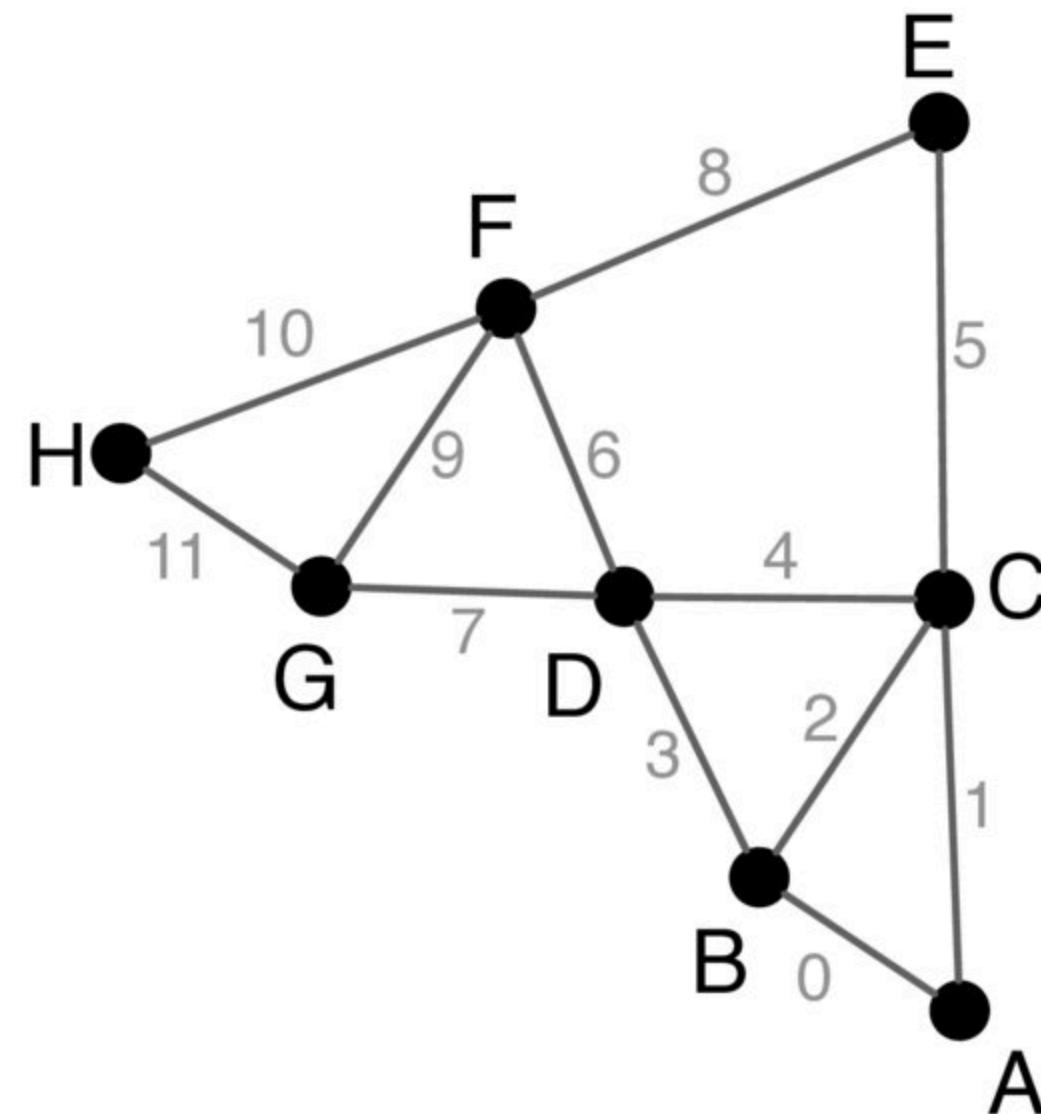
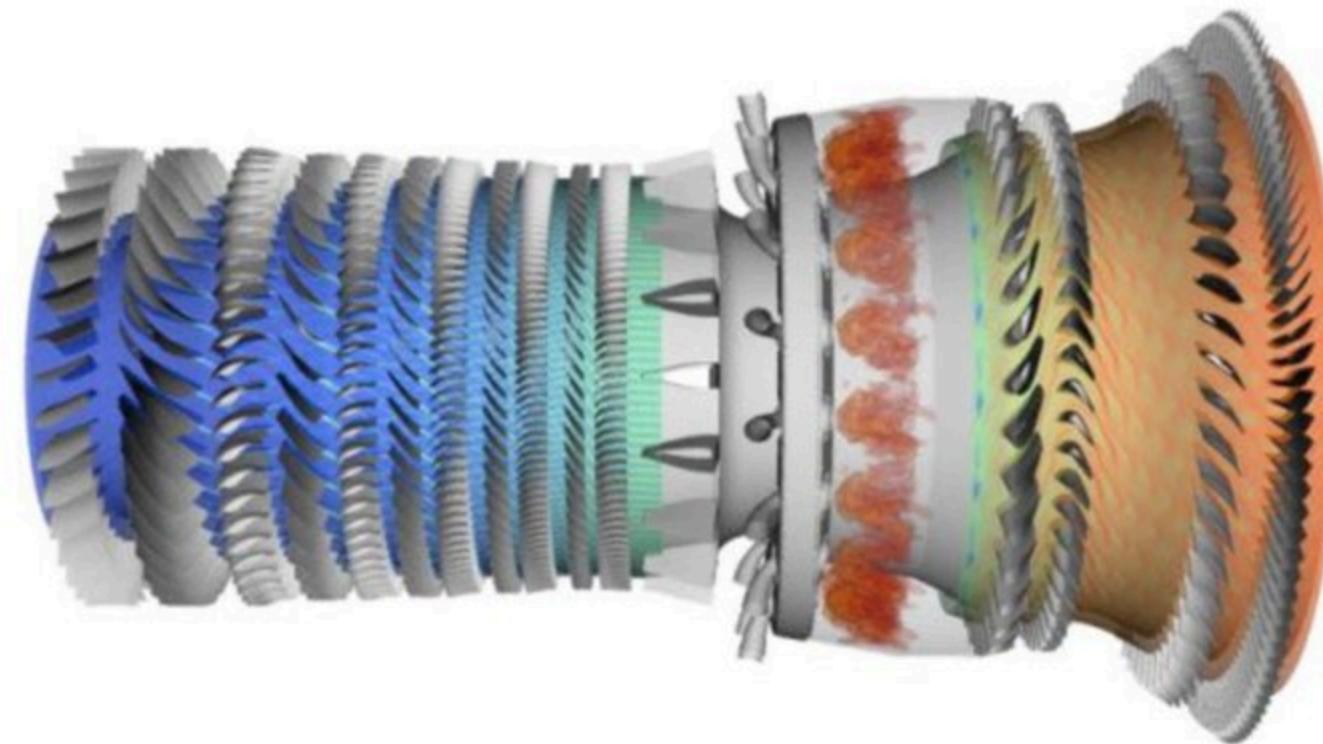
Pat Hanrahan and Zach DeVito (Stanford)

<http://liszt.stanford.edu/>

What a Liszt program does

A Liszt program is run on a mesh:

A Liszt program computes the value of fields defined on mesh faces, edges, or vertices



Liszt program: heat conduction on mesh

Program computes the value of fields defined on meshes

```
var i = 0;
while ( i < 1000 ) {
    Flux(vertices(mesh)) = 0.f;
    JacobiStep(vertices(mesh)) = 0.f;
    for (e <- edges(mesh)) {
        val v1 = head(e)
        val v2 = tail(e)
        val dP = Position(v1) - Position(v2)
        val dT = Temperature(v1) - Temperature(v2)
        val step = 1.0f/(length(dP))
        Flux(v1) += dT*step
        Flux(v2) -= dT*step
        JacobiStep(v1) += step
        JacobiStep(v2) += step
    }
    i += 1
}
```

Given edge, loop body accesses/modifies field values at adjacent mesh vertices

Set flux for all vertices to 0.f;

Independently, for each edge in the mesh



Access value of field at mesh vertex v2

Color key:

Fields

Mesh

Topology functions

Iteration over set

Liszt programming

- A Liszt program describes operations on fields of an abstract mesh representation
- Application specifies type of mesh (regular, irregular) and its topology
- Mesh representation is chosen by Liszt (not by the programmer)
 - Based on mesh type, program behavior, and target machine



Well, that's interesting. I write a program, and the compiler decides what data structure it should use based on what operations my code performs.

Compiling to parallel computers

Recall challenges you have faced in your assignments

- 1. Identify parallelism**
- 2. Identify data locality**
- 3. Reason about what synchronization is required**

Now consider how to automate this process in the Liszt compiler.

Key: determining program dependencies

1. Identify parallelism

- Absence of dependencies implies code can be executed in parallel

2. Identify data locality

- Partition data based on dependencies

3. Reason about required synchronization

- Synchronization is needed to respect dependencies (must wait until the values a computation depends on are known)

In general programs, compilers are unable to infer dependencies at global scale:

Consider: $a[f(i)] += b[i];$

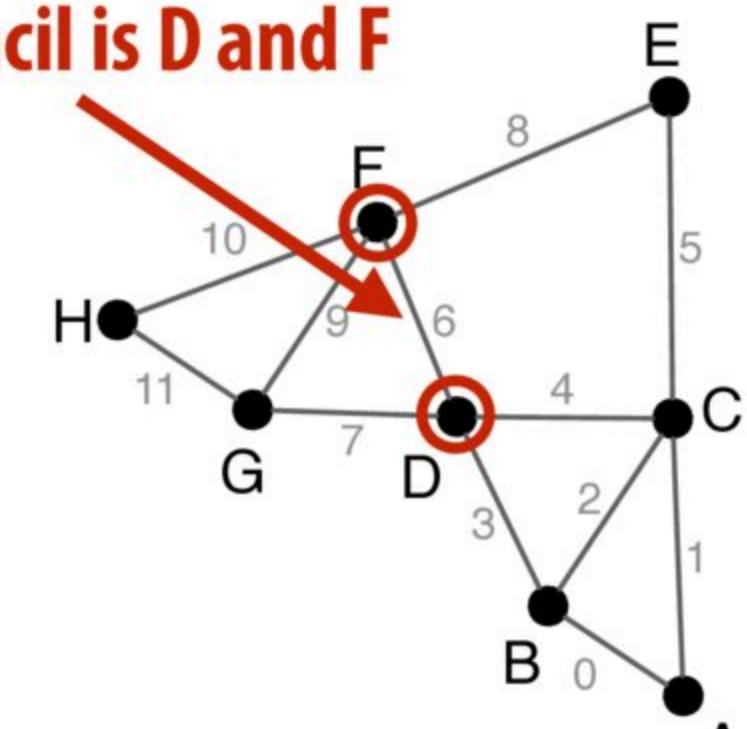
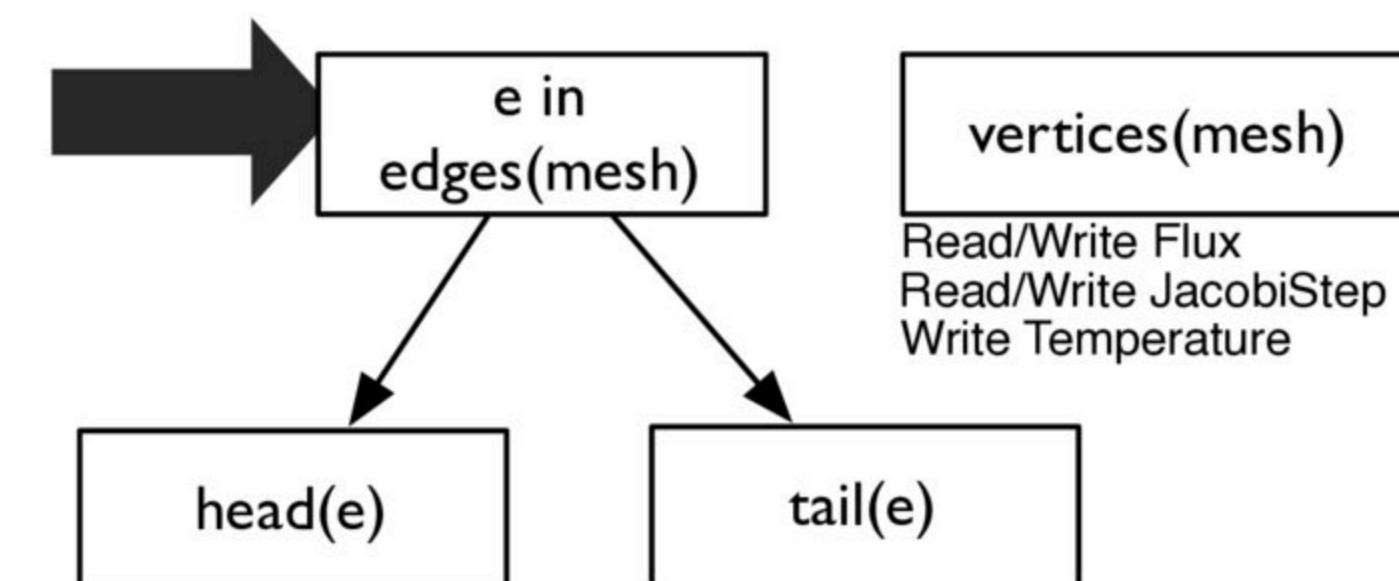
(must execute $f(i)$ to know if dependency exists across loop iterations i)

Liszt is constrained to allow dependency analysis

Statically analyze code to find stencil of each top-level `for` loop

- Extract nested mesh element reads
 - Extract operations on data at mesh elements

```
for (e <- edges(mesh)) {  
    val v1 = head(e)  
    val v2 = tail(e)  
    val dP = Position(v1) - Position(v2)  
    val dT = Temperature(v1) - Temperature(v2)  
    val step = 1.0f/length(dP))  
    Flux(v1) += dT*step  
    Flux(v2) -= dT*step  
    JacobiStep(v1) += step  
    JacobiStep(v2) += step  
}  
...
```



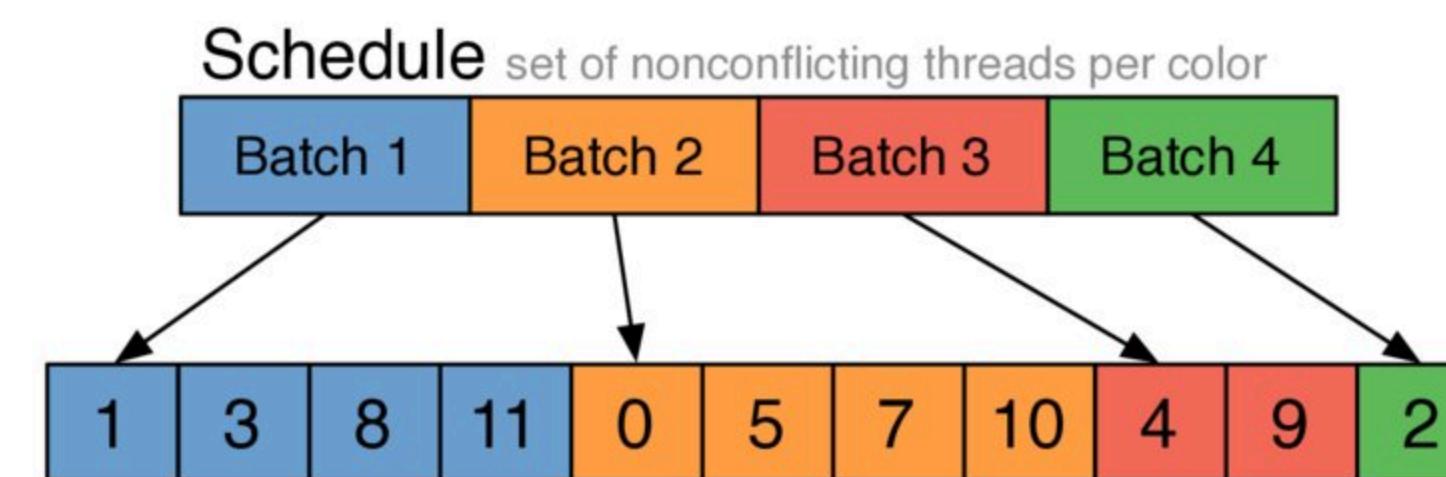
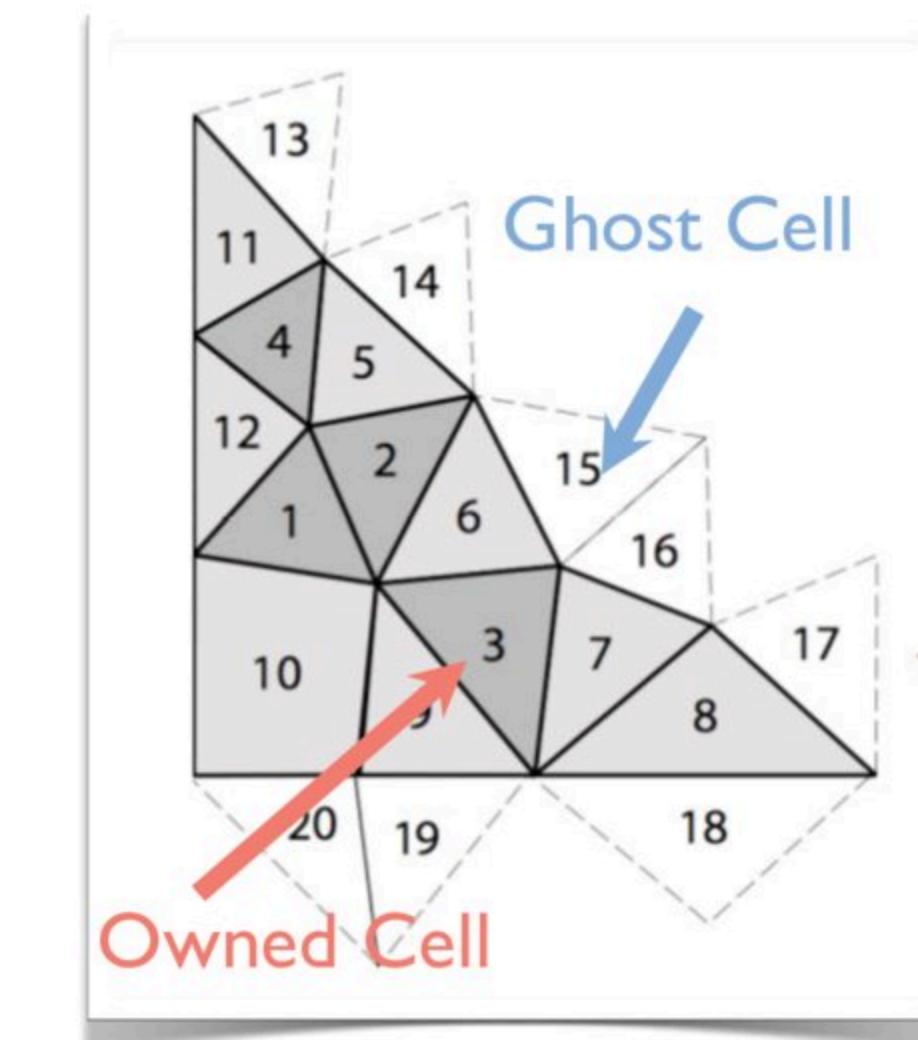
Edge 6's read stencil is D and F

Portable parallelism: compiler uses knowledge of dependencies to implement different parallel execution strategies

I'll discuss two strategies...

Strategy 1: mesh partitioning

Strategy 2: mesh coloring



Imagine compiling a Liszt program to a cluster

(multiple nodes, distributed address space)

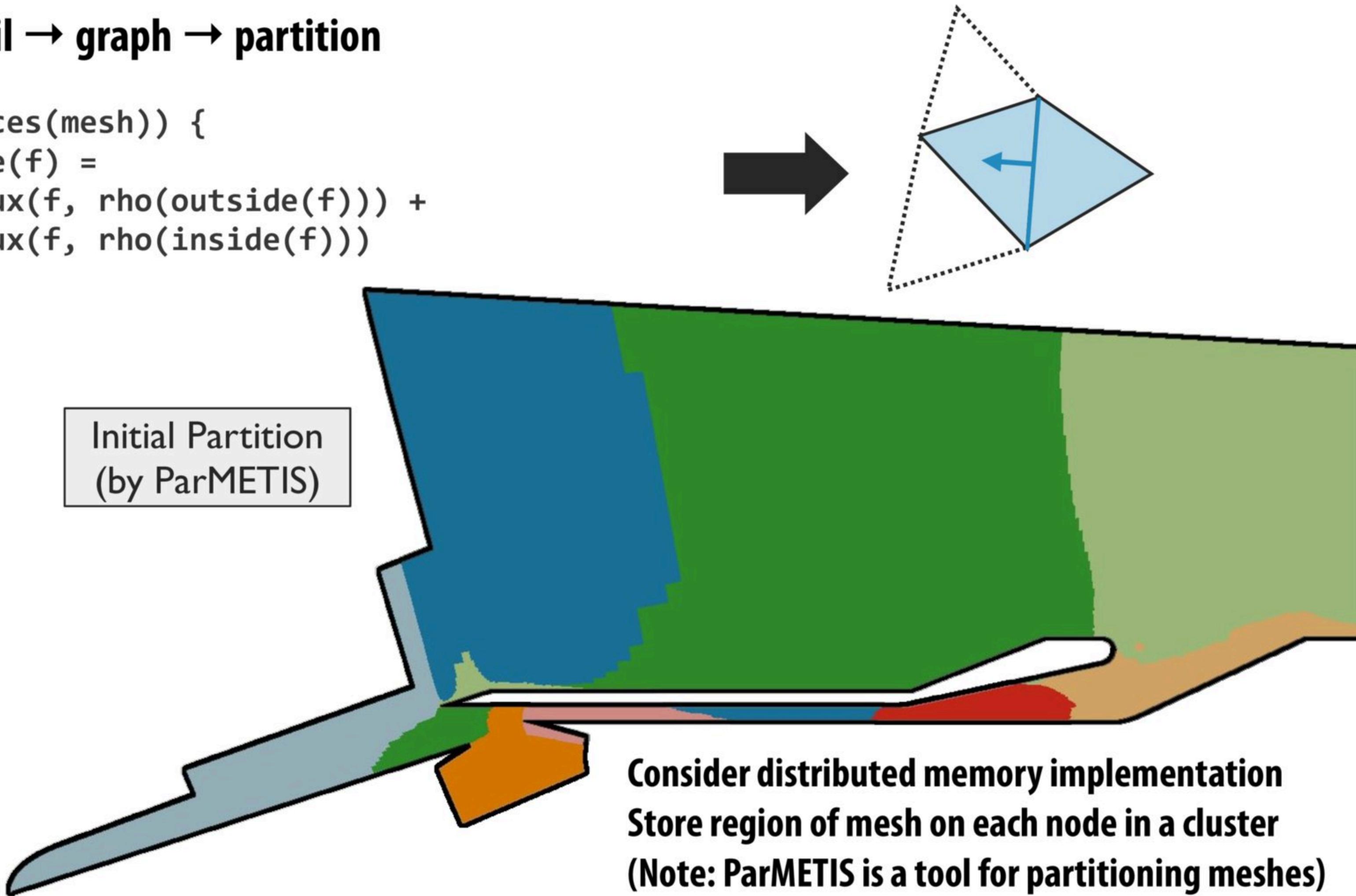
How might Liszt distribute a graph across these nodes?

- Must access mesh elements relative to some input vertex, edge, face, etc.)
- Notice how many operators return sets (e.g., “all edges of this face”)

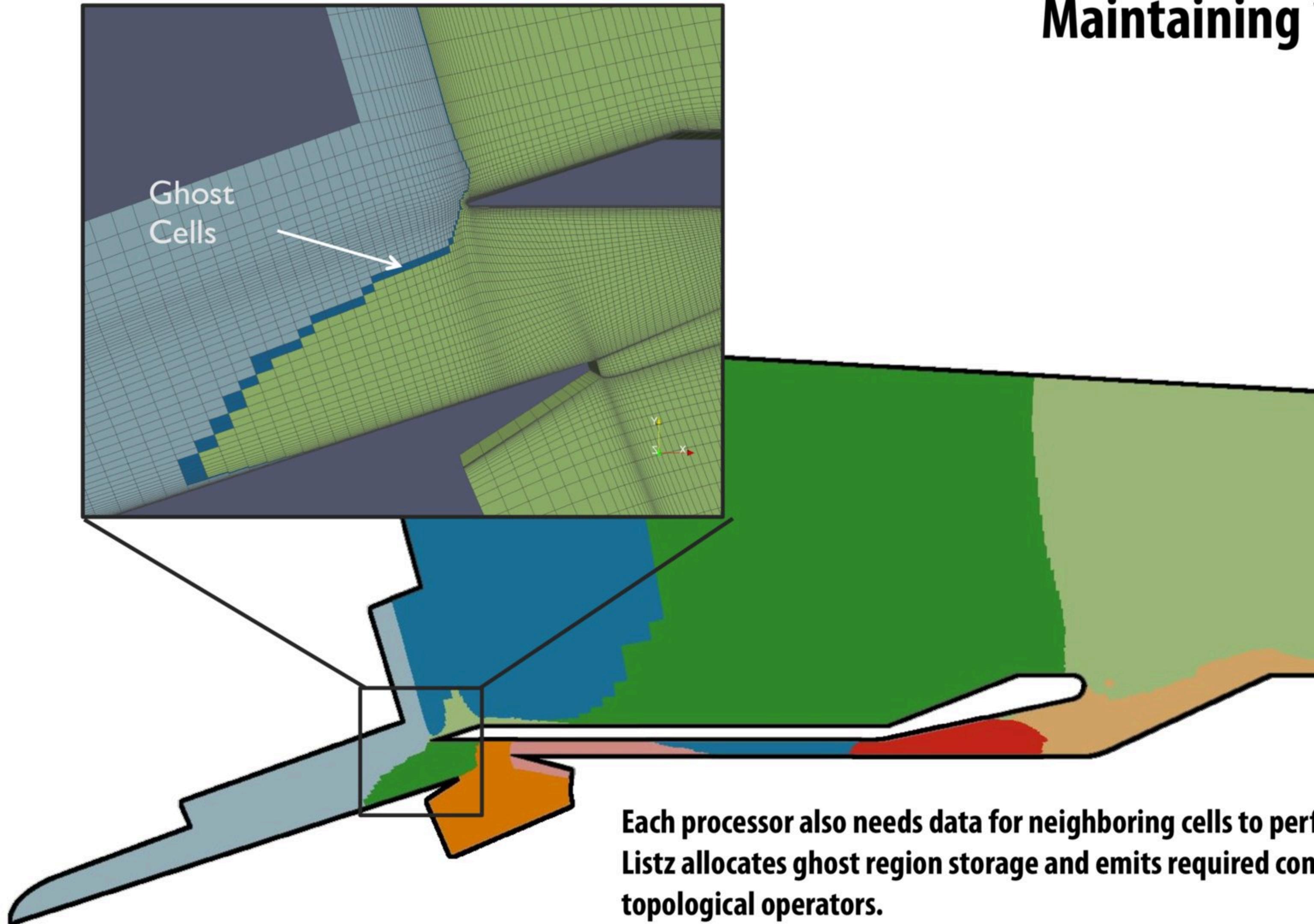
Distributed memory implementation of Liszt

Mesh + stencil → graph → partition

```
for(f <- faces(mesh)) {  
    rhoOutside(f) =  
        calc_flux(f, rho(outside(f))) +  
        calc_flux(f, rho(inside(f)))  
}
```



Maintaining 1-Level Ghost Cells



Imagine compiling a Lizst program to a GPU

- Used to access mesh elements relative to some input vertex, edge, face, etc.)
- Notice how many operators return sets (e.g., “all edges of this face”)

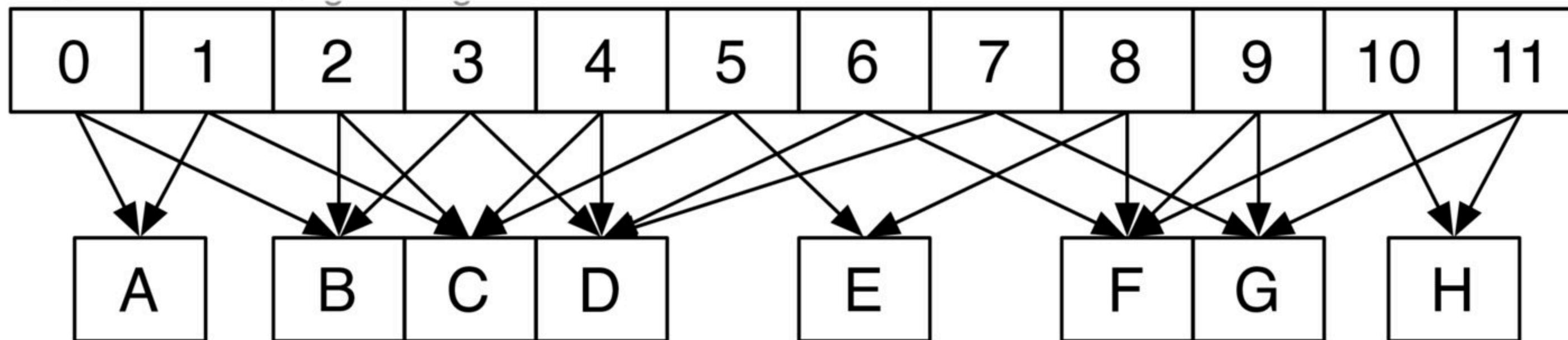
(single address space, many tiny threads)

GPU implementation: parallel reductions

In previous example, one region of mesh assigned per processor (or node in cluster)

On GPU, natural parallelization is one edge per CUDA thread

Edges (each edge assigned to 1 CUDA thread)



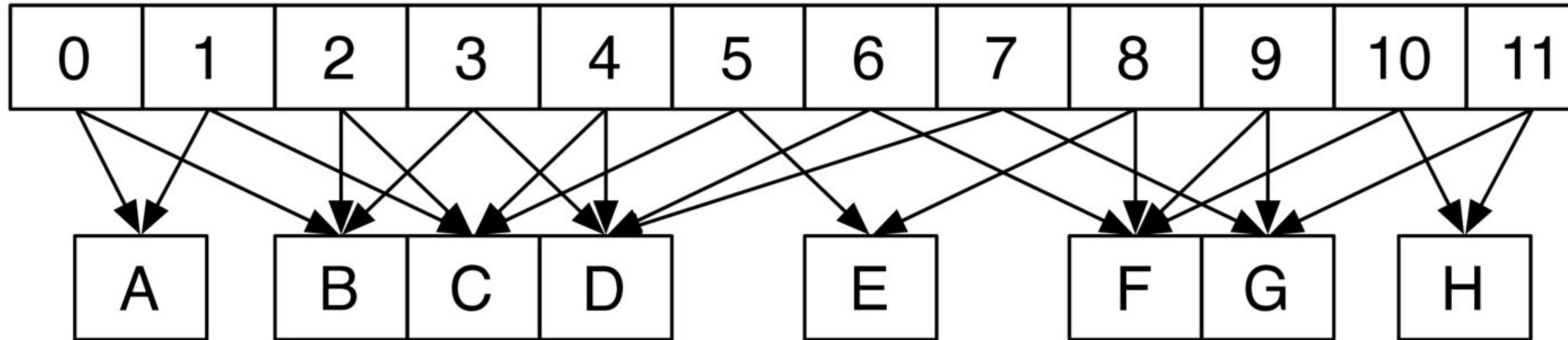
Flux field values (stored per vertex)

```
for (e <- edges(mesh)) {  
    ...  
    Flux(v1) += dT*step  
    Flux(v2) -= dT*step  
    ...  
}
```

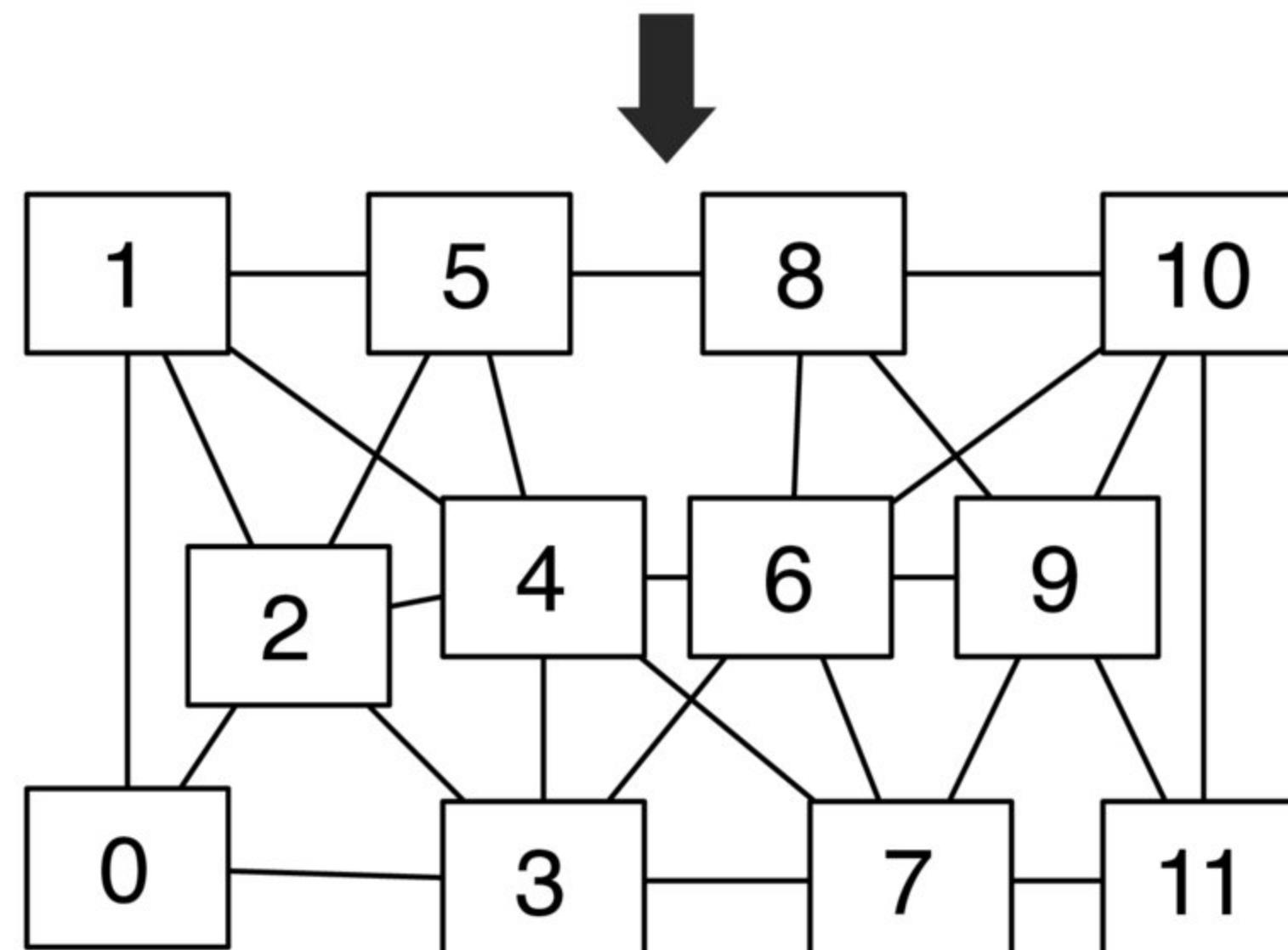
Different edges share a vertex: requires
atomic update of per-vertex field data

GPU implementation: conflict graph

Edges (each edge assigned to 1 CUDA thread)



Flux field values (per vertex)



Identify mesh edges with colliding writes
(lines in graph indicate presence of collision)

Can simply run program once to get this information.
(results remain valid for subsequent executions provided mesh does not change)