

(how to be l33t)

Lecture 6:

~~Performance Optimization~~ Part II:
Locality, Communication, and Contention

Parallel Programming
Stanford CS149, Fall 2020

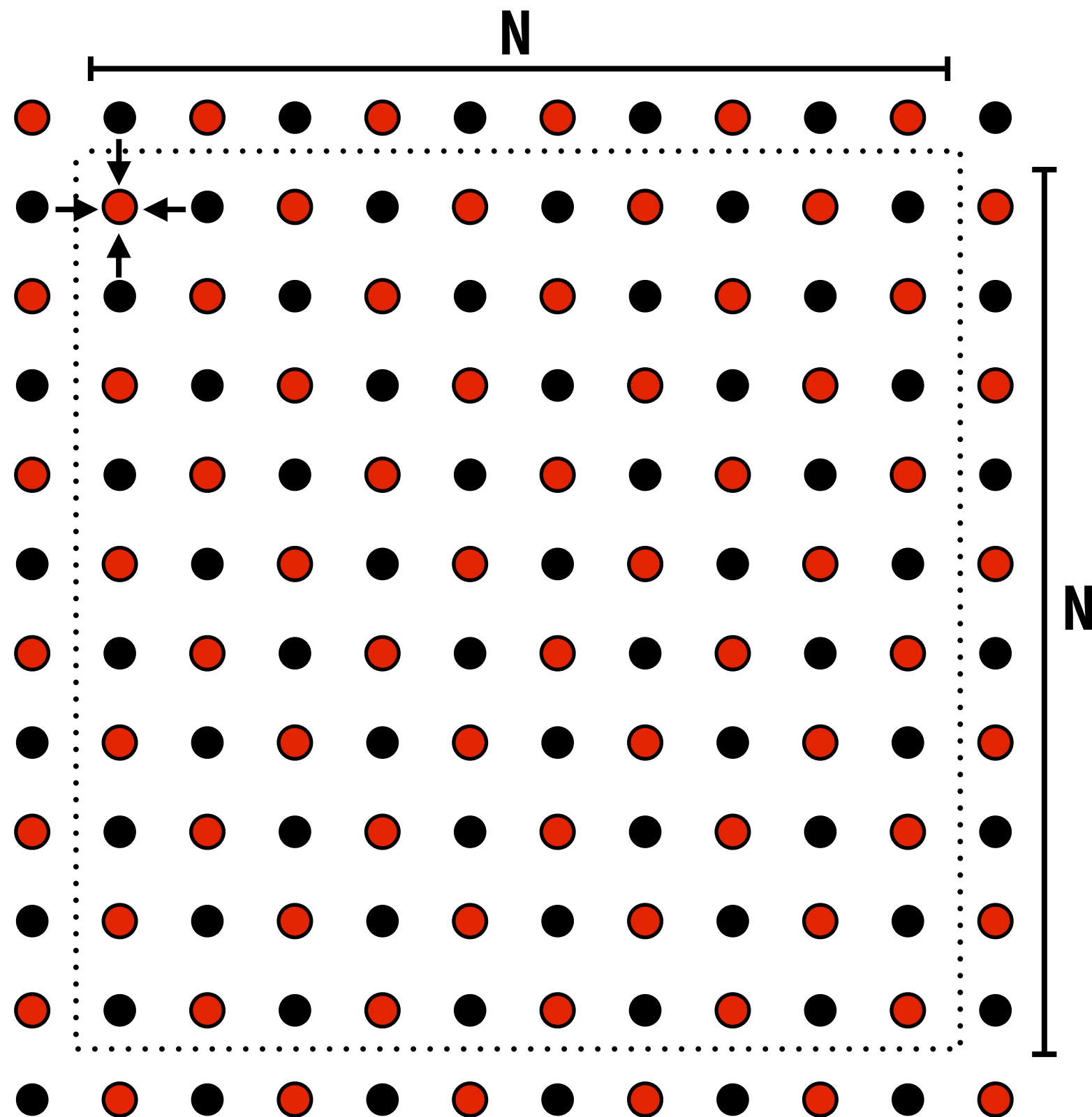
Today: more parallel program optimization

- **Last lecture: strategies for assigning work to workers (threads, processors, etc.)**
 - **Goal: achieving good workload balance while also minimizing overhead**
 - **Discussed tradeoffs between static and dynamic work assignment**
 - **Reminder: keep it simple (implement, analyze, then tune/optimize if required)**
- **Today: strategies for minimizing communication costs**

Warm up:
communication using message passing
(since it makes communication explicit)

Recall the grid-based solver example

In previous lectures we expressed this parallel program using data parallel and SPMD programming abstractions

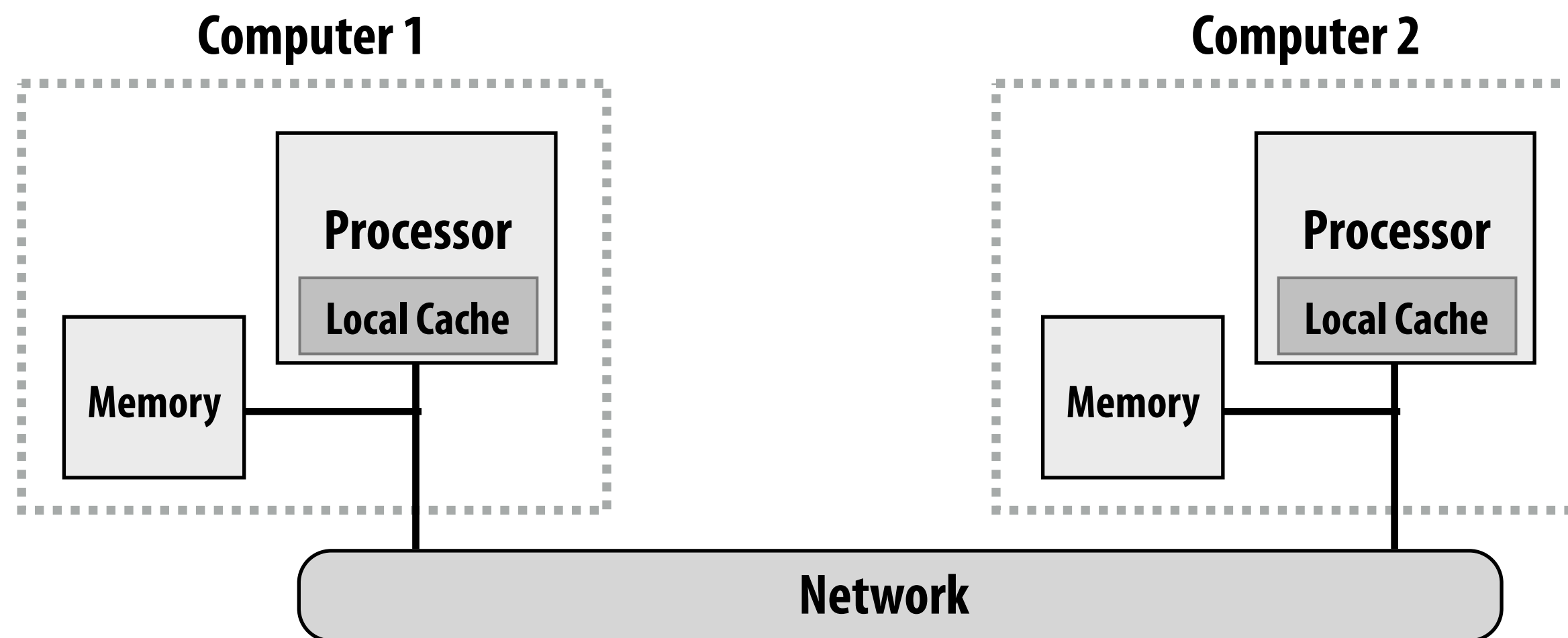


```
int N;
float* A = allocate(n+2, n+2);

void solve(float* A) {
  bool done = false;
  float diff = 0.f;
  while (!done) {
    for_all (red cells (i,j)) {
      float prev = A[i,j];
      A[i,j] = 0.2f * (A[i-1,j] + A[i,j-1] + A[i,j] +
                    A[i+1,j] + A[i,j+1]);
      reduceAdd(diff, abs(A[i,j] - prev));
    }
    if (diff/(N*N) < TOLERANCE)
      done = true;
  }
}
```

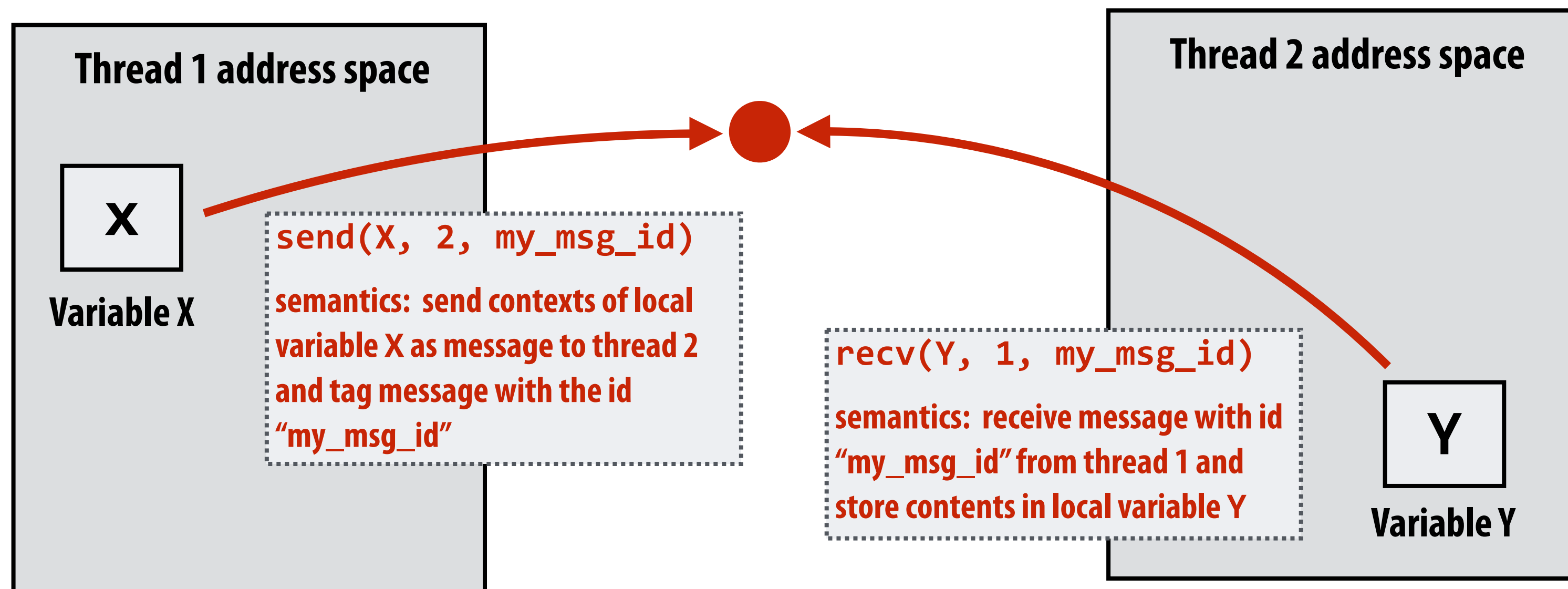

Let's think about expressing a parallel grid solver with communication via messages

One possible message passing machine configuration:
a cluster of two machines



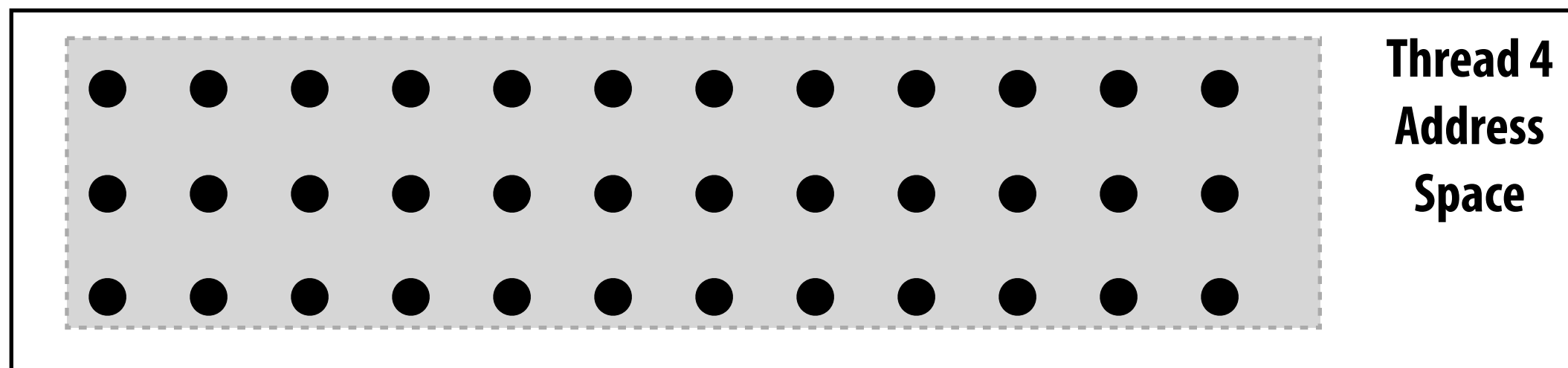
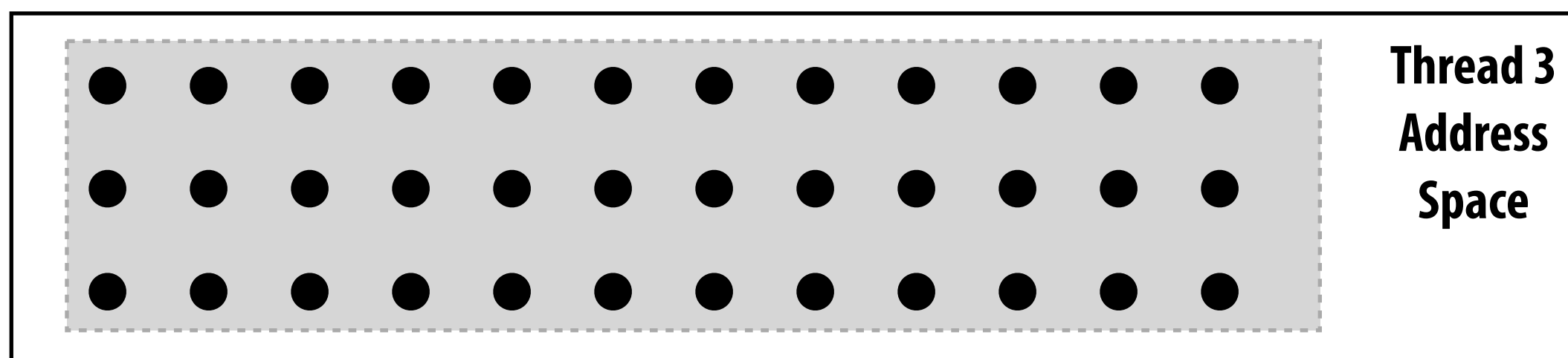
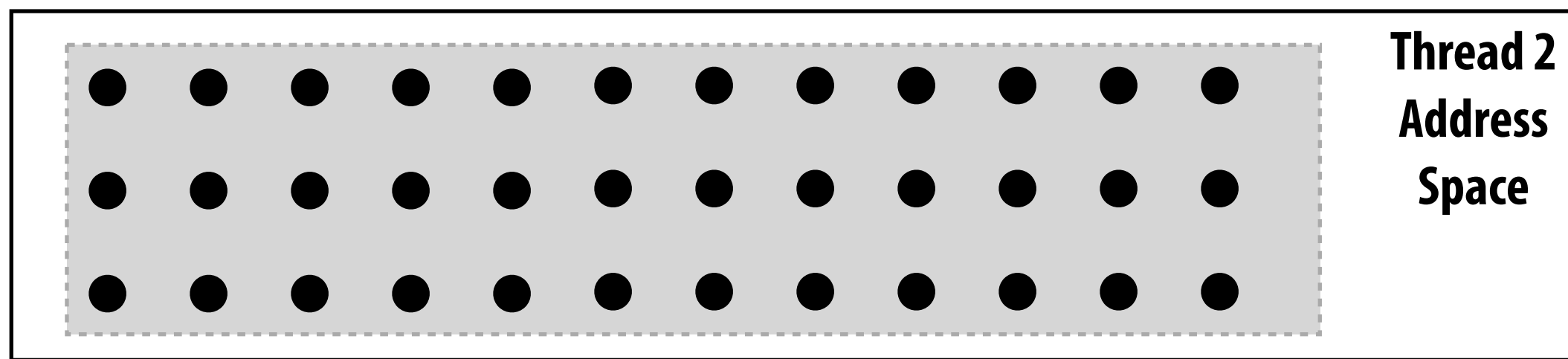
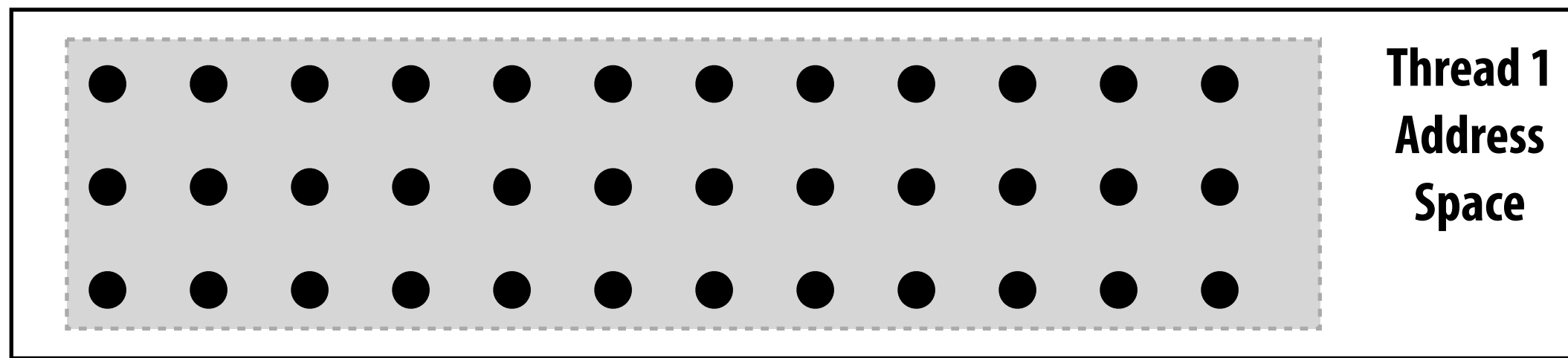
Review: message passing model

- Threads operate within their own private address spaces
- Threads communicate by sending/receiving messages
 - send: specifies recipient, buffer to be transmitted, and optional message identifier (“tag”)
 - receive: sender, specifies buffer to store data, and optional message identifier
 - Sending messages is the only way to exchange data between threads 1 and 2
 - Why?



(Communication operations shown in red)

Message passing model: each thread operates in its own address space

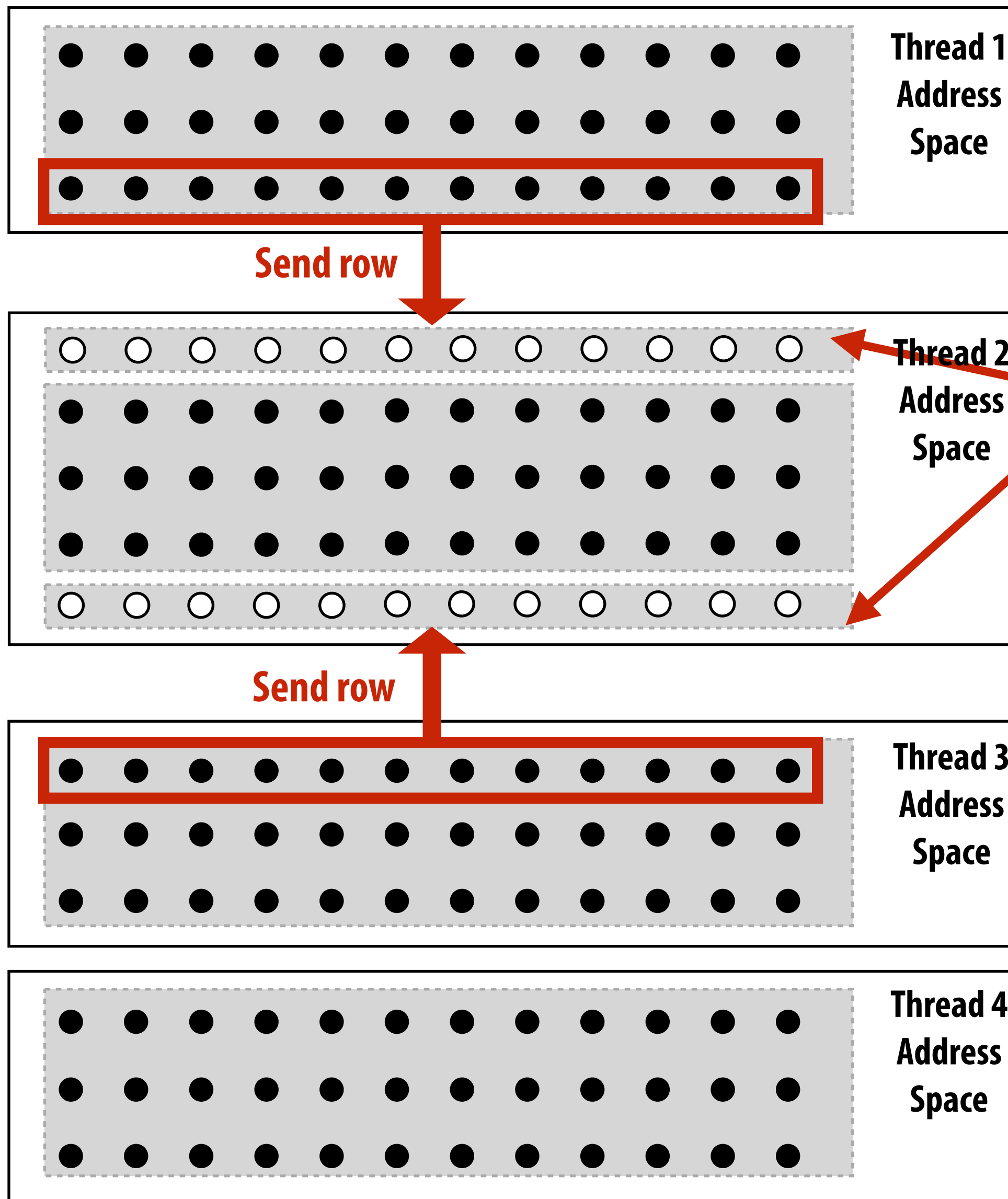


In this figure: four threads

The grid data is partitioned into four allocations, each residing in one of the four unique thread address spaces

(four per-thread private arrays)

Data replication is now required to correctly execute the program



Example:

After processing of red cells is complete, thread 1 and thread 3 send one row of data to thread 2 (thread 2 requires up-to-date red cell information to update black cells in the next phase)

"Ghost cells" are grid cells replicated from a remote address space. It's common to say that information in ghost cells is "owned" by other threads.

Thread 2 logic:

```
float* local_data = allocate(N+2, rows_per_thread+2);

int tid = get_thread_id();
int bytes = sizeof(float) * (N+2);

// receive ghost row cells (white dots)
recv(&local_data[0], bytes, tid-1);
recv(&local_data[rows_per_thread+1], bytes, tid+1);

// Thread 2 now has data necessary to perform
// its future computation
```

Message passing solver

Similar structure to shared address space solver, but now communication is explicit in message sends and receives

Send and receive ghost rows to “neighbor threads”

Perform computation

(just like in shared address space version of solver)

All threads send local `my_diff` to thread 0

Thread 0 computes global diff, evaluates termination predicate and sends result back to all other threads

```
int N;
int tid = get_thread_id();
int rows_per_thread = N / get_num_threads();

float* localA = allocate(rows_per_thread+2, N+2);

// assume localA is initialized with starting values
// assume MSG_ID_ROW, MSG_ID_DONE, MSG_ID_DIFF are constants used as msg ids

////////////////////////////////////

void solve() {
    bool done = false;
    while (!done) {

        float my_diff = 0.0f;

        if (tid != 0)
            send(&localA[1,0], sizeof(float)*(N+2), tid-1, MSG_ID_ROW);
        if (tid != get_num_threads()-1)
            send(&localA[rows_per_thread,0], sizeof(float)*(N+2), tid+1, MSG_ID_ROW);

        if (tid != 0)
            recv(&localA[0,0], sizeof(float)*(N+2), tid-1, MSG_ID_ROW);
        if (tid != get_num_threads()-1)
            recv(&localA[rows_per_thread+1,0], sizeof(float)*(N+2), tid+1, MSG_ID_ROW);

        for (int i=1; i<rows_per_thread+1; i++) {
            for (int j=1; j<n+1; j++) {
                float prev = localA[i,j];
                localA[i,j] = 0.2 * (localA[i-1,j] + localA[i,j] + localA[i+1,j] +
                                   localA[i,j-1] + localA[i,j+1]);
                my_diff += fabs(localA[i,j] - prev);
            }
        }

        if (tid != 0) {
            send(&mydiff, sizeof(float), 0, MSG_ID_DIFF);
            recv(&done, sizeof(bool), 0, MSG_ID_DONE);
        } else {
            float remote_diff;
            for (int i=1; i<get_num_threads()-1; i++) {
                recv(&remote_diff, sizeof(float), i, MSG_ID_DIFF);
                my_diff += remote_diff;
            }
            if (my_diff/(N*N) < TOLERANCE)
                done = true;
            for (int i=1; i<get_num_threads()-1; i++)
                send(&done, sizeof(bool), i, MSG_ID_DONE);
        }
    }
}
```

Notes on message passing example

■ Computation

- Array indexing is relative to local address space

■ Communication:

- Performed by sending and receiving messages
- Bulk transfer: communicate entire rows at a time

■ Synchronization:

- Performed by sending and receiving messages
- Consider how to implement mutual exclusion, barriers, flags using messages

Synchronous (blocking) send and receive

- **send(): call returns when sender receives acknowledgement that message data resides in address space of receiver**
- **recv(): call returns when data from received message is copied into address space of receiver and acknowledgement sent back to sender**

Sender:

Receiver:

Call SEND(foo)

Copy data from buffer 'foo' in sender's address space into network buffer

Send message



Receive ack



SEND() returns

Call RECV(bar)

Receive message

Copy data into buffer 'bar' in receiver's address space

Send ack

RECV() returns

As implemented on the prior slide, there is a big problem with our message passing solver if it uses synchronous send/recv!

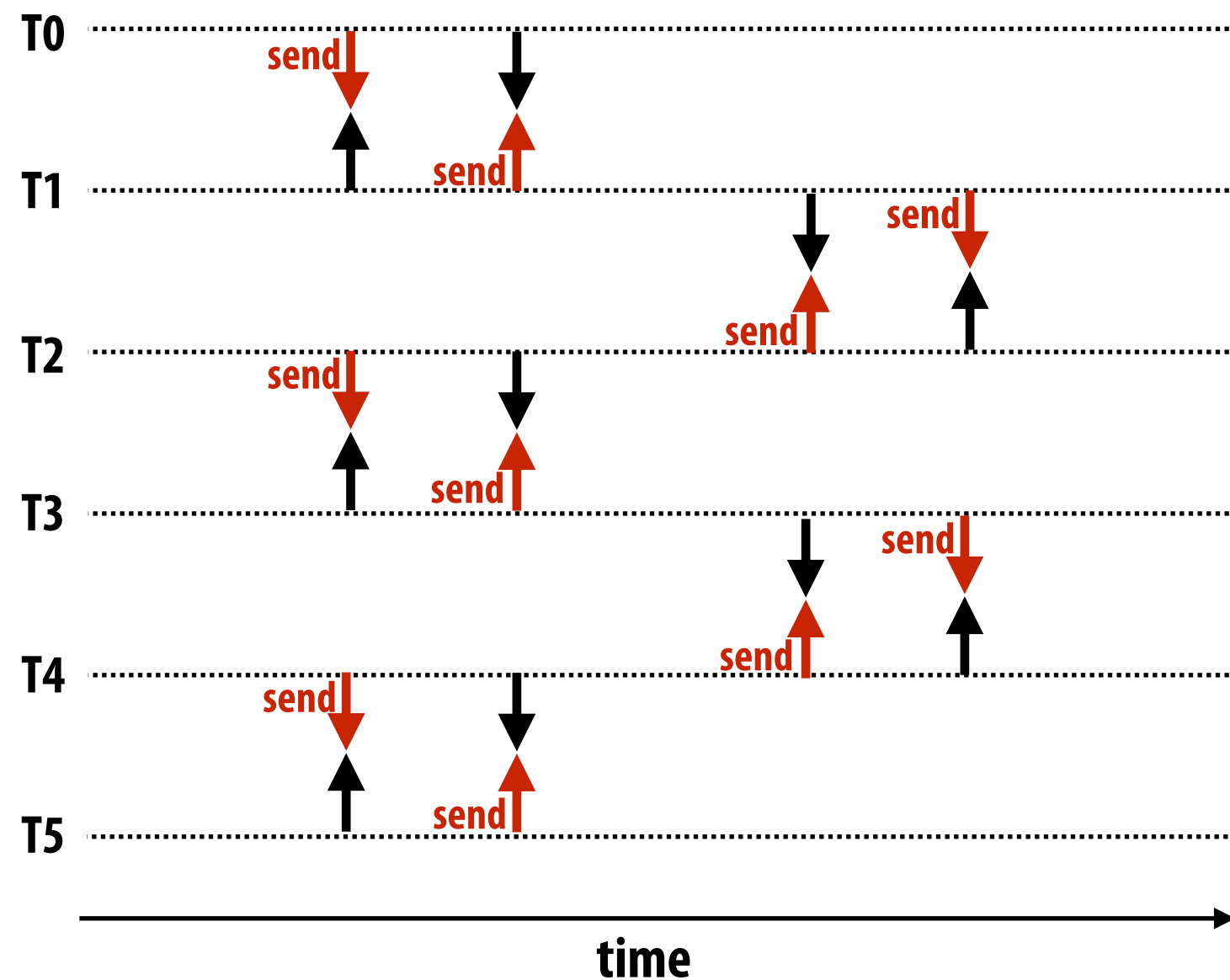
Why?

How can we fix it?

(while still using synchronous send/recv)

Message passing solver (fixed to avoid deadlock)

Send and receive ghost rows to "neighbor threads"
Even-numbered threads send, then receive
Odd-numbered thread rcv, then send



```

int N;
int tid = get_thread_id();
int rows_per_thread = N / get_num_threads();

float* localA = allocate(rows_per_thread+2, N+2);

// assume localA is initialized with starting values
// assume MSG_ID_ROW, MSG_ID_DONE, MSG_ID_DIFF are constants used as msg ids

////////////////////////////////////

void solve() {
    bool done = false;
    while (!done) {

        float my_diff = 0.0f;

        if (tid % 2 == 0) {
            sendDown(); rcvDown();
            sendUp();   rcvUp();
        } else {
            rcvUp();   sendUp();
            rcvDown(); sendDown();
        }

        for (int i=1; i<rows_per_thread-1; i++) {
            for (int j=1; j<n+1; j++) {
                float prev = localA[i,j];
                localA[i,j] = 0.2 * (localA[i-1,j] + localA[i,j] + localA[i+1,j] +
                                     localA[i,j-1] + localA[i,j+1]);
                my_diff += fabs(localA[i,j] - prev);
            }
        }

        if (tid != 0) {
            send(&mydiff, sizeof(float), 0, MSG_ID_DIFF);
            rcv(&done, sizeof(bool), 0, MSG_ID_DONE);
        } else {
            float remote_diff;
            for (int i=1; i<get_num_threads()-1; i++) {
                rcv(&remote_diff, sizeof(float), i, MSG_ID_DIFF);
                my_diff += remote_diff;
            }
            if (my_diff/(N*N) < TOLERANCE)
                done = true;
            if (int i=1; i<gen_num_threads()-1; i++)
                send(&done, sizeof(bool), i, MSD_ID_DONE);
        }
    }
}

```

Non-blocking asynchronous send/recv

- **send(): call returns immediately**
 - Buffer provided to send() cannot be modified by calling thread since message processing occurs concurrently with thread execution
 - Calling thread can perform other work while waiting for message to be sent
- **recv(): posts intent to receive in the future, returns immediately**
 - Use checksend(), checkrecv() to determine actual status of send/receipt
 - Calling thread can perform other work while waiting for message to be received

Sender:

Call SEND(foo)
SEND returns handle h1

Copy data from 'foo' into network buffer
Send message

Call CHECKSEND(h1) // if message sent, now safe for thread to modify 'foo'

Receiver:

Call RECV(bar)
RECV(bar) returns handle h2

Receive message
Messaging library copies data into 'bar'
Call CHECKRECV(h2)
// if received, now safe for thread
// to access 'bar'

RED TEXT = executes concurrently with application thread

When I talk about communication, I'm not just referring to messages between machines. (e.g., in a datacenter)

More examples:

Communication between cores on a chip

Communication between a core and its cache

Communication between a core and memory

Think of a parallel system as an extended memory hierarchy

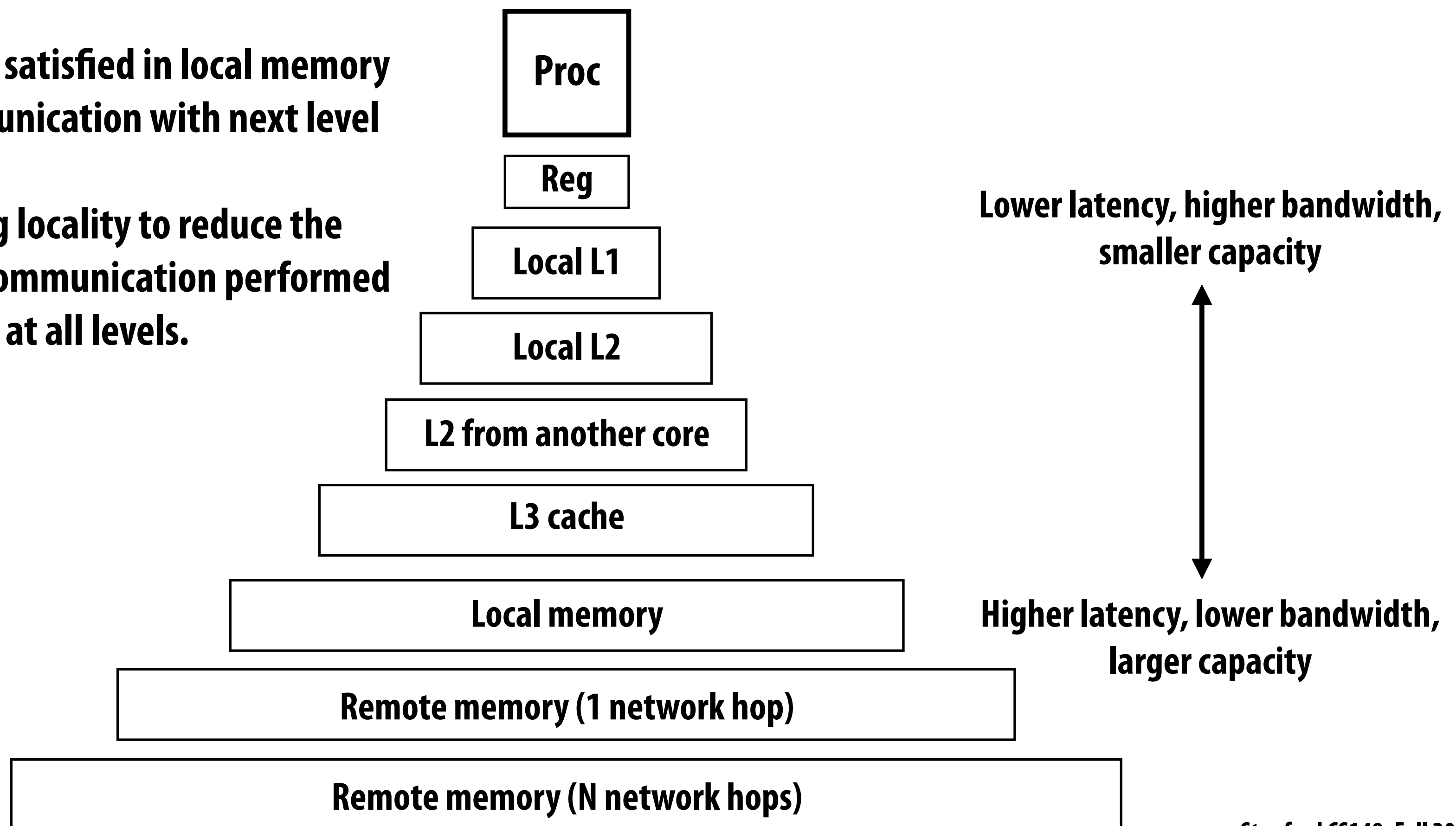
I want you to think of “communication” very generally:

- Communication between a processor and its cache
- Communication between processor and memory (e.g., memory on same machine)
- Communication between processor and a remote memory (e.g., memory on another node in the cluster, accessed by sending a network message)

Accesses not satisfied in local memory cause communication with next level

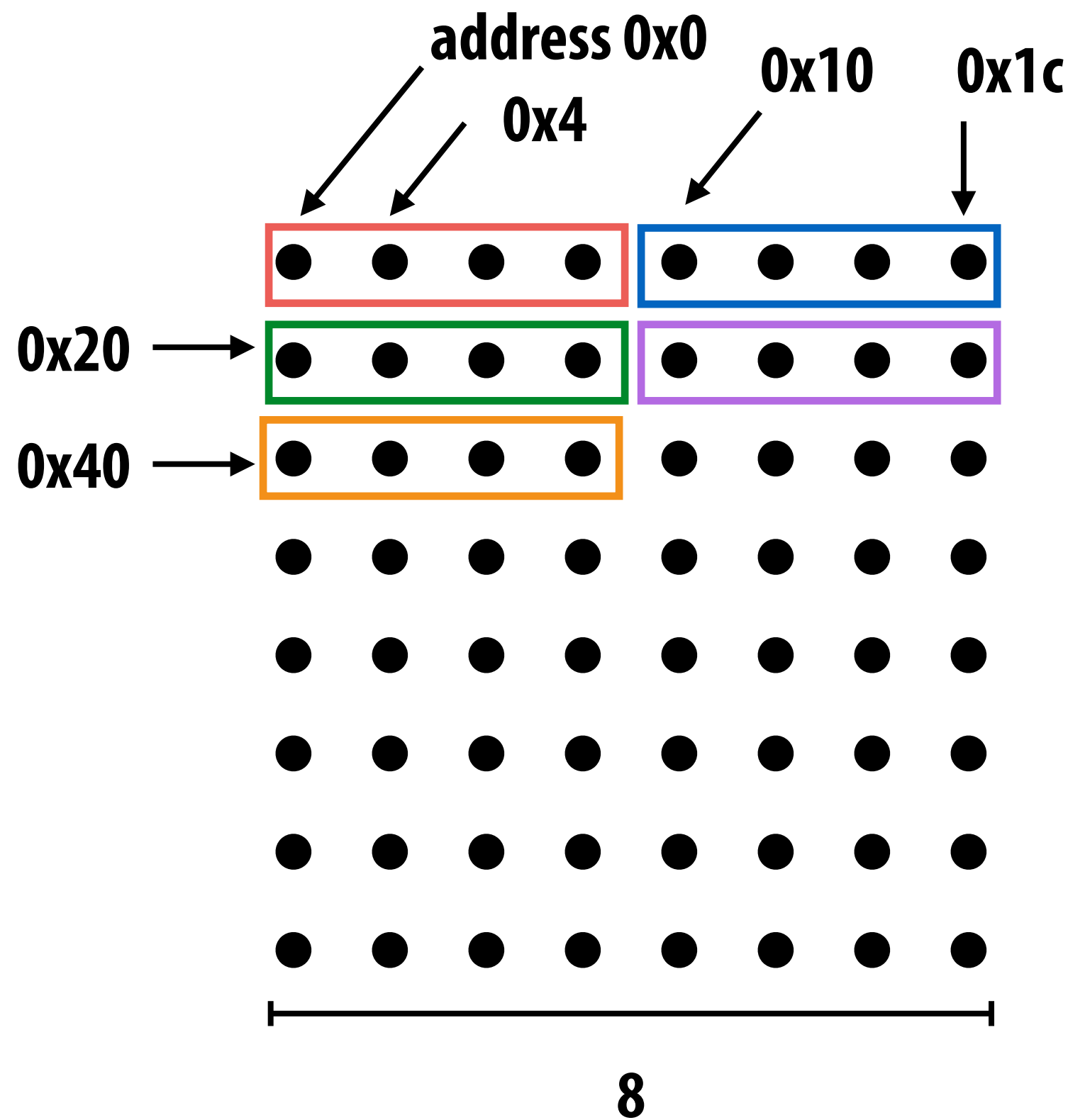
So managing locality to reduce the amount of communication performed is important at all levels.

View from one processor



Review of caches

Cache review



Consider 4-byte elements

Consider a cache with 16-byte cache lines

and a capacity of 32 bytes

(2 lines fit in cache)

Least recently used (LRU) replacement

Address accessed	Cache state (after load is complete)		
0x0	0x0 ●●●●		"cold miss"
0x4	0x0 ●●●●		hit
0x8	0x0 ●●●●		hit
0xc	0x0 ●●●●		hit
0x10	0x0 ●●●●	0x10 ●●●●	cold miss
0x14	0x0 ●●●●	0x10 ●●●●	hit
0x18	0x0 ●●●●	0x10 ●●●●	hit
0x1c	0x0 ●●●●	0x10 ●●●●	hit
0x20	0x20 ●●●●	0x10 ●●●●	cold miss (evict 0x0)
0x24	0x20 ●●●●	0x10 ●●●●	hit
0x28	0x20 ●●●●	0x10 ●●●●	hit
0x2c	0x20 ●●●●	0x10 ●●●●	hit
0x30	0x20 ●●●●	0x30 ●●●●	cold miss (evict 0x10)
0x34	0x20 ●●●●	0x30 ●●●●	hit
0x38	0x20 ●●●●	0x30 ●●●●	hit
0x3c	0x20 ●●●●	0x30 ●●●●	hit
0x40	0x40 ●●●●	0x30 ●●●●	cold miss (evict 0x20)

Review: latency vs throughput

Latency

The amount of time needed for an operation to complete.

A memory load that misses the cache has a latency of 200 cycles

A packet takes 20 ms to be sent from my computer to Google

Asking a question on Piazza gets a response in 10 minutes

Throughput

The rate at which operations are performed.

Memory can provide data to the processor at 25 GB/sec (memory bandwidth)

A communication link can send 10 million messages per second

The TAs answer 50 questions per day on Piazza

A simple model of communication

Example: sending a N-bit message



Link can communicate B (bits/clock)

Latency of first bit to travel to destination = T clocks

How long does it take for an entire message to be sent from processor 1 to processor 2?

One possible answer: $T_0 \times (N / B)$

Another possible answer: $T_0 + N/B$

??

Let's talk about people leaving SF...



San Francisco fog vs. South Bay sun

When it looks like this in SF

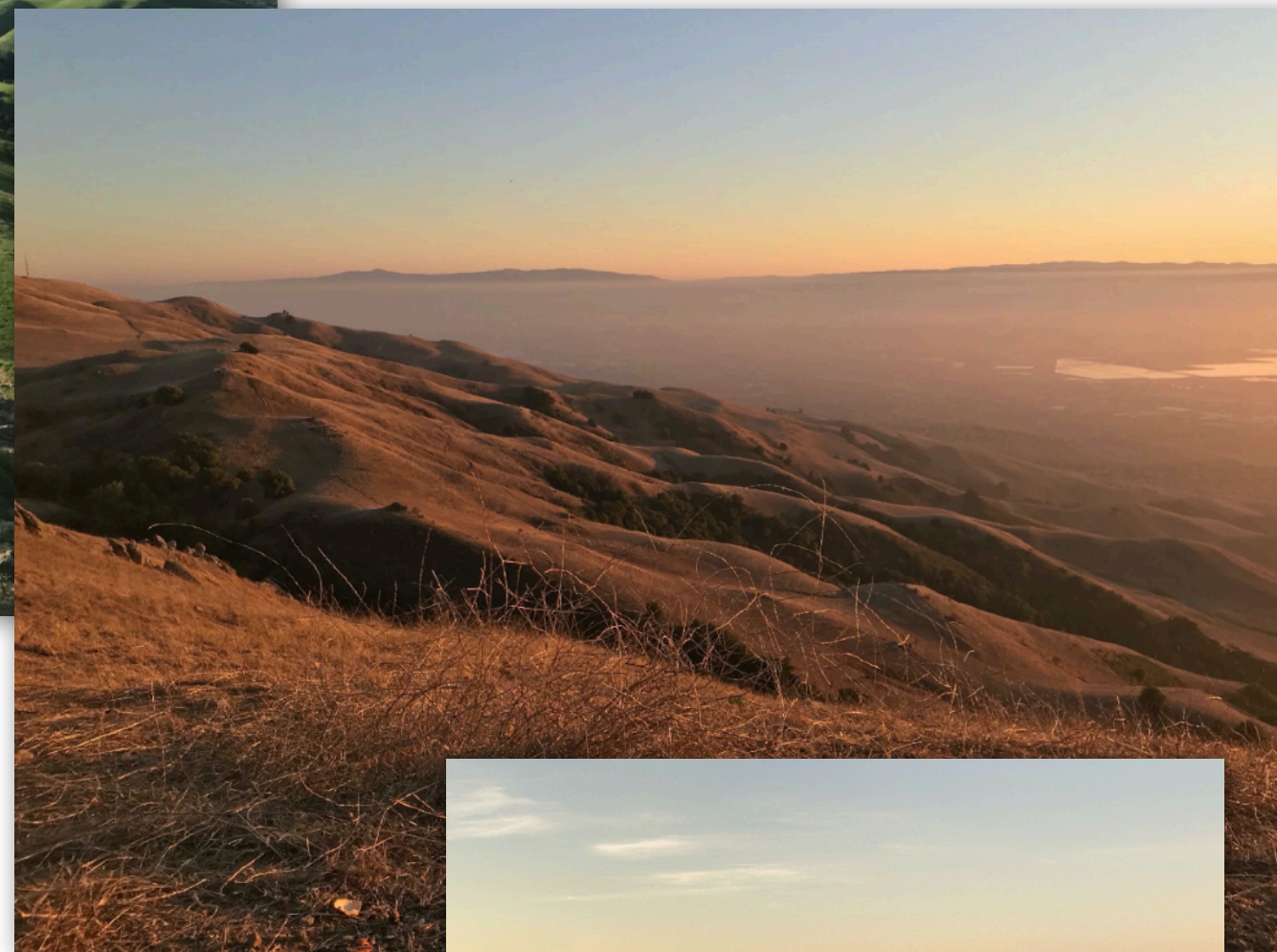


It looks like this at Stanford



Great social distancing opportunities

- Quick plug:
 - Kayvon's guide to local bay area hikes
 - http://graphics.stanford.edu/~kayvonf/misc/local_hikes.pdf

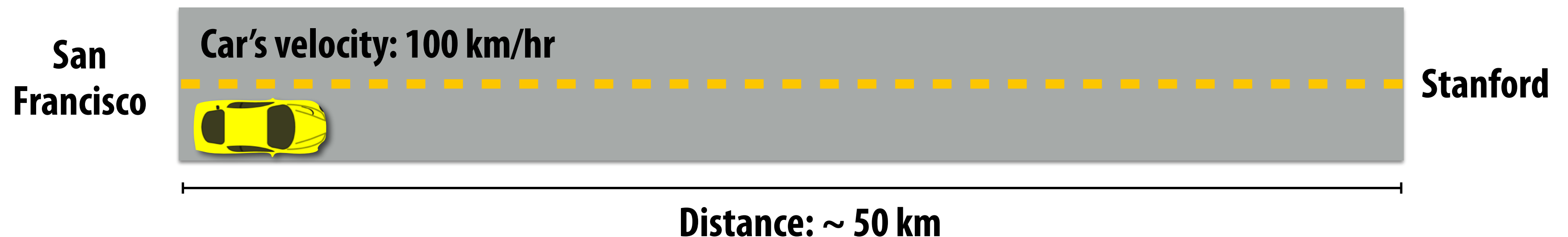


Everyone wants to get to back to the South Bay!

(Latency vs. throughput review)

Assume only one car in a lane of the highway at once.

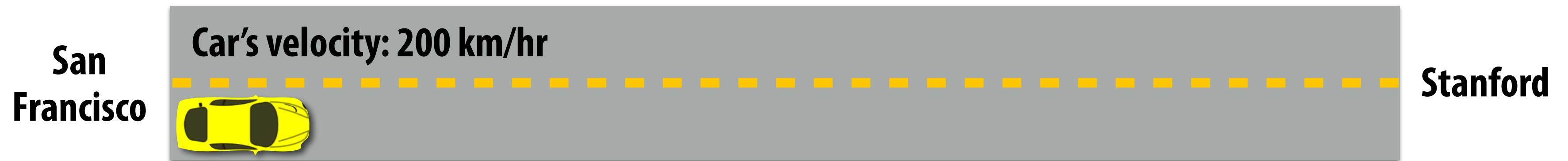
When car on highway reaches Stanford, the next car leaves San Francisco.



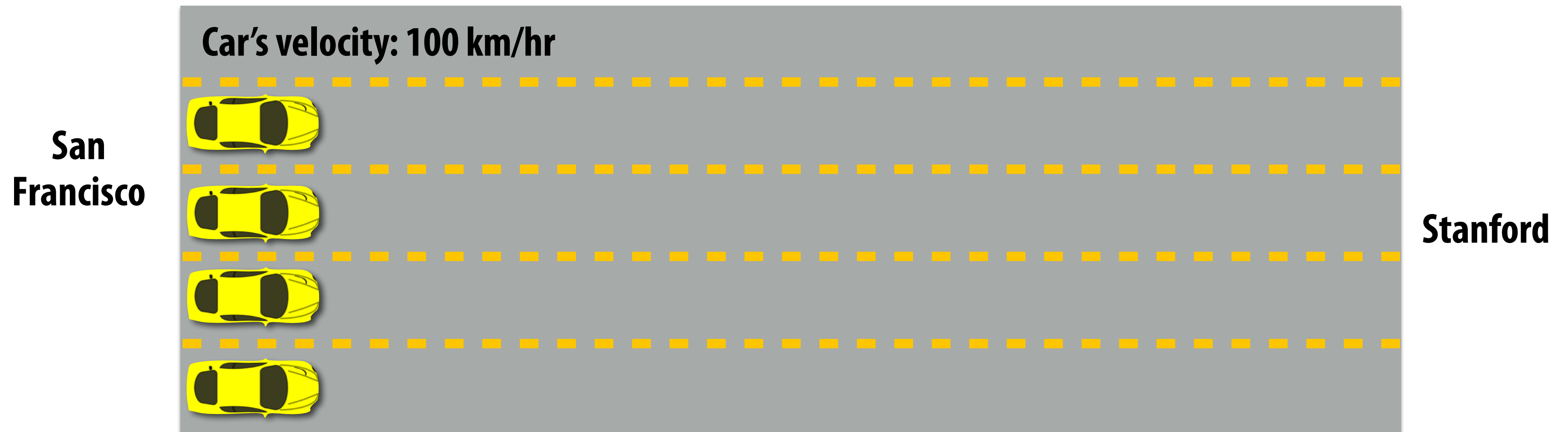
Latency of driving from San Francisco to Stanford: 0.5 hours

Throughput: 2 cars per hour

Improving throughput

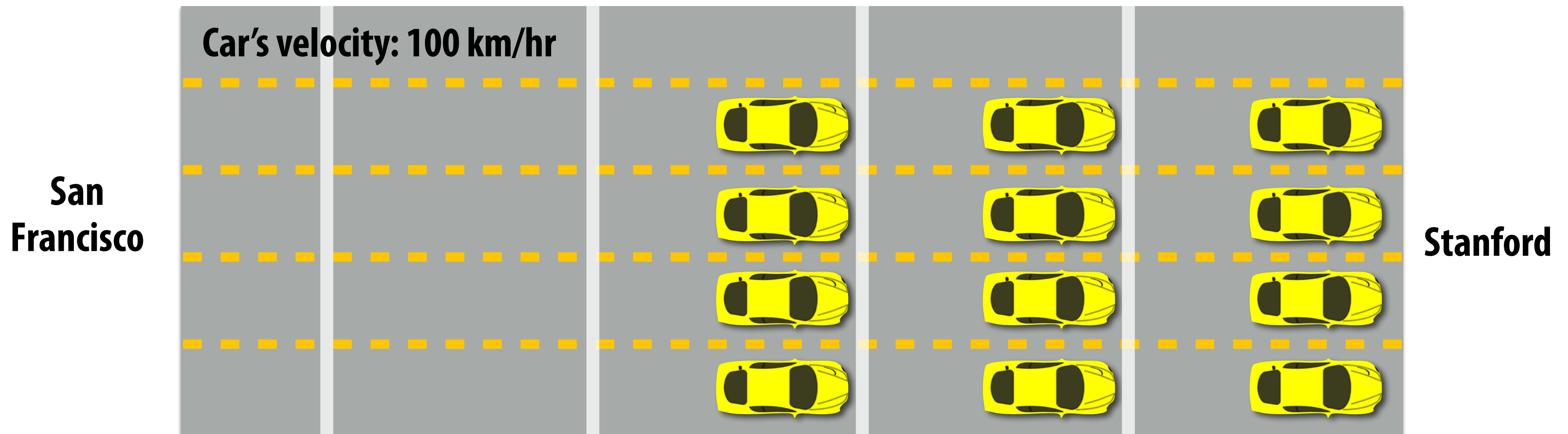


Approach 1: drive faster!
Throughput = 4 cars per hour



Approach 2: build more lanes!
Throughput = 8 cars per hour (2 cars per hour per lane)

Now caravan of 12 cars allowed on highway at once!



Latency of first car to arrive at Stanford = 0.5 hrs

Four cars from caravan arrive every 1/100th of an hour

Total time for 12-car caravan: 0.5 hrs + 3/100th of an hour

"Effective car throughput" = 12 cars / 0.53 hours = ~ 23 cars / hour

Cars spaced out by 1 km

Processor 1

Processor 2

Message size = N bits

Link can communicate B (bits/clock)

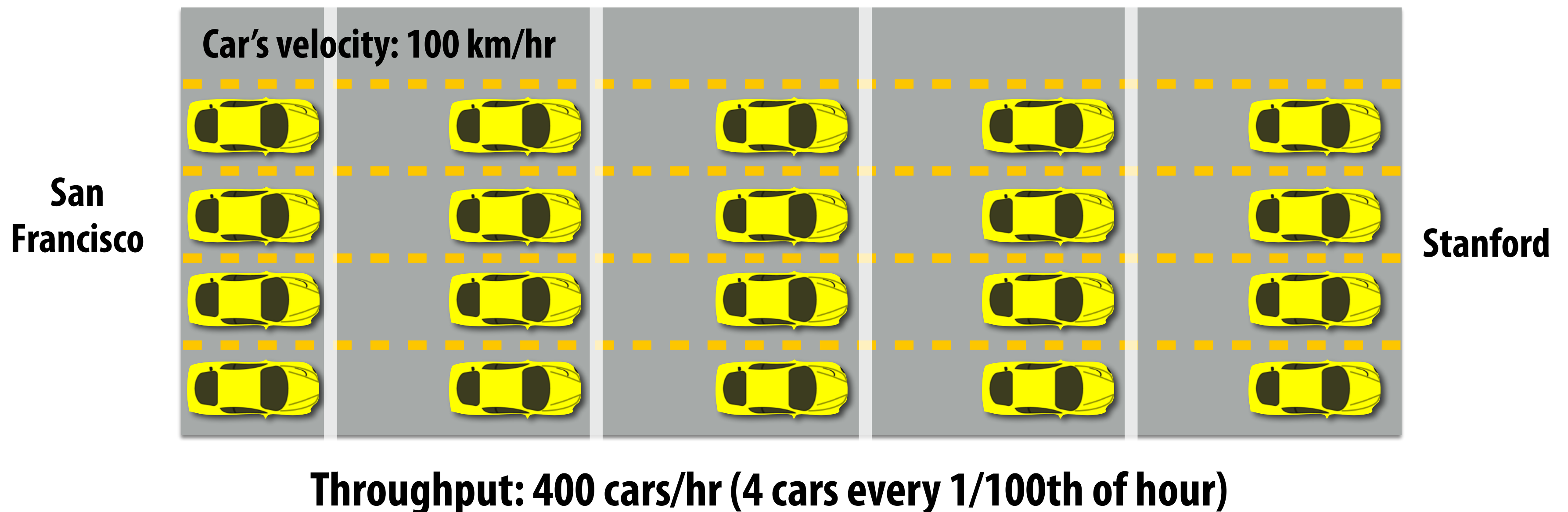
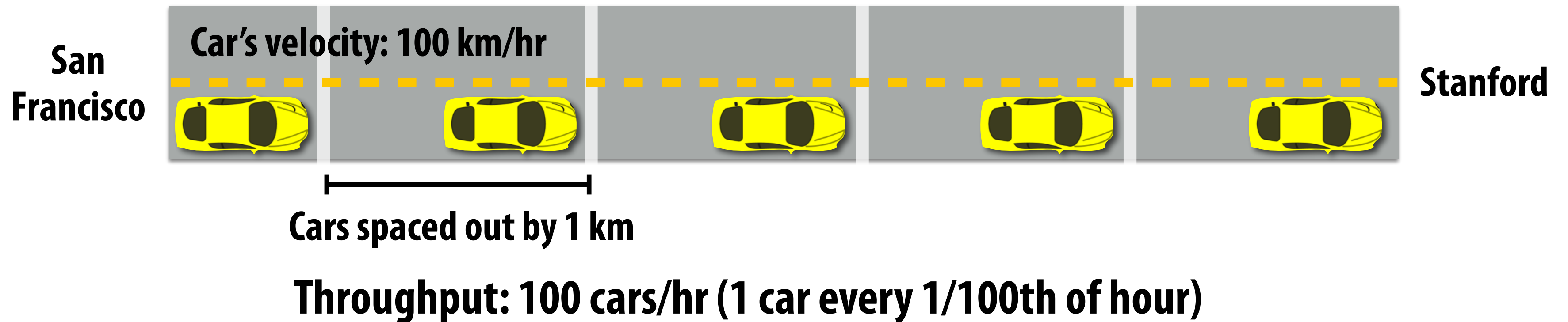
Latency of first bit to travel to destination = T clocks

Time to send message = $T + N/B$ clocks

If message must arrive before the link can send the next message, then

"Effective bandwidth" = N bits / ($T + N/B$) clocks

Using the highway more efficiently



Pipelining

Example: doing your laundry

Operation: do your laundry

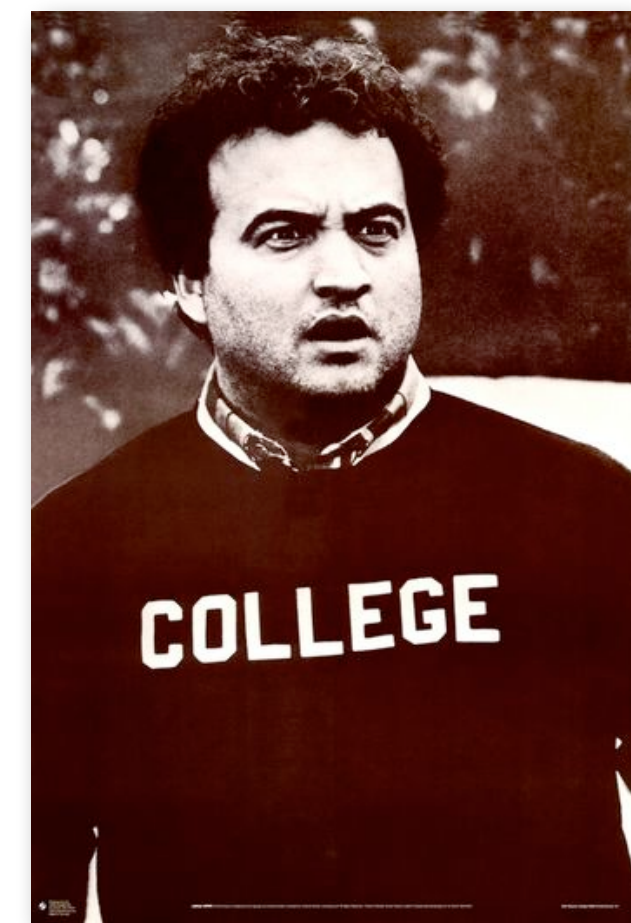
1. Wash clothes
2. Dry clothes
3. Fold clothes



Washer
45 min



Dryer
60 min



College Student
15 min

Latency of completing 1 load of laundry = 2 hours

Increasing laundry throughput

Goal: maximize throughput of many loads of laundry

**One approach: duplicate execution resources:
use two washers, two dryers, and call a friend**



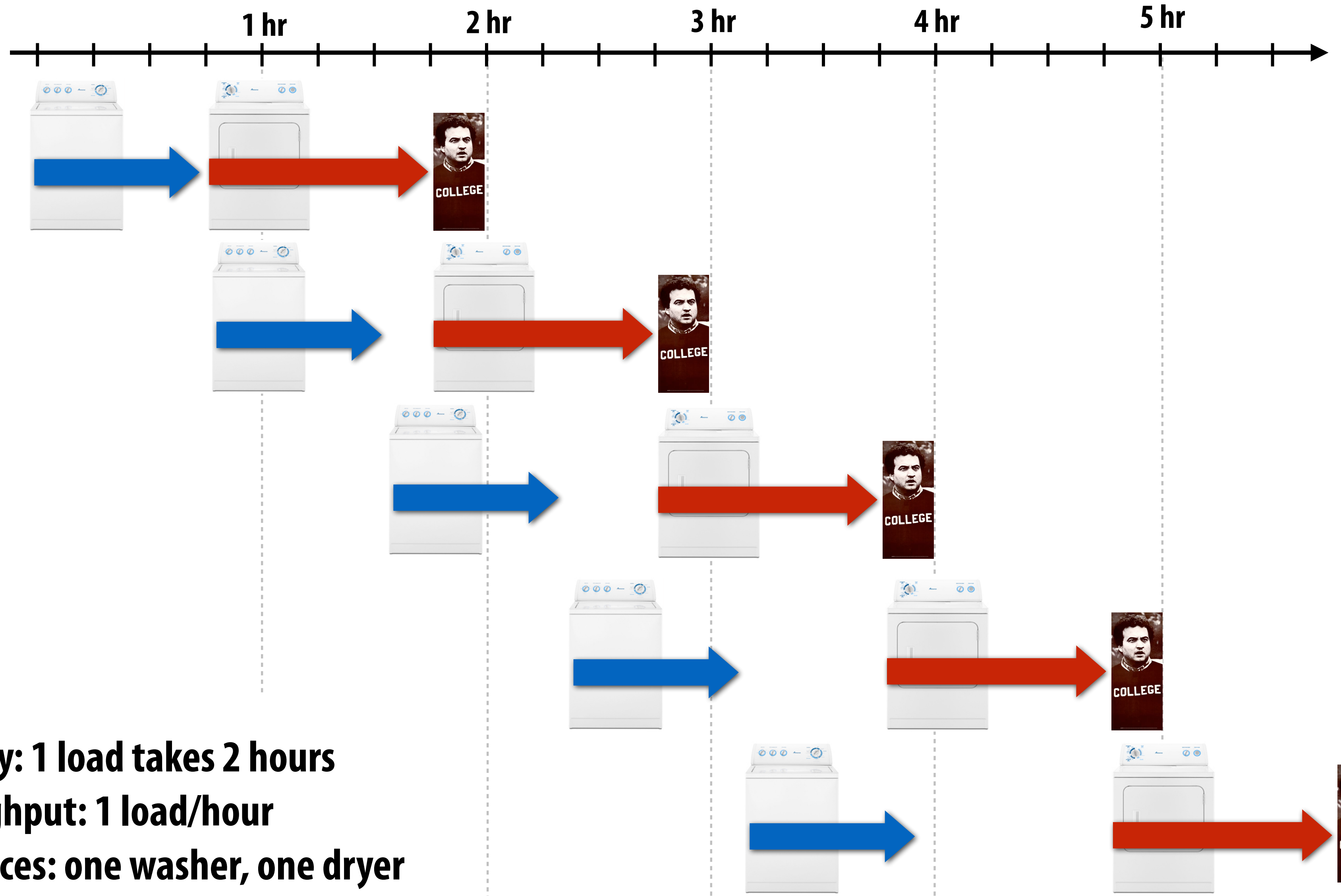
Latency of completing 2 loads of laundry = 2 hours

Throughput increases by 2x: 1 load/hour

Number of resources increased by 2x: two washers, two dryers

Pipelining

Goal: maximize throughput of doing many loads of laundry



Latency: 1 load takes 2 hours

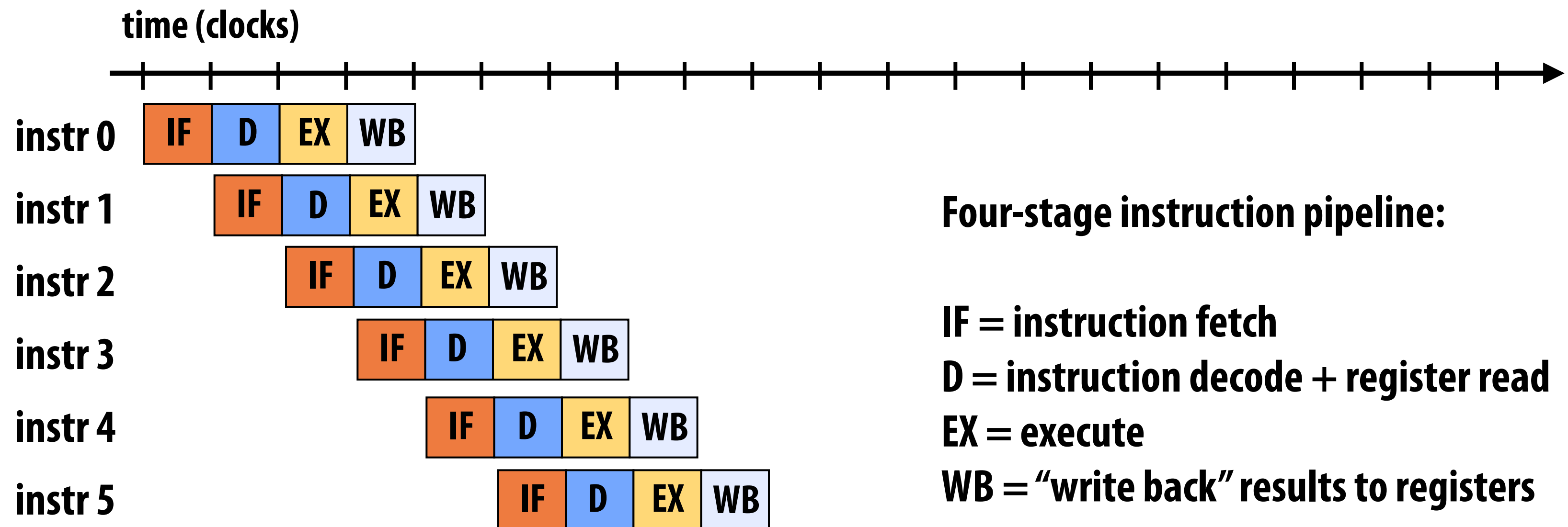
Throughput: 1 load/hour

Resources: one washer, one dryer

Example: an instruction pipeline

Break execution of each instruction down into several smaller steps

Enables higher clock frequency (only a simple, short operation is done by each part of pipeline each clock)



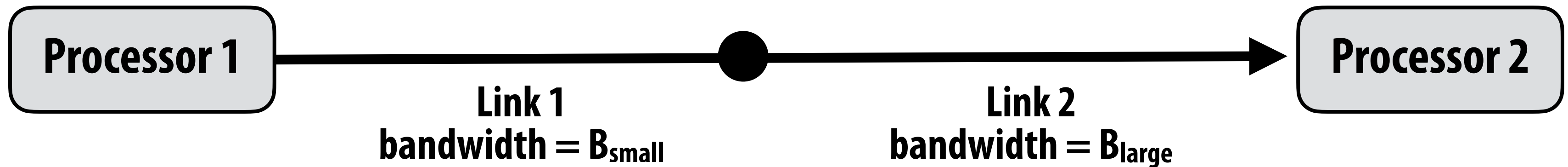
Latency: 1 instruction takes 4 cycles

Throughput: 1 instruction per cycle

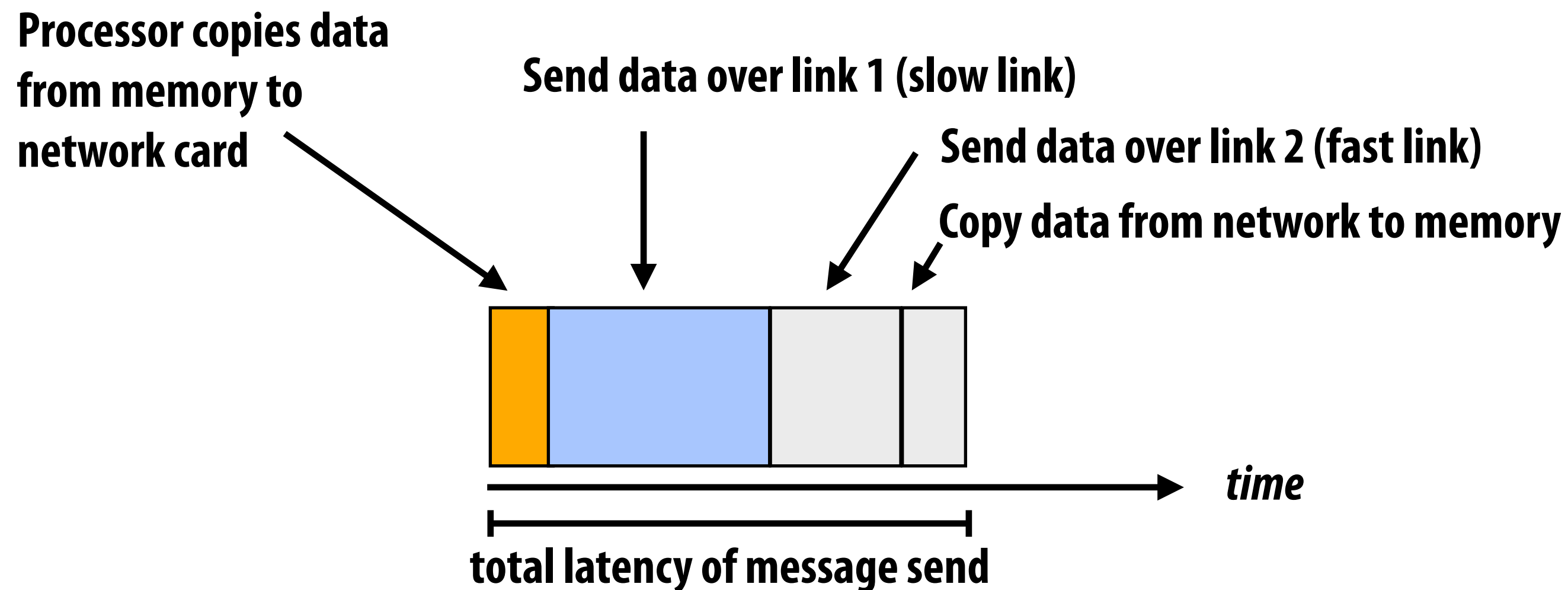
(Yes, care must be taken to ensure program correctness when back-to-back instructions are dependent.)

Intel Core i7 pipeline is variable length (it depends on the instruction) ~20 stages

A general model of msg communication



Steps in sending a N-bit message



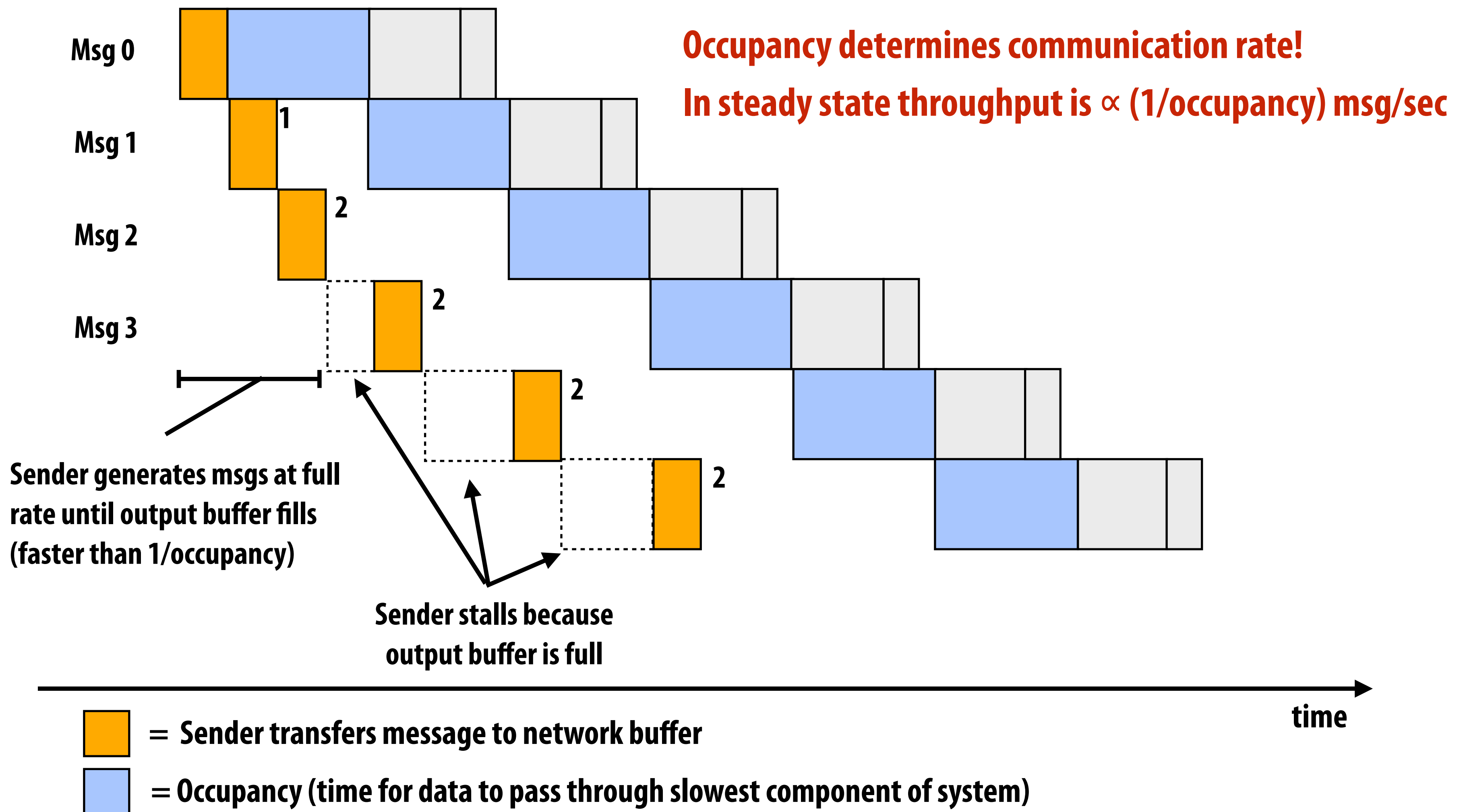
 = Occupancy (time for data to pass through slowest component of system)

Pipelined communication of messages

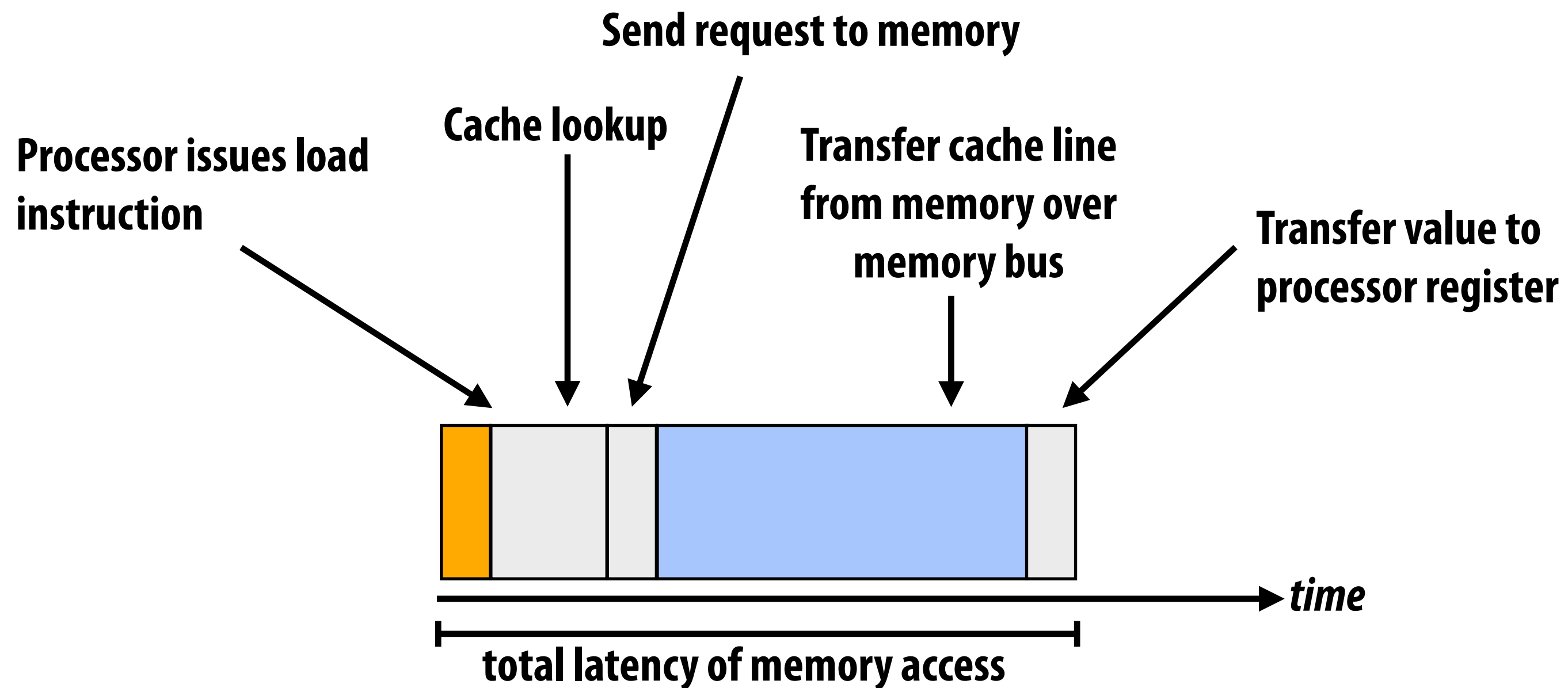
Assume processor copies messages into buffer for network to transmit from

Buffer is fixed size and can hold at most two messages

(numbers indicate number of msgs in buffer after insert)



Another example: CPU to memory communication

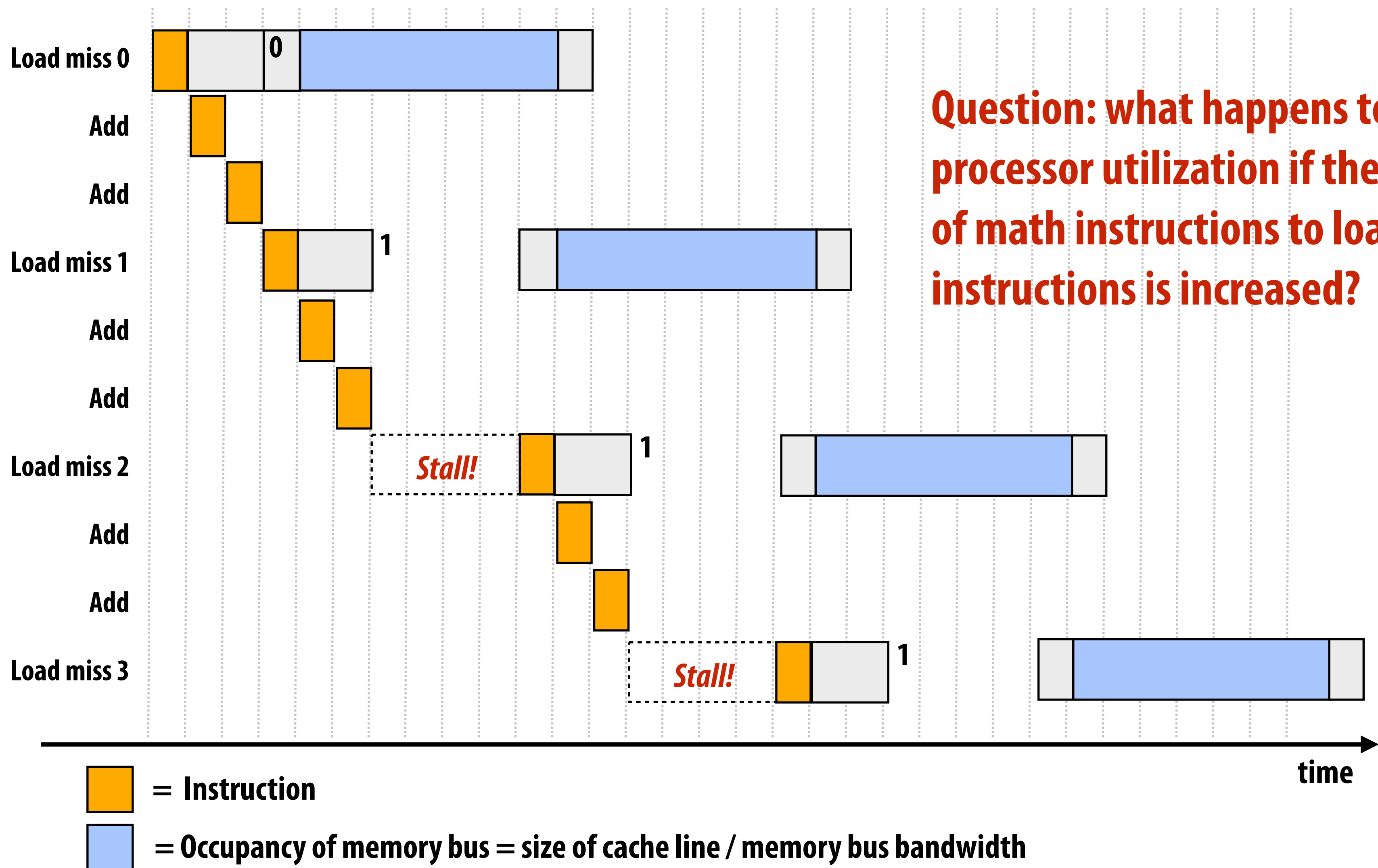


 = Time to send cache line over memory bus

Example: memory bound instruction stream!

Cache miss results in request added to size=1 memory request buffer

Notice: memory fully utilized, but processor's utilization drops to rate determined by BW



Question: what happens to processor utilization if the ratio of math instructions to load instructions is increased?

Good questions about the previous slide

- **How do you tell from the figure that the memory bus is fully utilized?**
- **How would you illustrate higher memory latency (keep in mind memory requests are pipelined and memory bus bandwidth is not changed)?**
- **How would the figure change if memory bus bandwidth was increased?**
- **Would there still be processor stalls if the ratio of math instructions to load instructions was significantly increased? Why?**

Communication-to-computation ratio

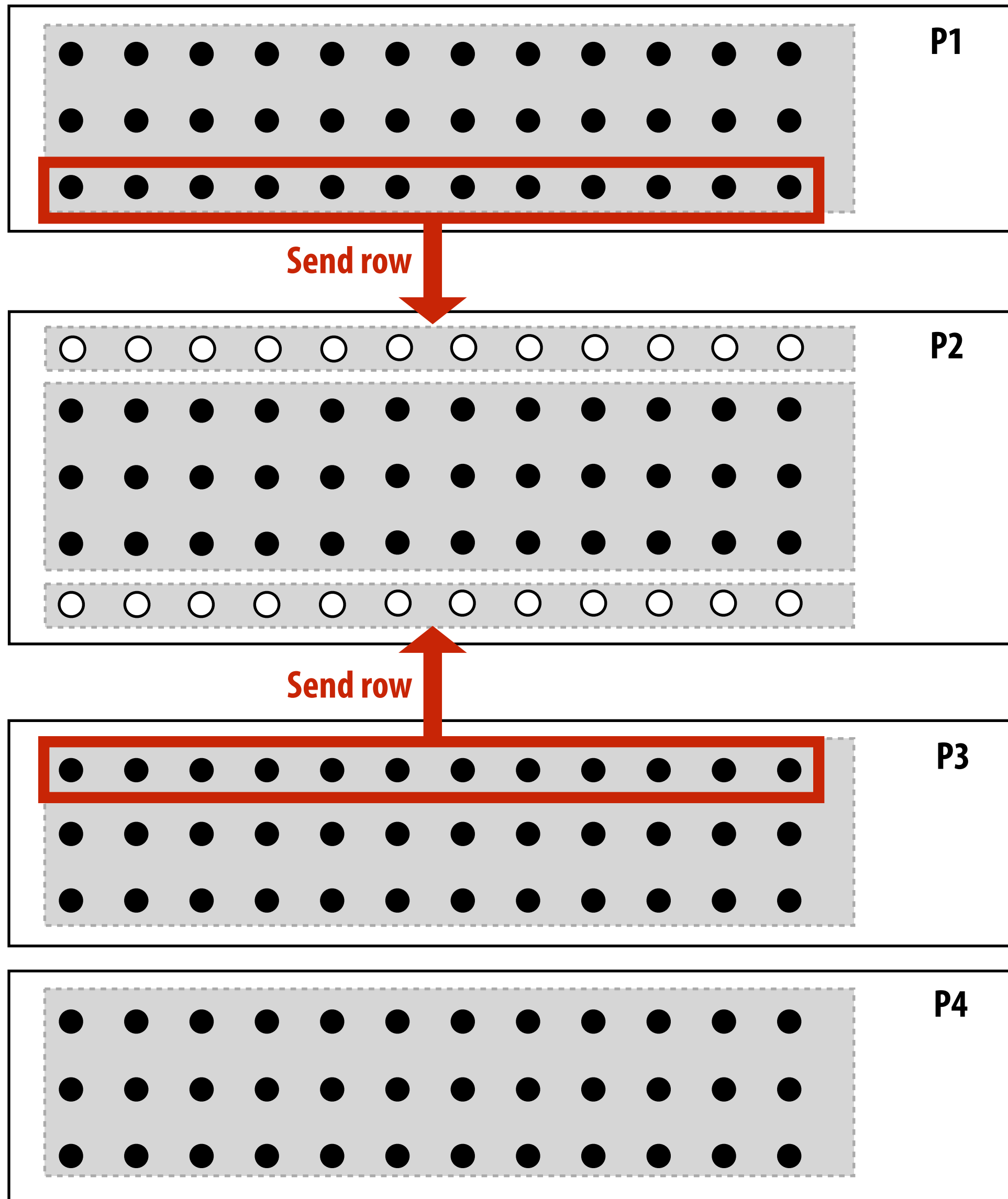
amount of communication (e.g., bytes)

amount of computation (e.g., instructions)

- If denominator is the execution time of computation, ratio gives average bandwidth requirement of code
- **“Arithmetic intensity”** = $1 / \text{communication-to-computation ratio}$
 - I find arithmetic intensity a more intuitive quantity, since higher is better.
 - It also sounds cooler
- High arithmetic intensity (low communication-to-computation ratio) is required to efficiently utilize modern parallel processors since the ratio of compute capability to available bandwidth is high (recall element-wise vector multiply example from the end of lecture 2)

Two reasons for communication: inherent vs. artifactual communication

Inherent communication



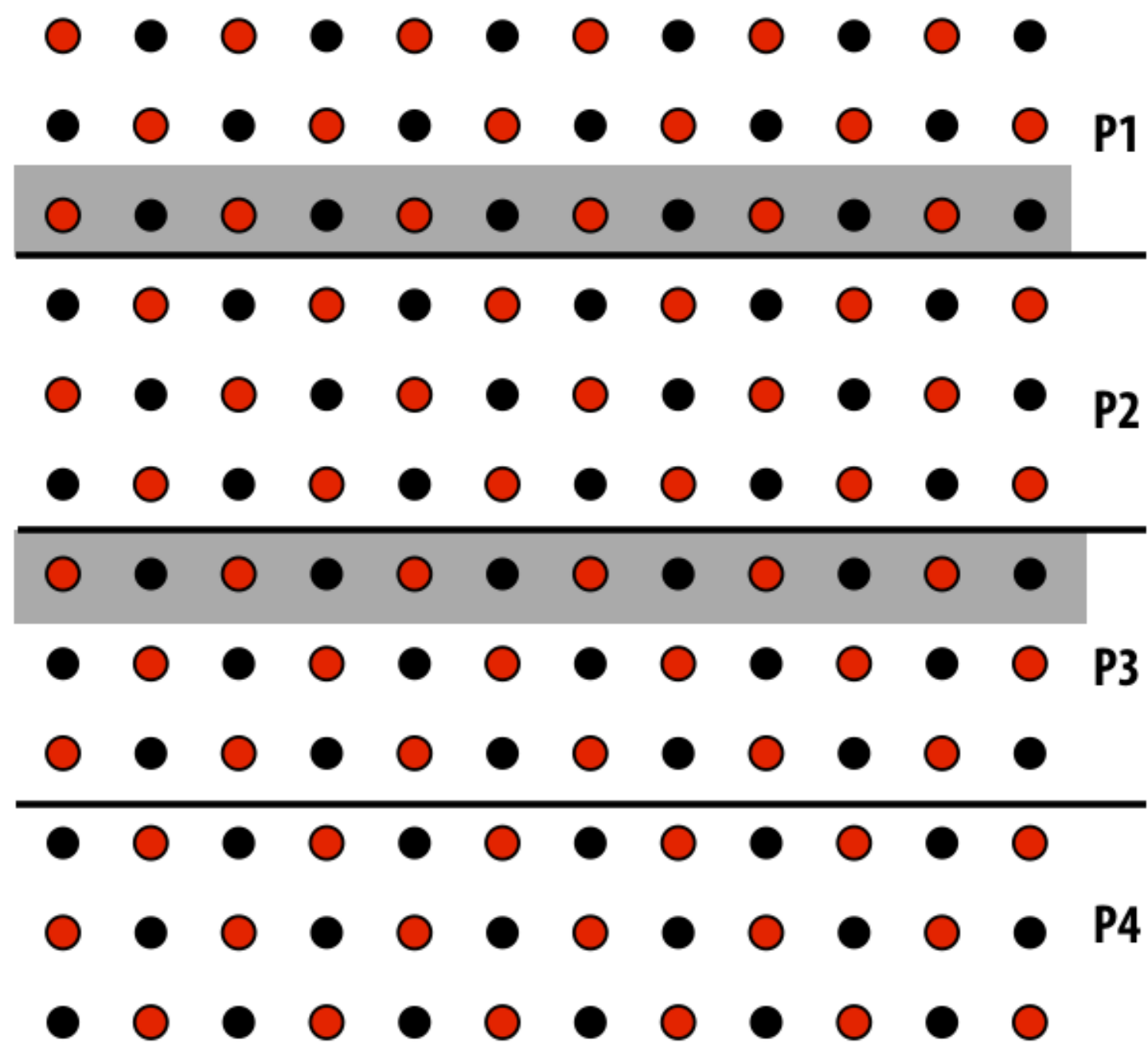
Communication that must occur in a parallel algorithm. The communication is fundamental to the algorithm.

In our messaging passing example at the start of class, sending ghost rows was inherent communication

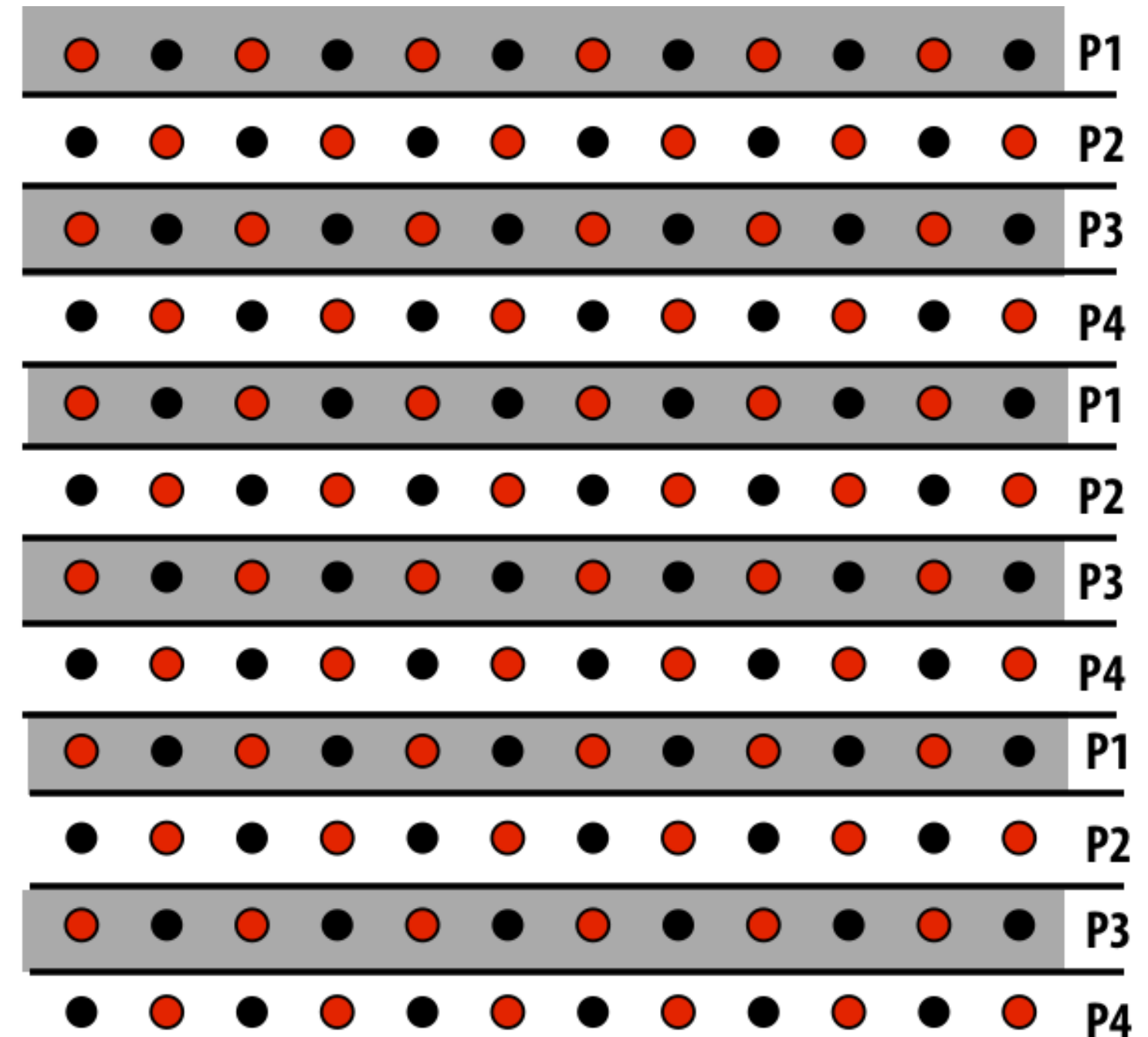
Reducing inherent communication

Good assignment decisions can reduce inherent communication
(increase arithmetic intensity)

1D blocked assignment: N x N grid



1D interleaved assignment: N x N grid

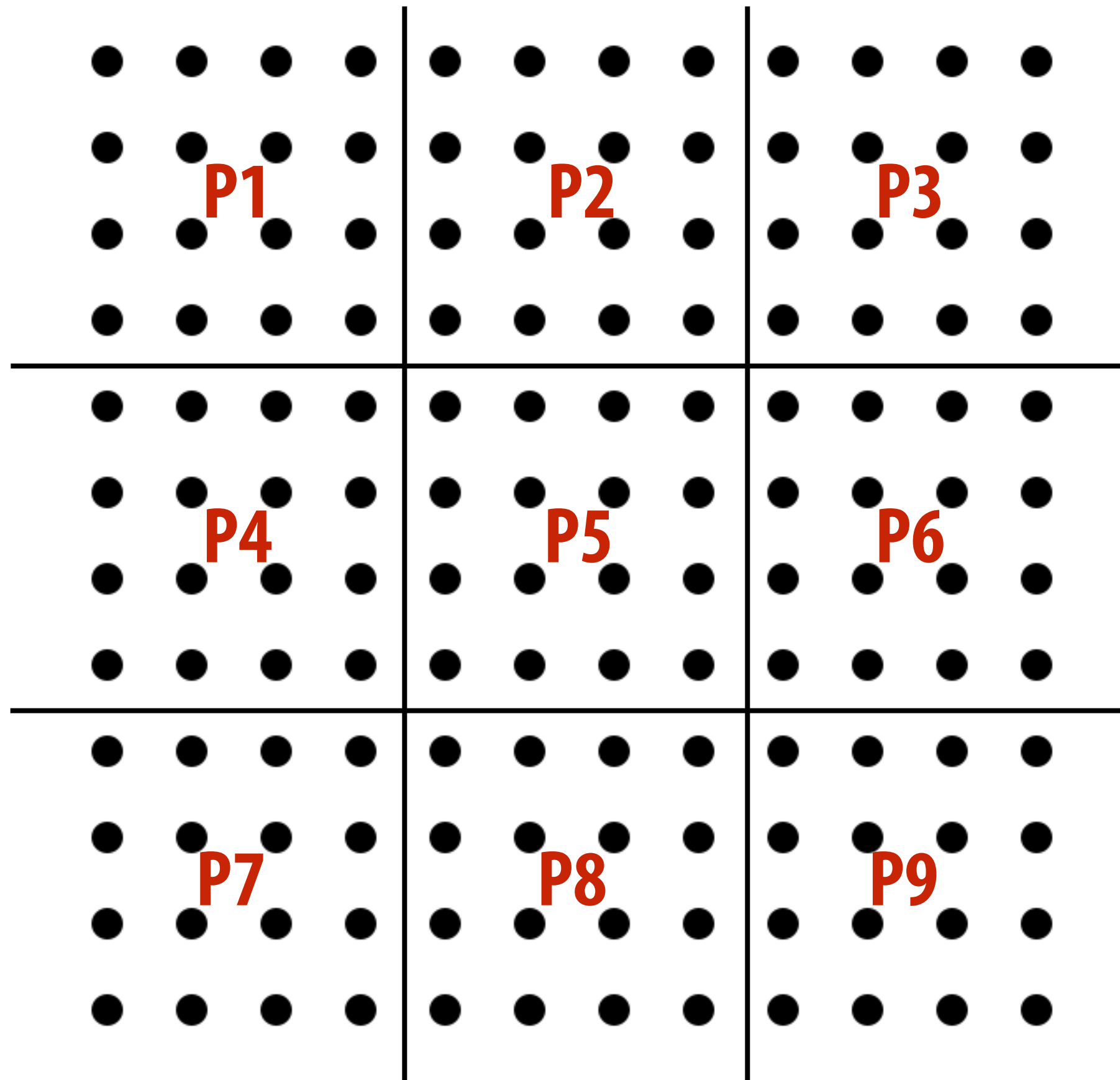


$$\frac{\text{elements computed (per processor)} \approx N^2/P}{\text{elements communicated (per processor)} \approx 2N} \propto N/P$$

$$\frac{\text{elements computed}}{\text{elements communicated}} = 1/2$$

Reducing inherent communication

2D blocked assignment: $N \times N$ grid



N^2 elements

P processors

elements computed:
(per processor)

elements communicated:
(per processor)

arithmetic intensity:

$$\frac{N^2}{P}$$

$$\propto \frac{N}{\sqrt{P}}$$

$$\frac{N}{\sqrt{P}}$$

Asymptotically better communication scaling than 1D blocked assignment

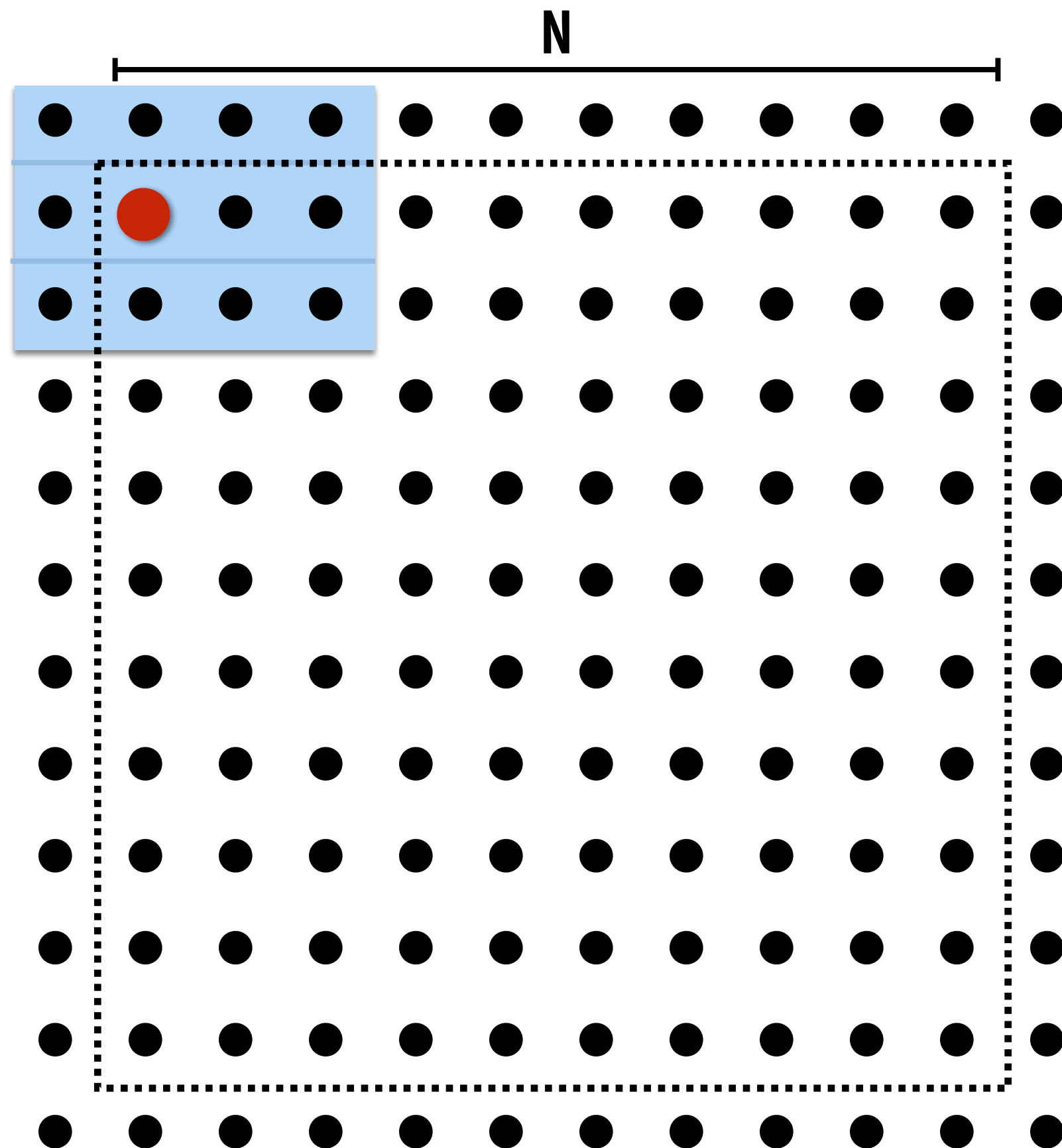
Communication costs increase sub-linearly with P

Assignment captures 2D locality of algorithm

Artifactual communication

- **Inherent communication: information that fundamentally must be moved between processors to carry out the algorithm given the specified assignment (assumes unlimited capacity caches, minimum granularity transfers, etc.)**
- **Artifactual communication: all other communication (artifactual communication results from practical details of system implementation)**

Data access in grid solver: row-major traversal



Assume row-major grid layout.

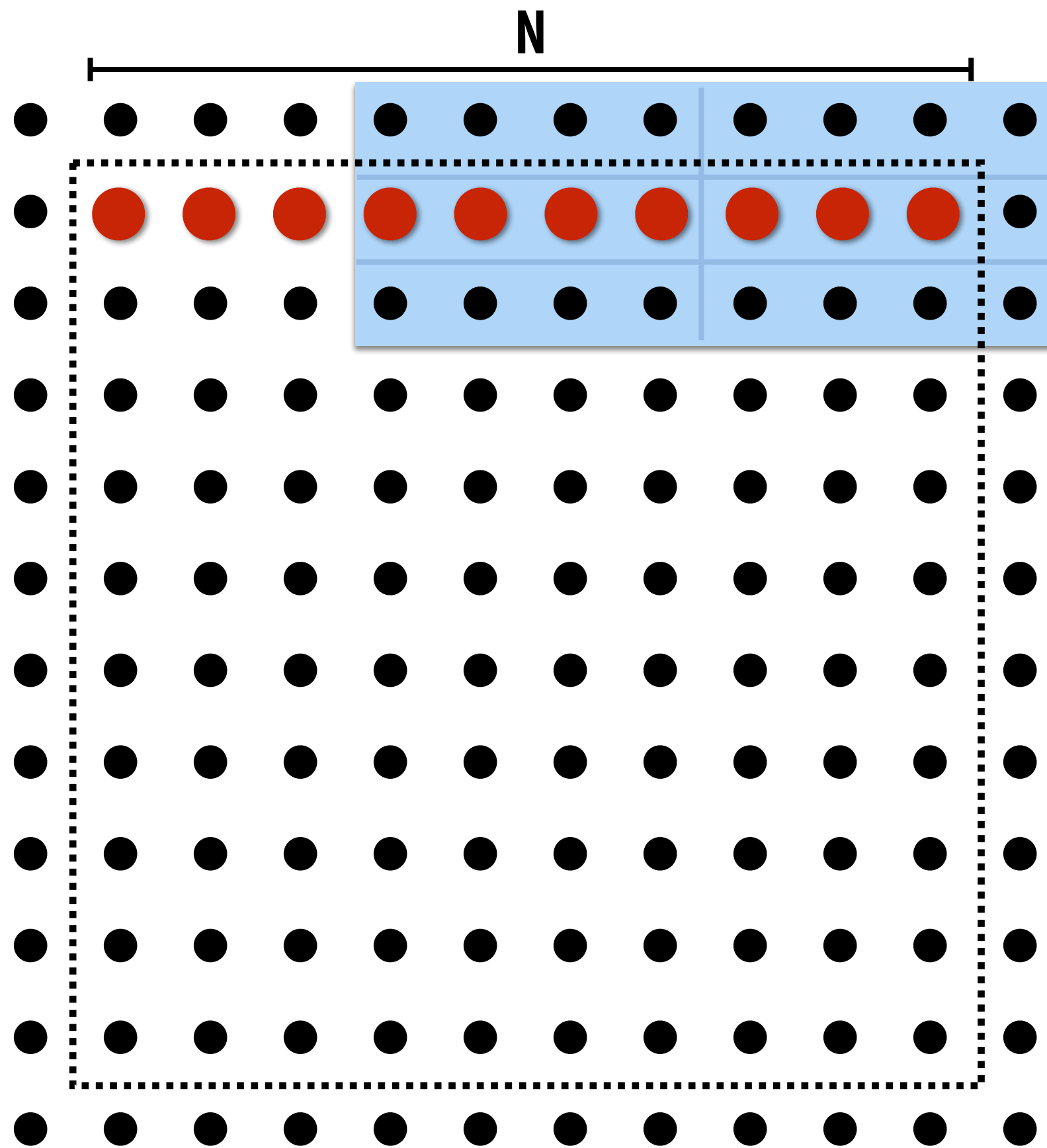
Assume cache line is 4 grid elements.

Cache capacity is 24 grid elements (6 lines)

Recall grid solver application.

Blue elements show data that is in cache after update to red element.

Data access in grid solver: row-major traversal



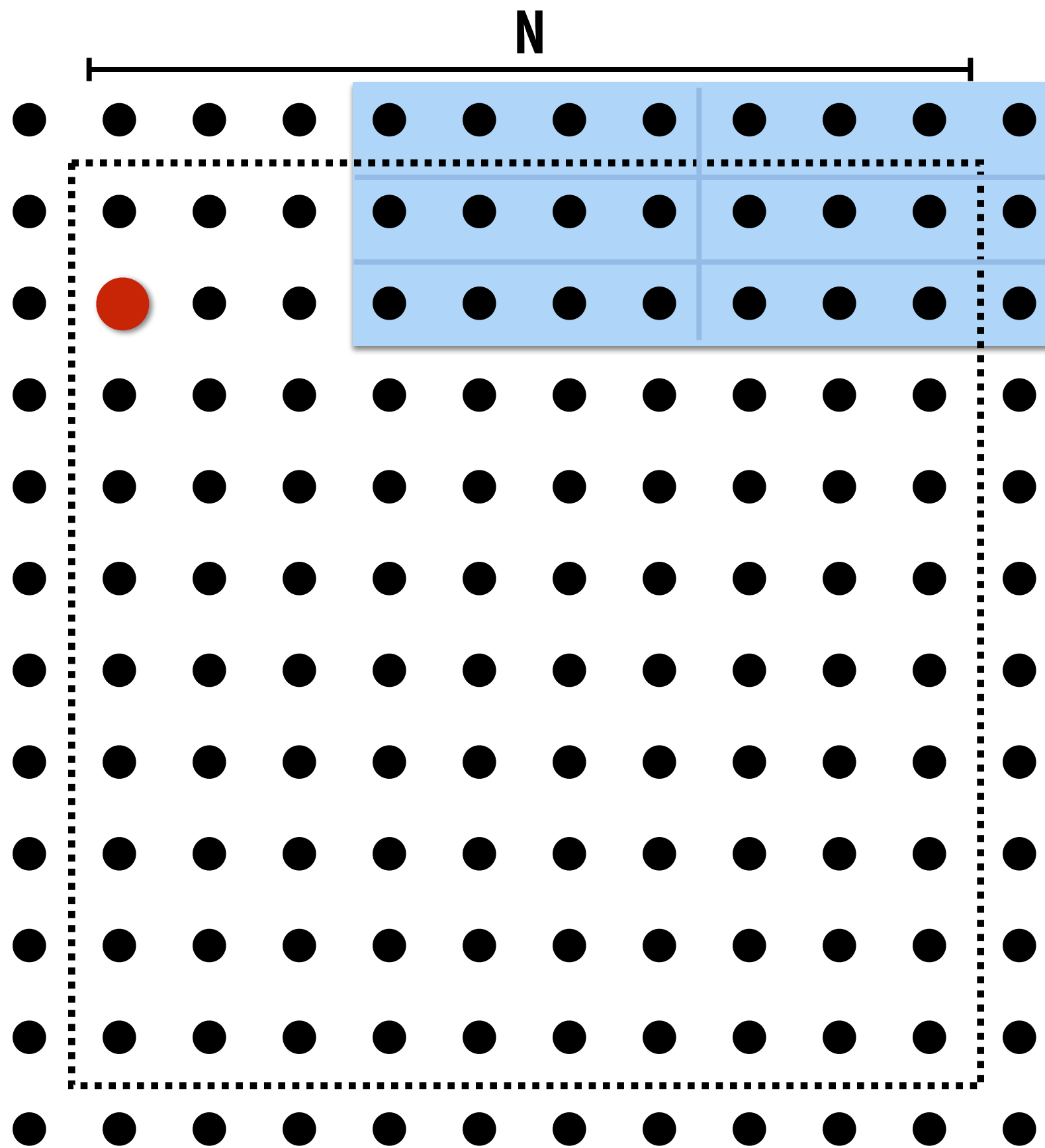
Assume row-major grid layout.

Assume cache line is 4 grid elements.

Cache capacity is 24 grid elements (6 lines)

Blue elements show data in cache at end of processing first row.

Problem with row-major traversal: long time between accesses to same data



Assume row-major grid layout.

Assume cache line is 4 grid elements.

Cache capacity is 24 grid elements (6 lines)

Although elements (0,2) and (0,1) had been accessed previously, they are no longer present in cache at start of processing row 2.

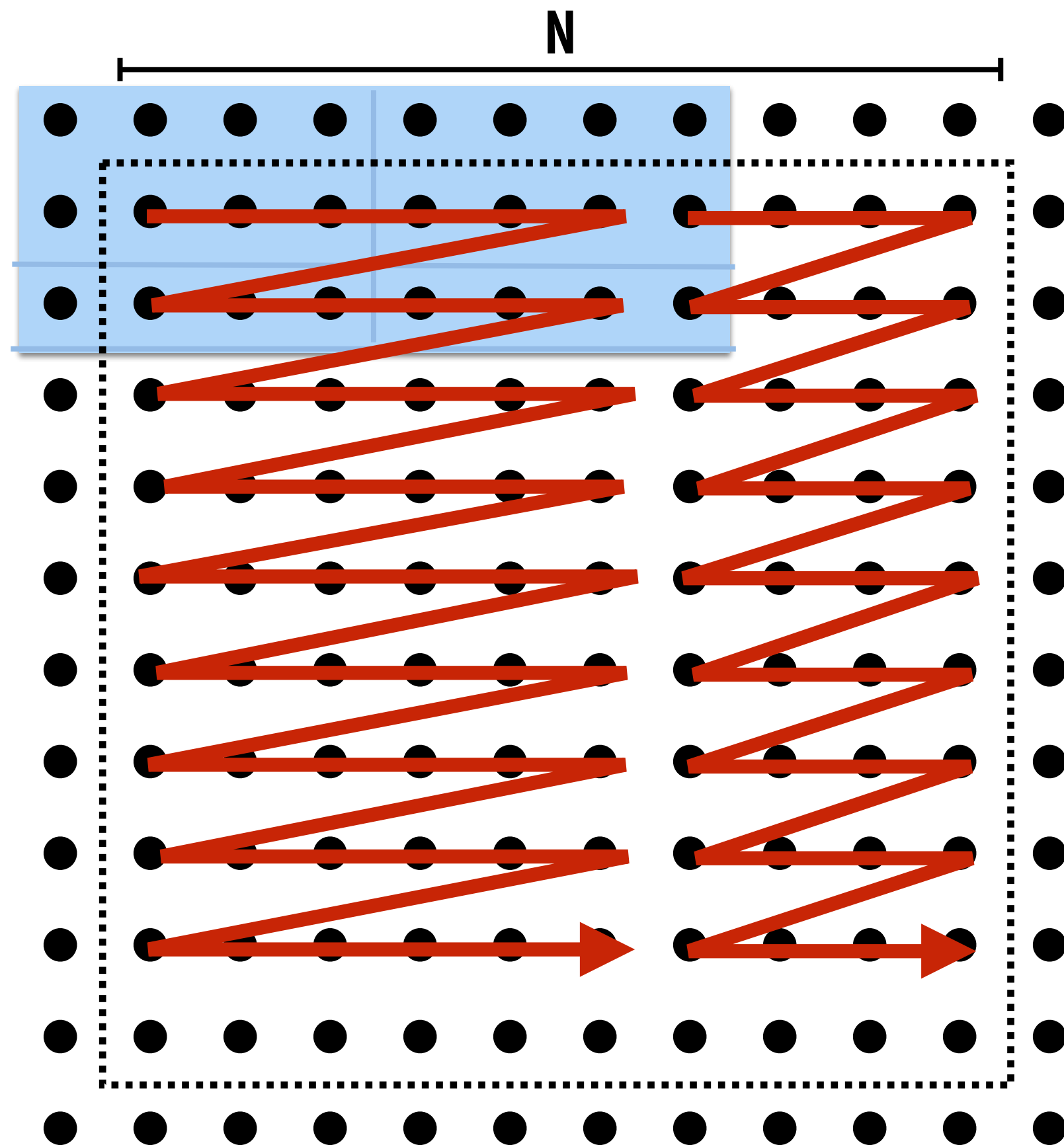
This program loads three lines for every four elements of output.

Artifactual communication examples

- **System might have a minimum granularity of data transfer (result: system must communicate more data than what is needed)**
 - **Program loads one 4-byte float value but entire 64-byte cache line must be transferred from memory (16x more communication than necessary)**
- **System operation might result in unnecessary communication:**
 - **Program stores 16 consecutive 4-byte float values, and as a result the entire 64-byte cache line is loaded from memory, entirely overwritten, then subsequently stored to memory (2x overhead... load was unnecessary)**
- **Finite replication capacity (the same data communicated to processor multiple times because cache is too small to retain it between accesses)**

Techniques for reducing communication

Improving temporal locality by changing grid traversal order



Assume row-major grid layout.

Assume cache line is 4 grid elements.

Cache capacity is 24 grid elements (6 lines)

“Blocked” iteration order

(diagram shows state of cache after finishing work from first row of first block)

Now load two cache lines for every six elements of output

Improving temporal locality by fusing loops

```
void add(int n, float* A, float* B, float* C) {  
    for (int i=0; i<n; i++)  
        C[i] = A[i] + B[i];  
}
```

Two loads, one store per math op
(arithmetic intensity = 1/3)

```
void mul(int n, float* A, float* B, float* C) {  
    for (int i=0; i<n; i++)  
        C[i] = A[i] * B[i];  
}
```

Two loads, one store per math op
(arithmetic intensity = 1/3)

```
float* A, *B, *C, *D, *E, *tmp1, *tmp2;
```

```
// assume arrays are allocated here
```

```
// compute E = D + ((A + B) * C)
```

```
add(n, A, B, tmp1);  
mul(n, tmp1, C, tmp2);  
add(n, tmp2, D, E);
```

Overall arithmetic intensity = 1/3

```
void fused(int n, float* A, float* B, float* C, float* D, float* E) {  
    for (int i=0; i<n; i++)  
        E[i] = D[i] + (A[i] + B[i]) * C[i];  
}
```

Four loads, one store per 3 math ops
(arithmetic intensity = 3/5)

```
// compute E = D + (A + B) * C  
fused(n, A, B, C, D, E);
```

Code on top is more modular (e.g, array-based math library like numPy in Python)

Code on bottom performs much better. Why?

Improve arithmetic intensity by sharing data

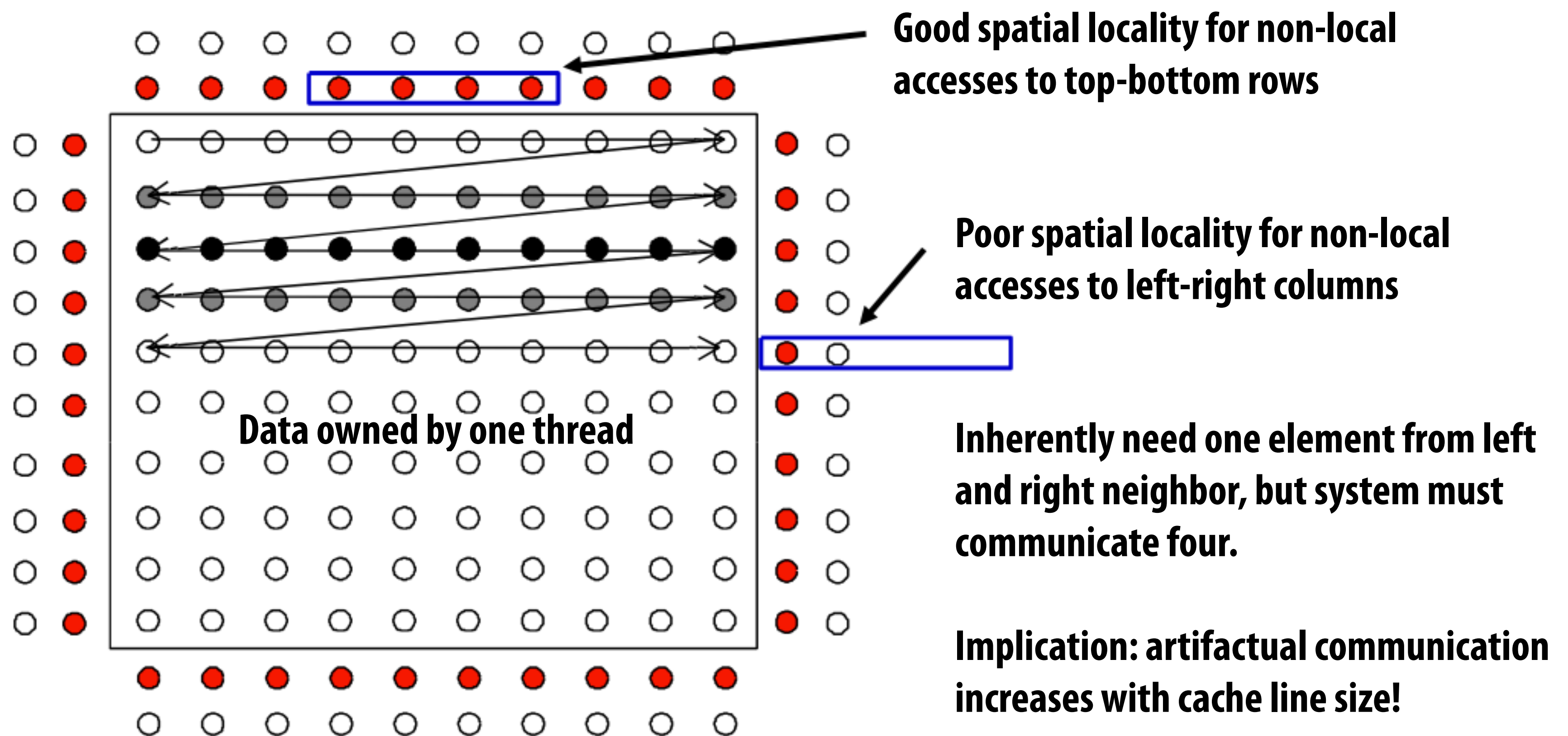
- **Exploit sharing: co-locate tasks that operate on the same data**
 - **Schedule threads working on the same data structure at the same time on the same processor**
 - **Reduces inherent communication**

Exploiting spatial locality

- **High granularity communication (e.g., a cache line) may introduce artifactual communication**
 - **If application has low spatial locality, system may transfer data that program did not need**

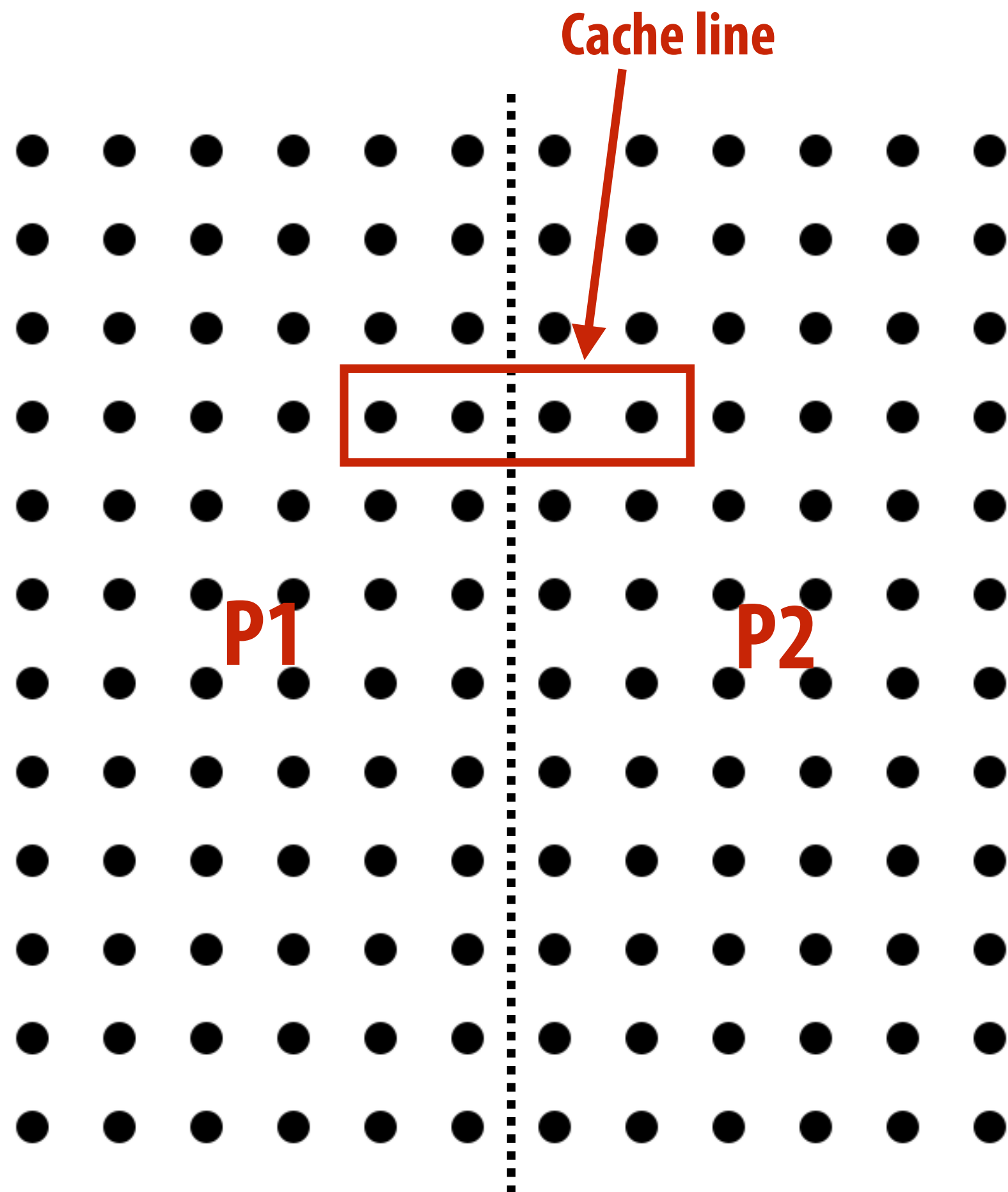
Artifactual communication due to comm. granularity

Consider 2D blocked assignment of data to processors described previously.
Assume: communication granularity is a cache line, and a cache line contains four elements



● = required elements assigned to other processors

Artifactual communication due to cache line communication granularity



Data partitioned in half by column. Partitions assigned to threads running on P1 and P2

Threads access their assigned elements (no inherent communication exists)

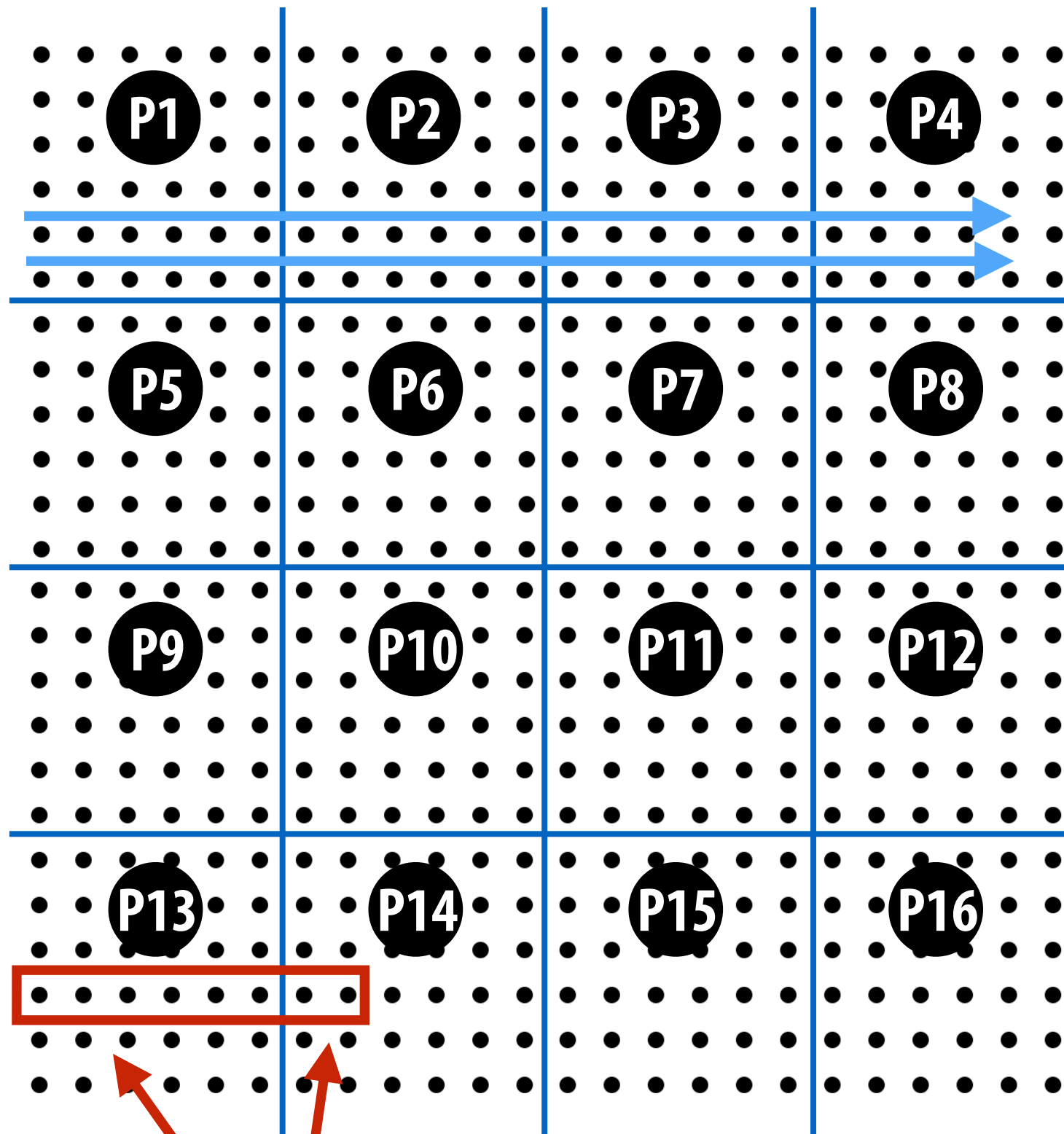
But data access on real machine triggers (artifactual) communication due to the cache line being written to by both processors *

*** further detail in the upcoming cache coherence lectures**

Reducing artifactual comm: blocked data layout

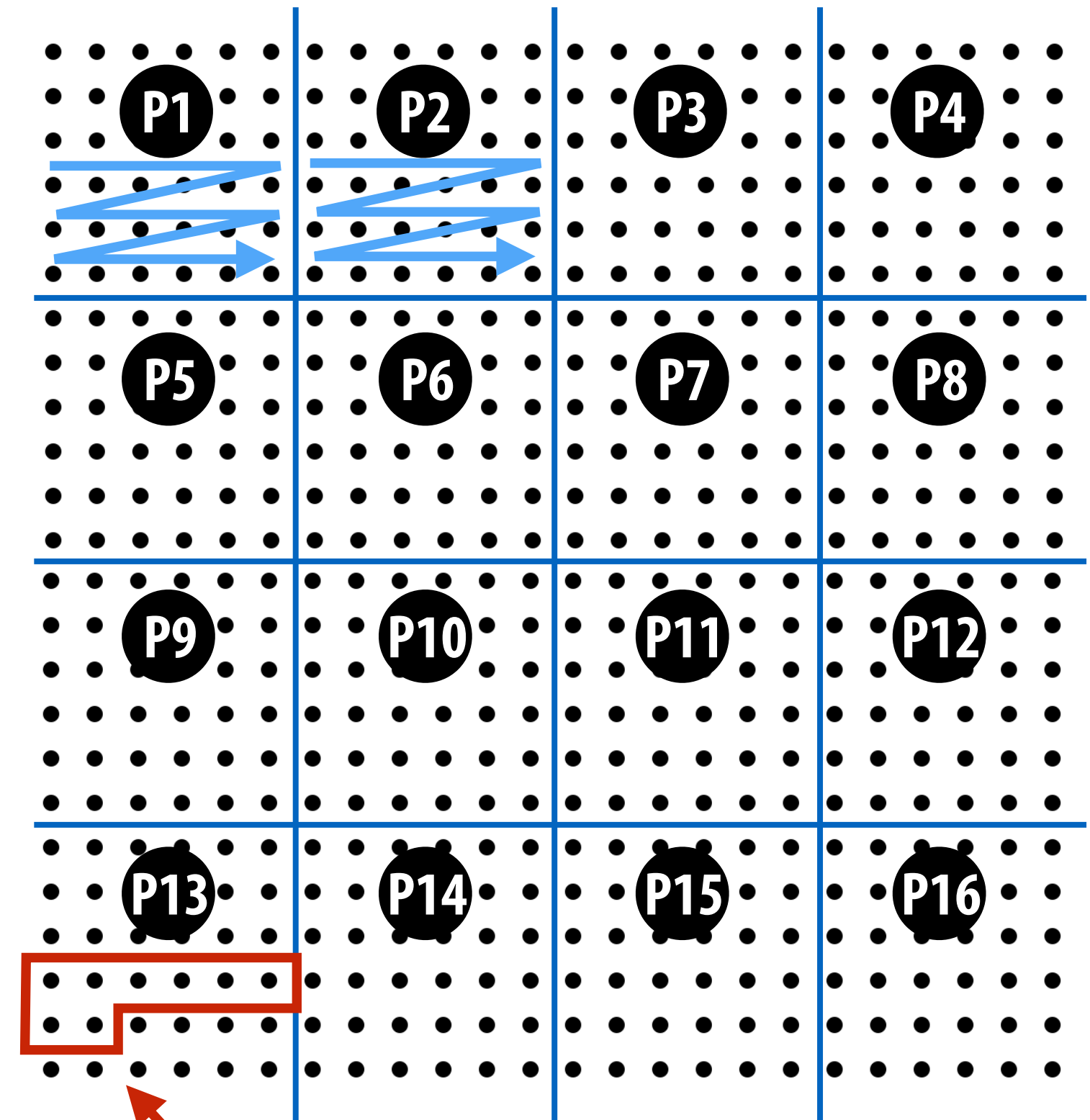
(Blue lines indicate consecutive memory addresses)

2D, row-major array layout



Consecutive addresses
straddle partition boundary

4D array layout (block-major)

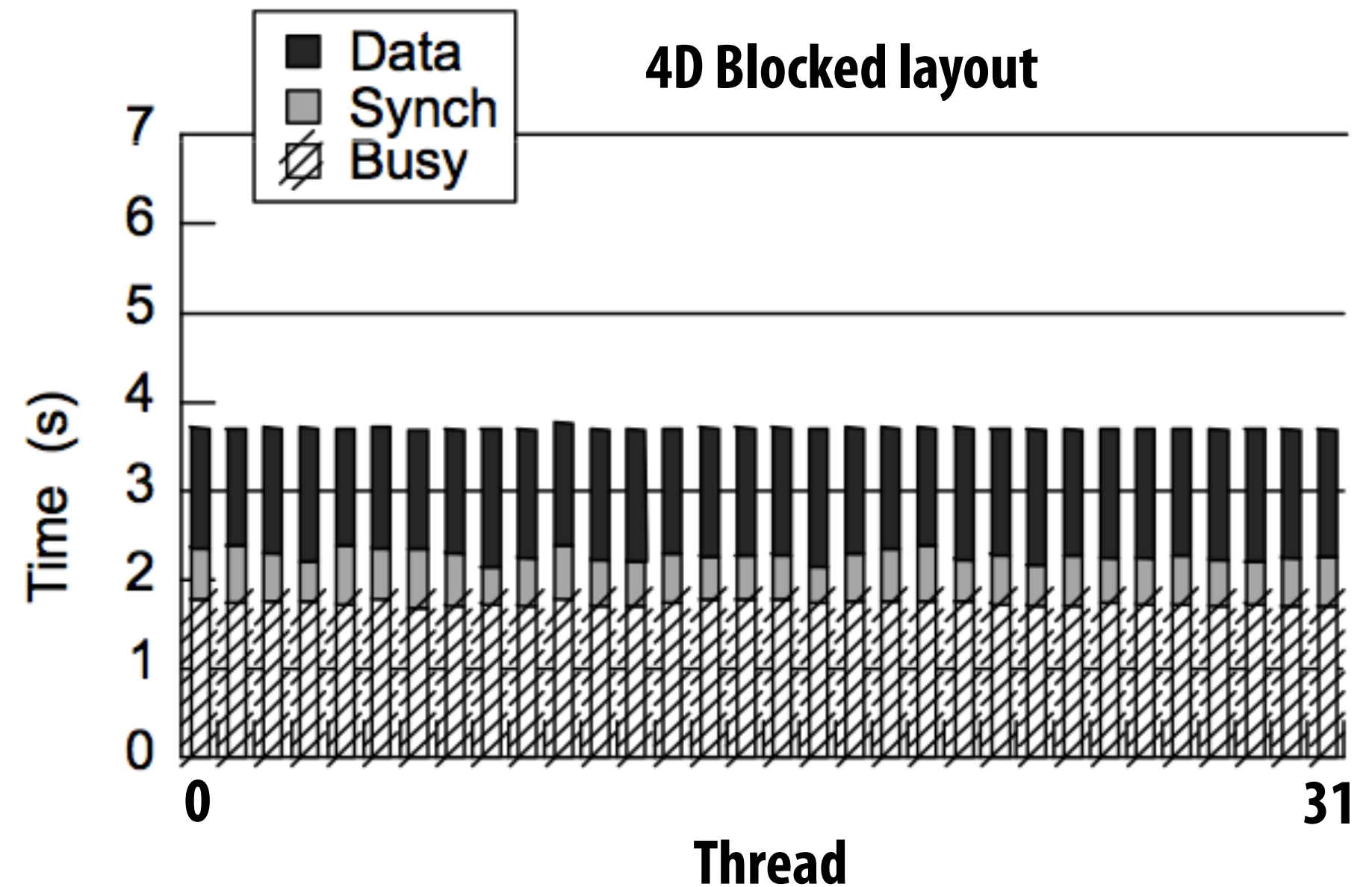
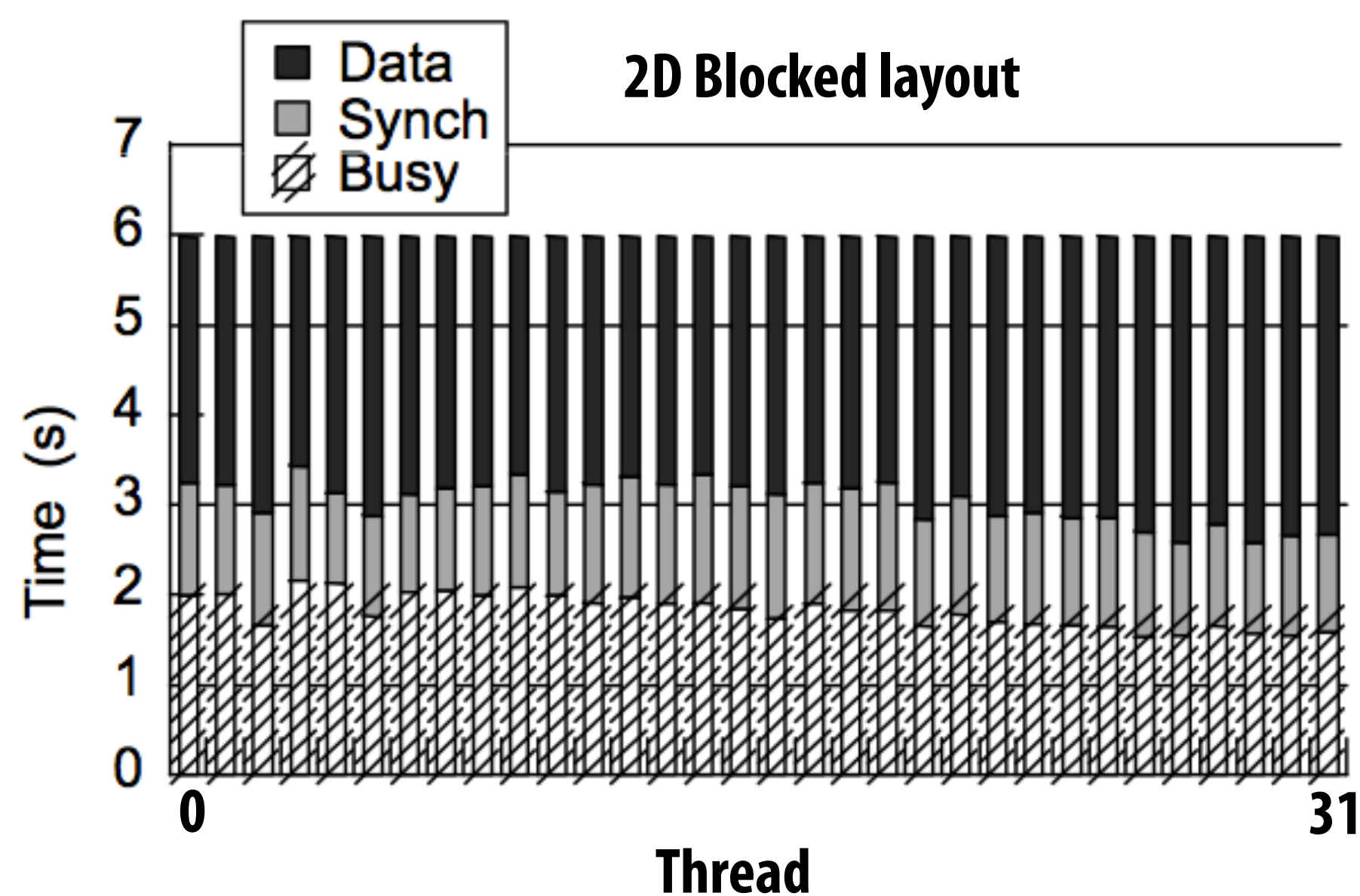


Consecutive addresses remain
within single partition

**Note: don't confuse blocked assignment of work to threads (true in both cases above)
with blocked data layout of data elements in the address space (only at right)**

Grid solver: execution time breakdown

Execution on 32-processor SGI Origin 2000 (1026 x 1026 grids)



■ Observations:

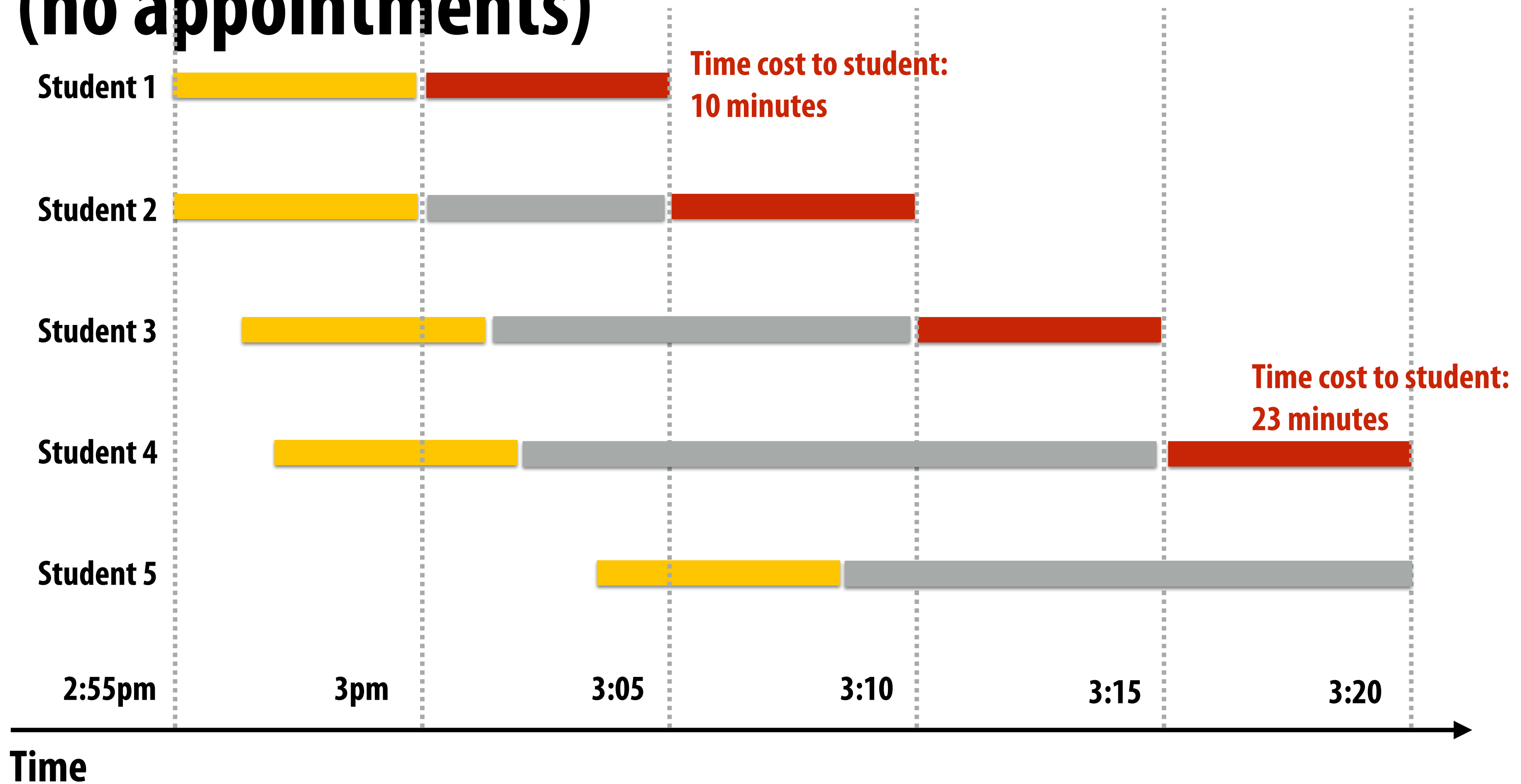
- Static assignment is sufficient (approximately equal busy time per thread)
- 4D blocking of grid reduces time spent on communication (reflected on graph as data wait time)
- Synchronization cost is largely due to waiting at barriers

Contention

Example: office hours from 3-3:20pm (no appointments)

- **Operation to perform: Professor Kayvon helps a student with a question**
- **Execution resource: Professor Kayvon**
- **Steps in operation:**
 1. **Student walks from Bytes Cafe to Kayvon's office (5 minutes)**
 2. **Student waits in line (if necessary)**
 3. **Student gets question answered with insightful answer (5 minutes)**

Example: office hours from 3-3:20pm (no appointments)



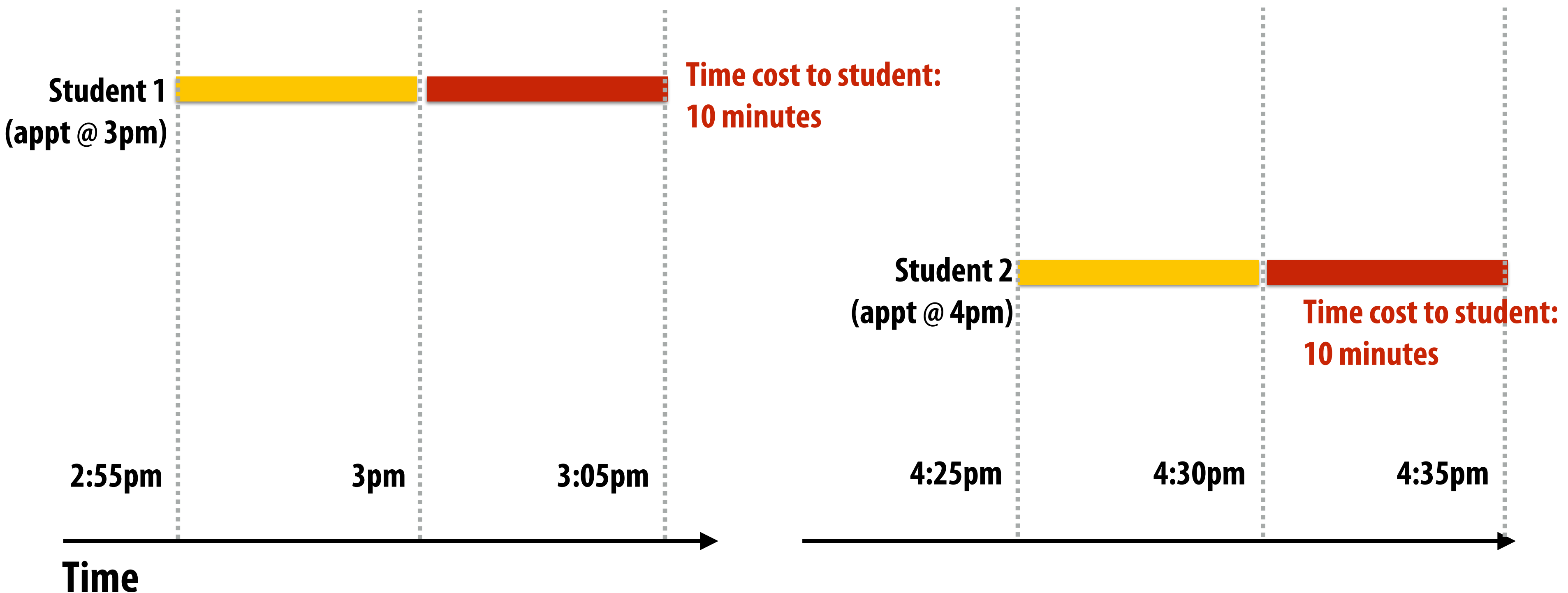
■ = Walk to Kayvon's office (5 minutes)

■ = Wait in line

■ = Get question answered

Problem: contention for shared resource results in longer overall operation times (and likely higher cost to students)

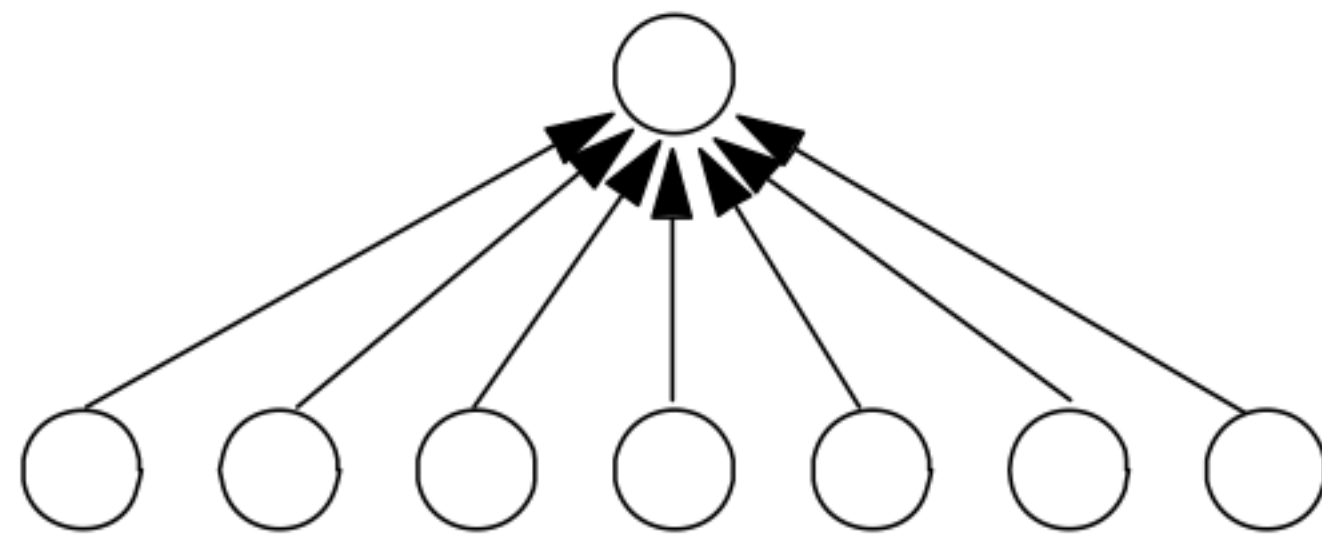
Example: two students make appointments to talk to me about course material (at 3pm and at 4:30pm)



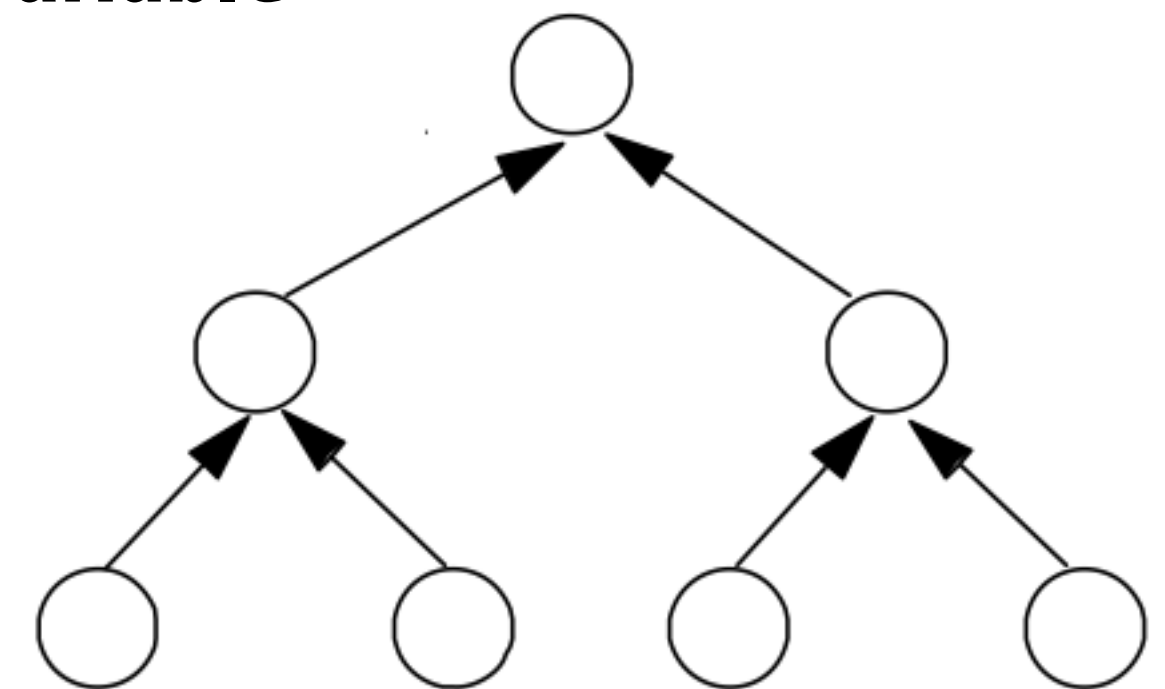
Contention

- **A resource can perform operations at a given throughput (number of transactions per unit time)**
 - **Memory, communication links, servers, TA's at office hours, etc.**
- **Contention occurs when many requests to a resource are made within a small window of time (the resource is a "hot spot")**

Example: updating a shared variable



Flat communication:
potential for high contention
(but low latency if no contention)

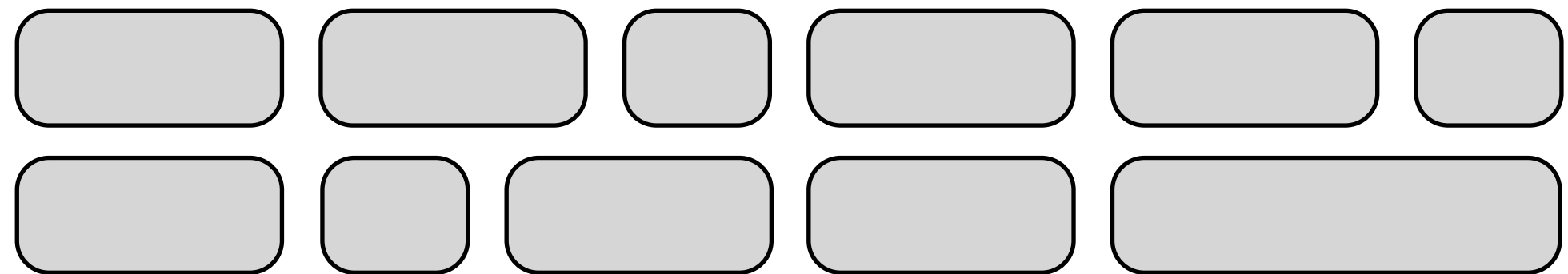


Tree structured communication:
reduces contention
(but higher latency under no contention)

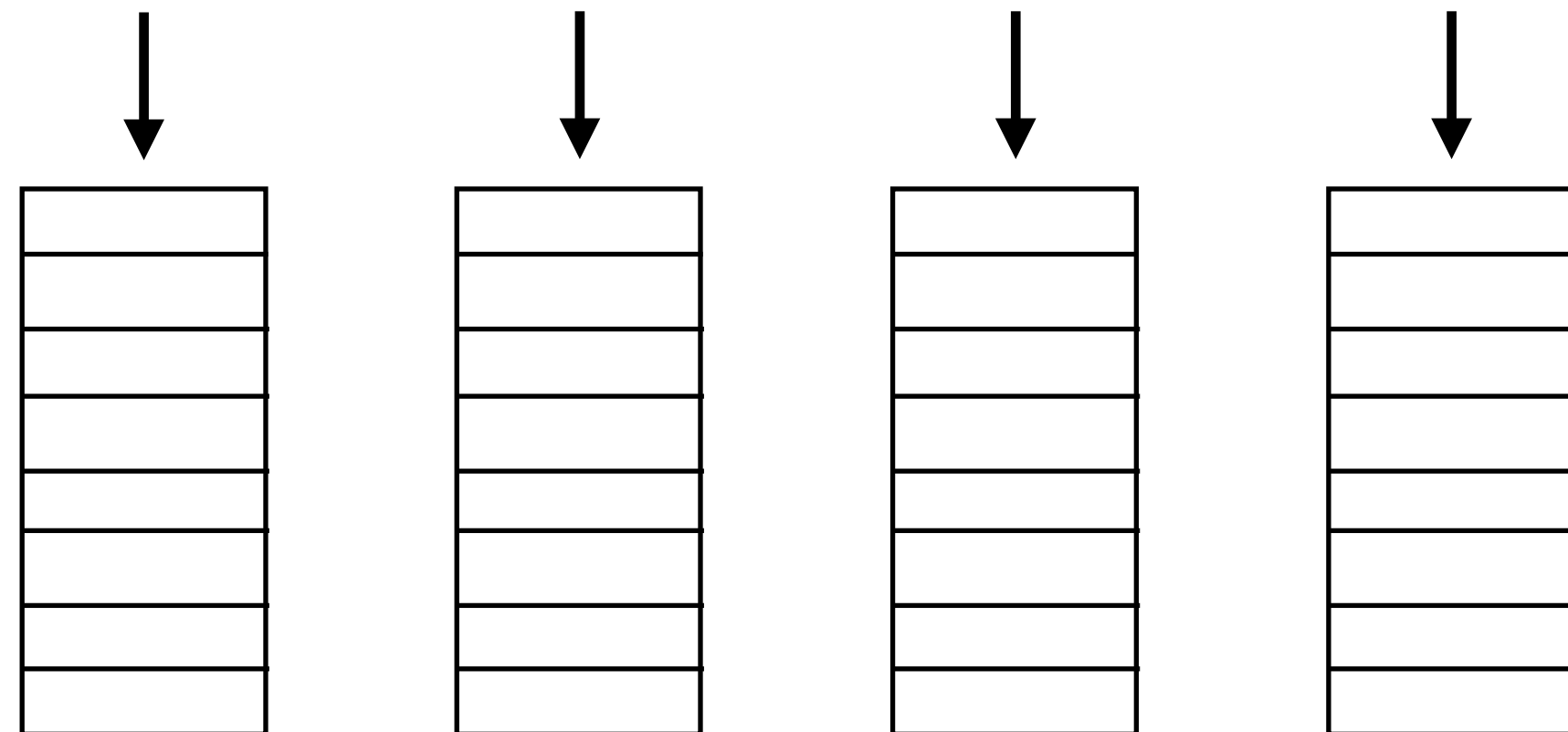
Example: distributed work queues reduce contention

(contention in access to single shared work queue)

Subproblems
(a.k.a. "tasks", "work to do")



Set of work queues
(In general, one per worker thread)



Worker threads:
Pull data from OWN work queue
Push new work to OWN work queue
(no contention when all processors have work to do)



When local work queue is empty...
STEAL work from random work queue
(synchronization okay since processor would have sat idle anyway)

Summary: reducing communication costs

- **Reduce overhead of communication to sender/receiver**
 - **Send fewer messages, make messages larger (amortize overhead)**
 - **Coalesce many small messages into large ones**
- **Reduce latency of communication**
 - **Application writer: restructure code to exploit locality**
 - **Hardware implementor: improve communication architecture**
- **Reduce contention**
 - **Replicate contended resources (e.g., local copies, fine-grained locks)**
 - **Stagger access to contended resources**
- **Increase communication/computation overlap**
 - **Application writer: use asynchronous communication (e.g., async messages)**
 - **HW implementor: pipelining, multi-threading, pre-fetching, out-of-order exec**
 - **Requires additional concurrency in application (more concurrency than number of execution units)**

**Here are some tricks for understanding
the performance of parallel software**

Remember:

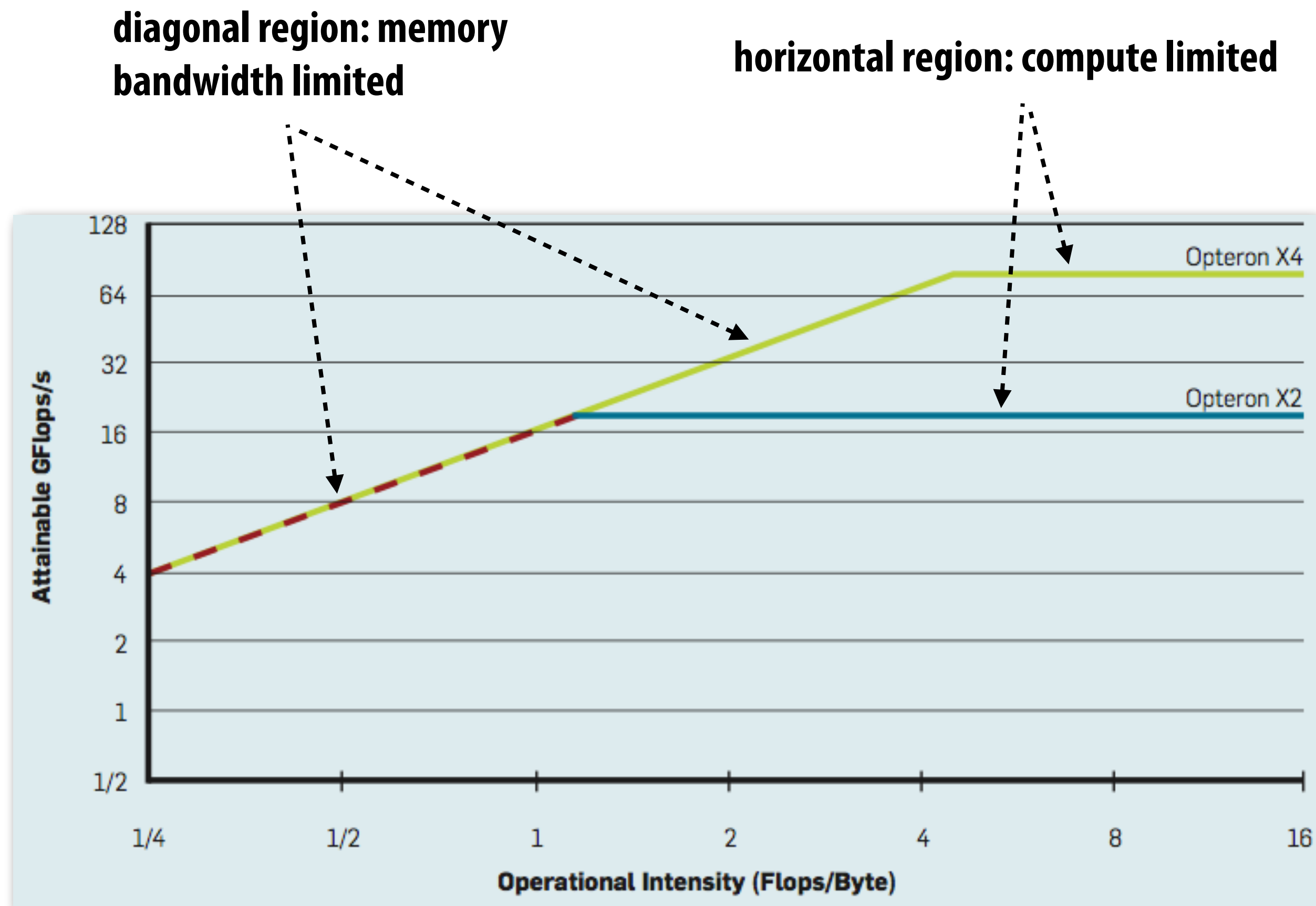
**Always, always, always try the simplest
parallel solution first, then **measure
performance** to see where you stand.**

A useful performance analysis strategy

- **Determine if your performance is limited by computation, memory bandwidth (or memory latency), or synchronization?**
- **Try and establish “high watermarks”**
 - **What’s the best you can do in practice?**
 - **How close is your implementation to a best-case scenario?**

Roofline model

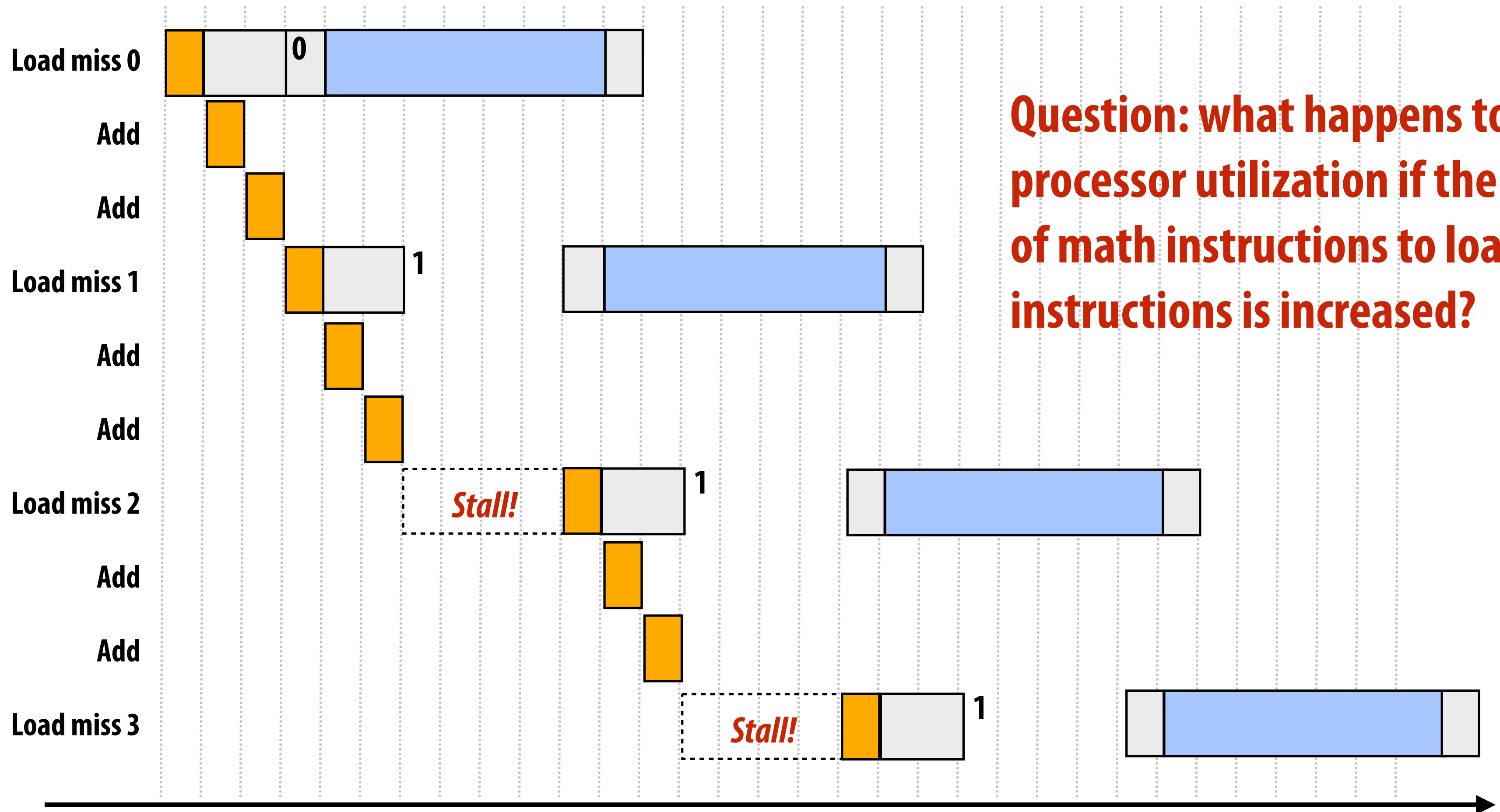
- Use microbenchmarks to compute peak performance of a machine as a function of arithmetic intensity of application
- Then compare application's performance to known peak values



Recall our bandwidth bound example

Cache misses result in request added to size=1 memory request buffer

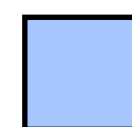
Notice: memory fully utilized, but processor's utilization drop to rate determined by BW



Question: what happens to processor utilization if the ratio of math instructions to load instructions is increased?



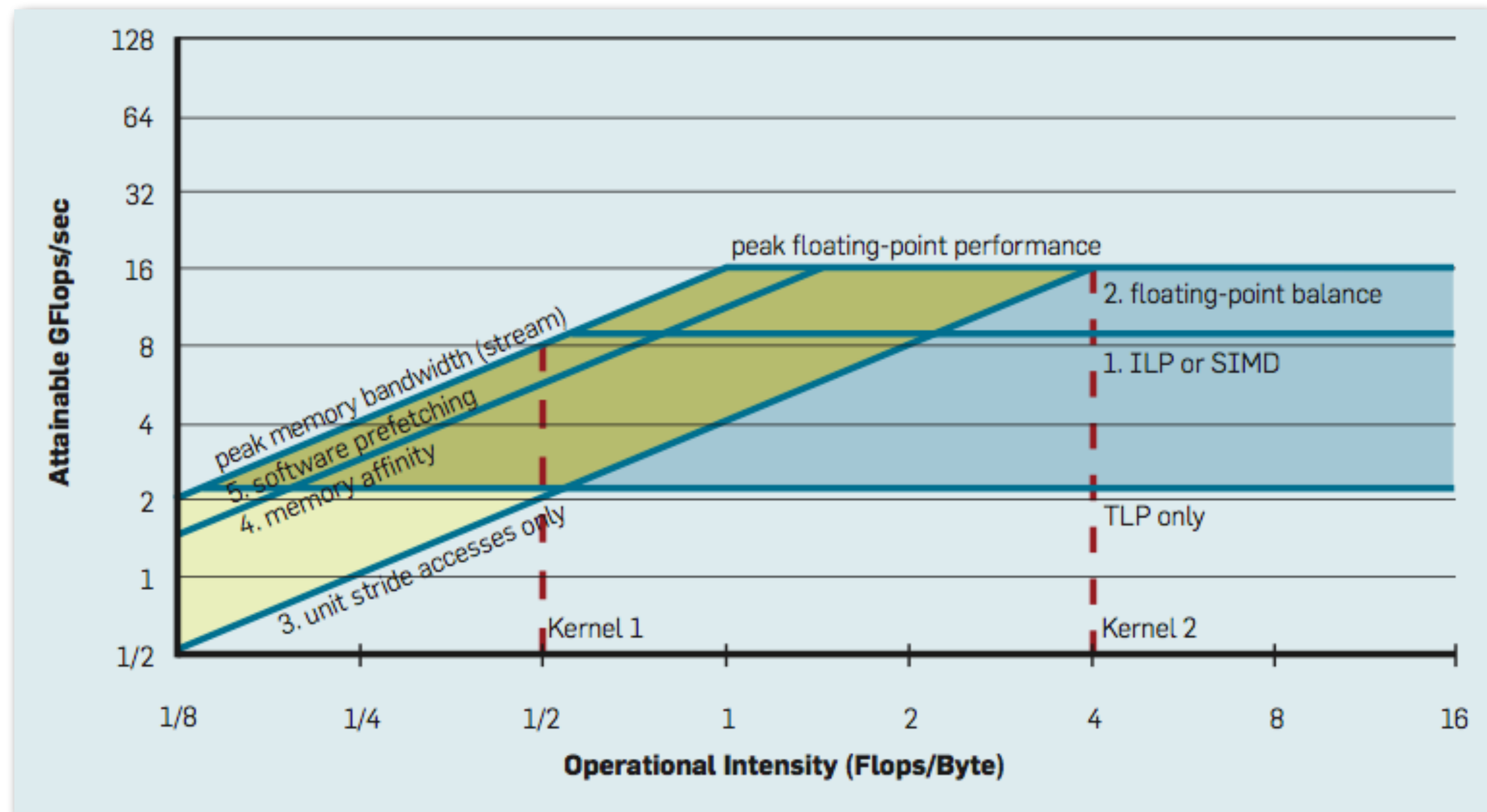
= Instruction



= Occupancy of memory bus = size of cache line / memory bus bandwidth

Roofline model: optimization regions

- Use various levels of optimization in benchmarks (e.g., best performance with and without using SIMD instructions)



Establishing high watermarks *

Add “math” (non-memory instructions)

Does execution time increase linearly with operation count as math is added?

(If so, this is evidence that code is instruction-rate limited)

Remove almost all math, but load same data

How much does execution-time decrease? If not much, suspect memory bottleneck

Change all array accesses to A[0]

How much faster does your code get?

(This establishes an upper bound on benefit of improving locality of data access)

Remove all atomic operations or locks

How much faster does your code get? (provided it still does approximately the same amount of work)

(This establishes an upper bound on benefit of reducing sync overhead.)

* Computation, memory access, and synchronization are almost never perfectly overlapped. As a result, overall performance will rarely be dictated entirely by compute or by bandwidth or by sync. Even so, the sensitivity of performance change to the above program modifications can be a good indication of dominant costs

Use profilers/performance monitoring tools

- Image at left is “CPU usage” from activity monitor in OS X while browsing the web in Chrome (my laptop has a quad-core Core i7 CPU)
 - Graph plots percentage of time OS has scheduled a process thread onto a processor execution context
 - Not very helpful for optimizing performance
- All modern processors have low-level event “performance counters”
 - Registers that count important details such as: instructions completed, clock ticks, L2/L3 cache hits/misses, bytes read from memory controller, etc.
- Example: Intel’s Performance Counter Monitor Tool provides a C++ API for accessing these registers.

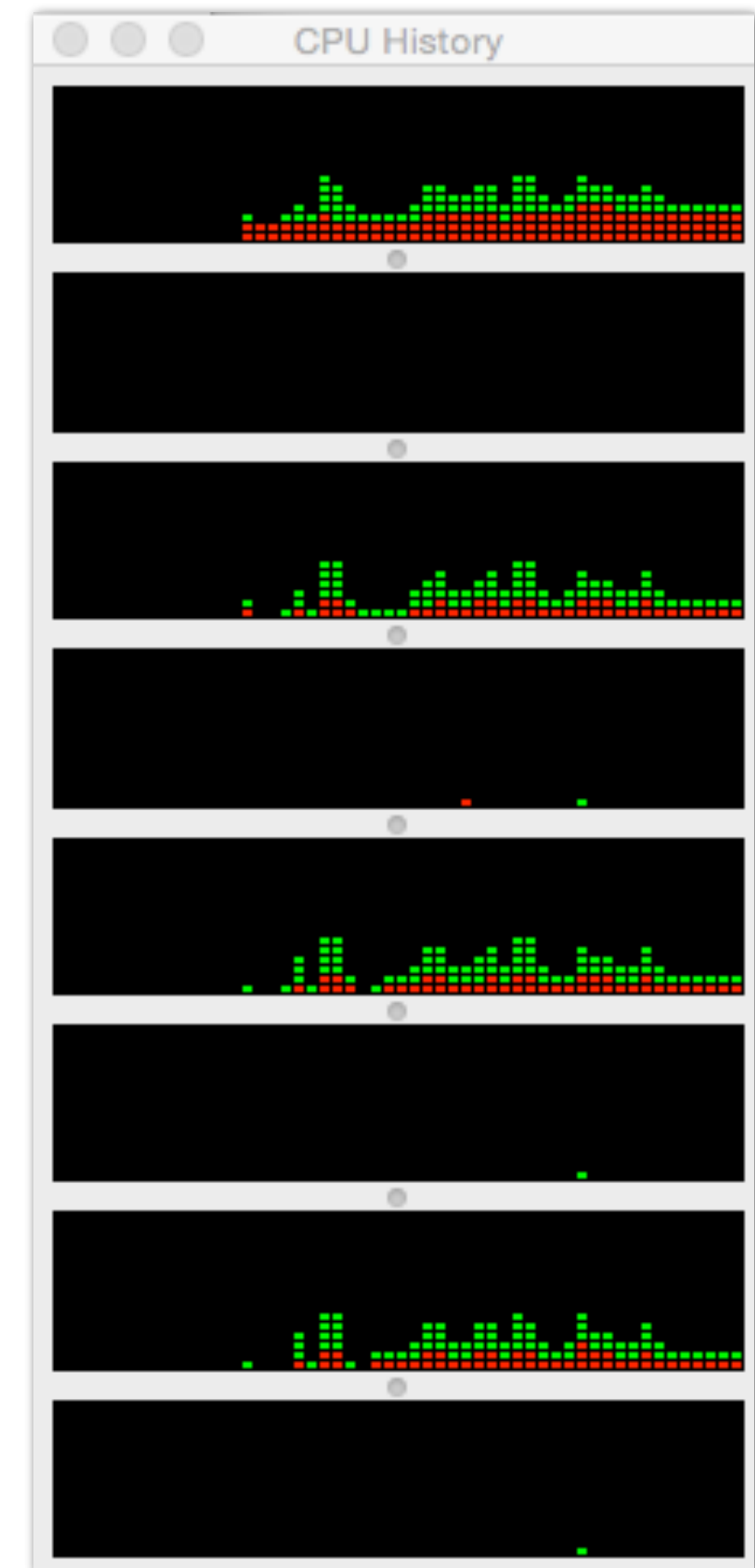
```
PCM *m = PCM::getInstance();
SystemCounterState begin = getSystemCounterState();

// code to analyze goes here

SystemCounterState end = getSystemCounterState();

printf("Instructions per clock: %f\n", getIPC(begin, end));
printf("L3 cache hit ratio: %f\n", getL3CacheHitRatio(begin, end));
printf("Bytes read: %d\n", getBytesReadFromMC(begin, end));
```

- Also see Intel VTune, PAPI, oprofile, etc.



Bonus slides:
**Understanding problem size issues can
very helpful when assessing program
performance**

You are hired by [insert your favorite chip company here].

**You walk in on day one, and your boss says
“All of our senior architects have decided to take the year off.
Your job is to lead the design of our next parallel processor.”**

What questions might you ask?

**Your boss selects the application that matters most to the company
“I want you to demonstrate good performance on this application.”**

How do you know if you have a good design?

■ **Absolute performance?**

- Often measured as wall clock time
- Another example: operations per second

■ **Speedup: performance improvement due to parallelism?**

- Execution time of sequential program / execution time on P processors
- Operations per second on P processors / operations per second of sequential program

■ **Efficiency?**

- Performance per unit resource
- e.g., operations per second per chip area, per dollar, per watt

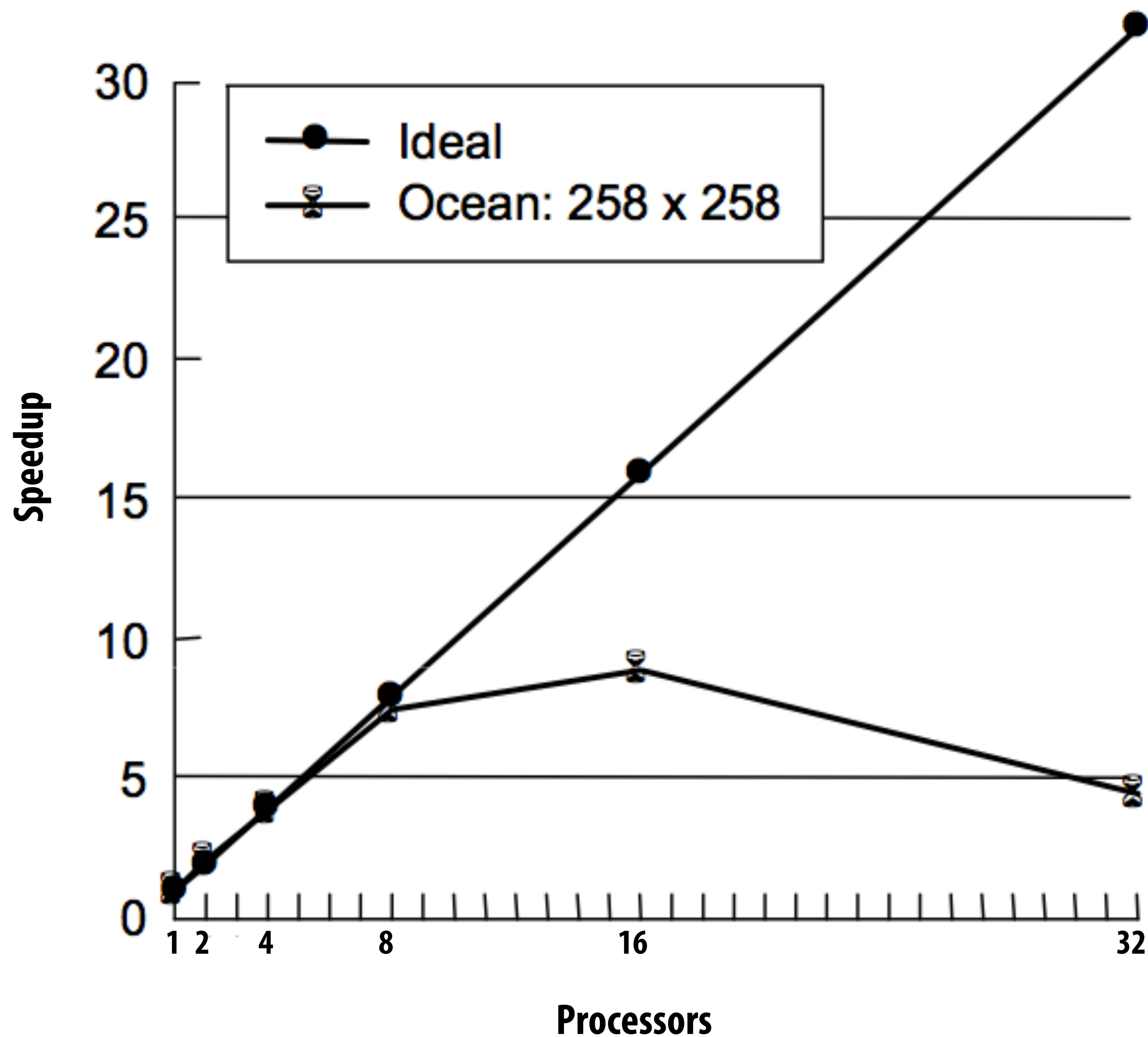
Measuring scaling

- **Consider the grid solver example from last week's class**
 - **We changed the algorithm to allow for parallelism**
 - **The new algorithm might converge more slowly, requiring more iterations of the solver**
- **Should speedup be measured against the performance of a parallel version of a program running on one processor, or the best sequential program?**

Common pitfall: compare parallel program speedup to parallel algorithm running on one core (easier to make yourself look good)

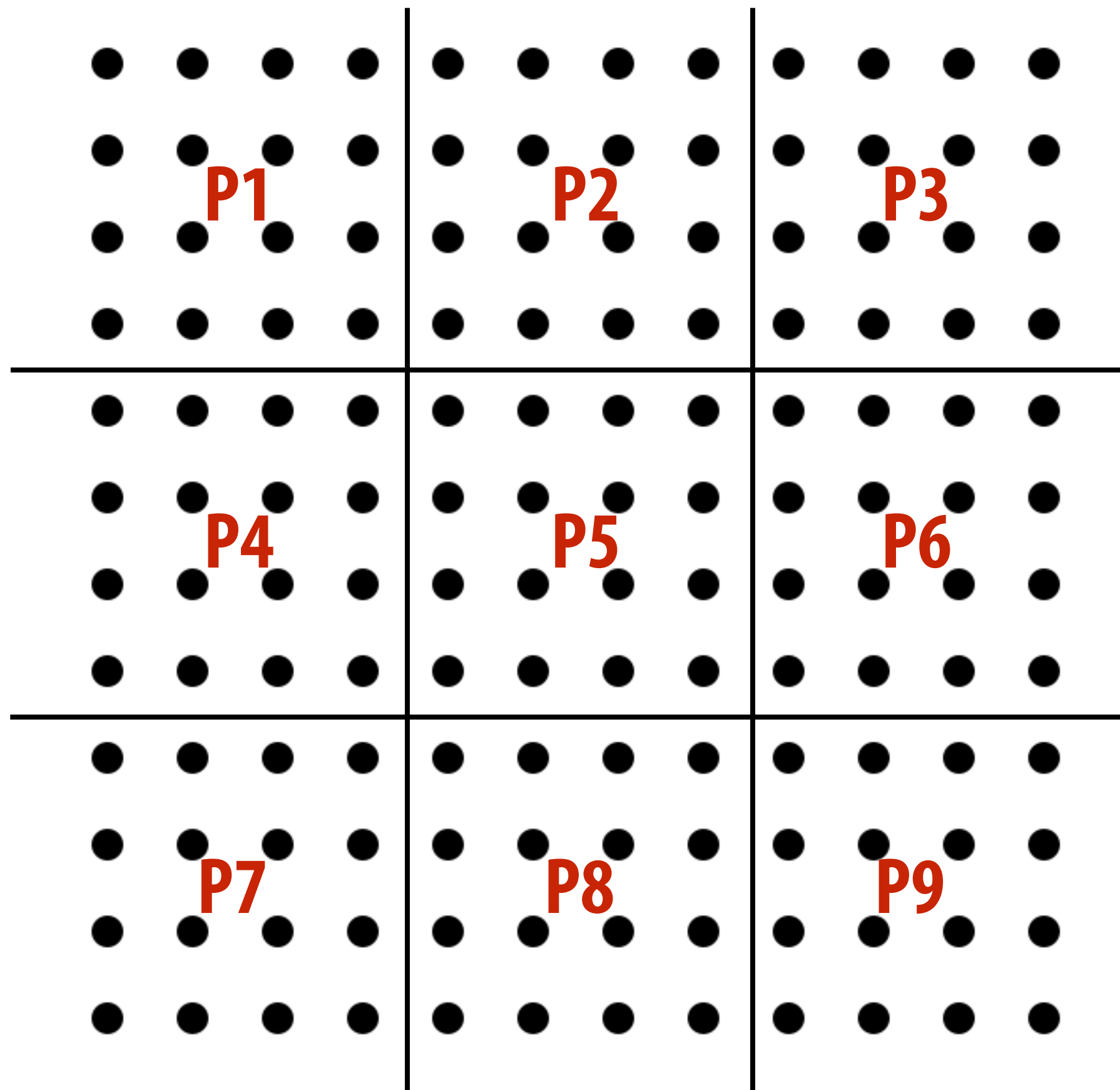
Speedup of solver application: 258 x 258 grid

Execution on 32 processor SGI Origin 2000



Remember: work assignment in solver

2D blocked assignment: $N \times N$ grid



N^2 elements

P processors

elements computed:
(per processor)

$$\frac{N^2}{P}$$

elements communicated:
(per processor)

$$\propto \frac{N}{\sqrt{P}}$$

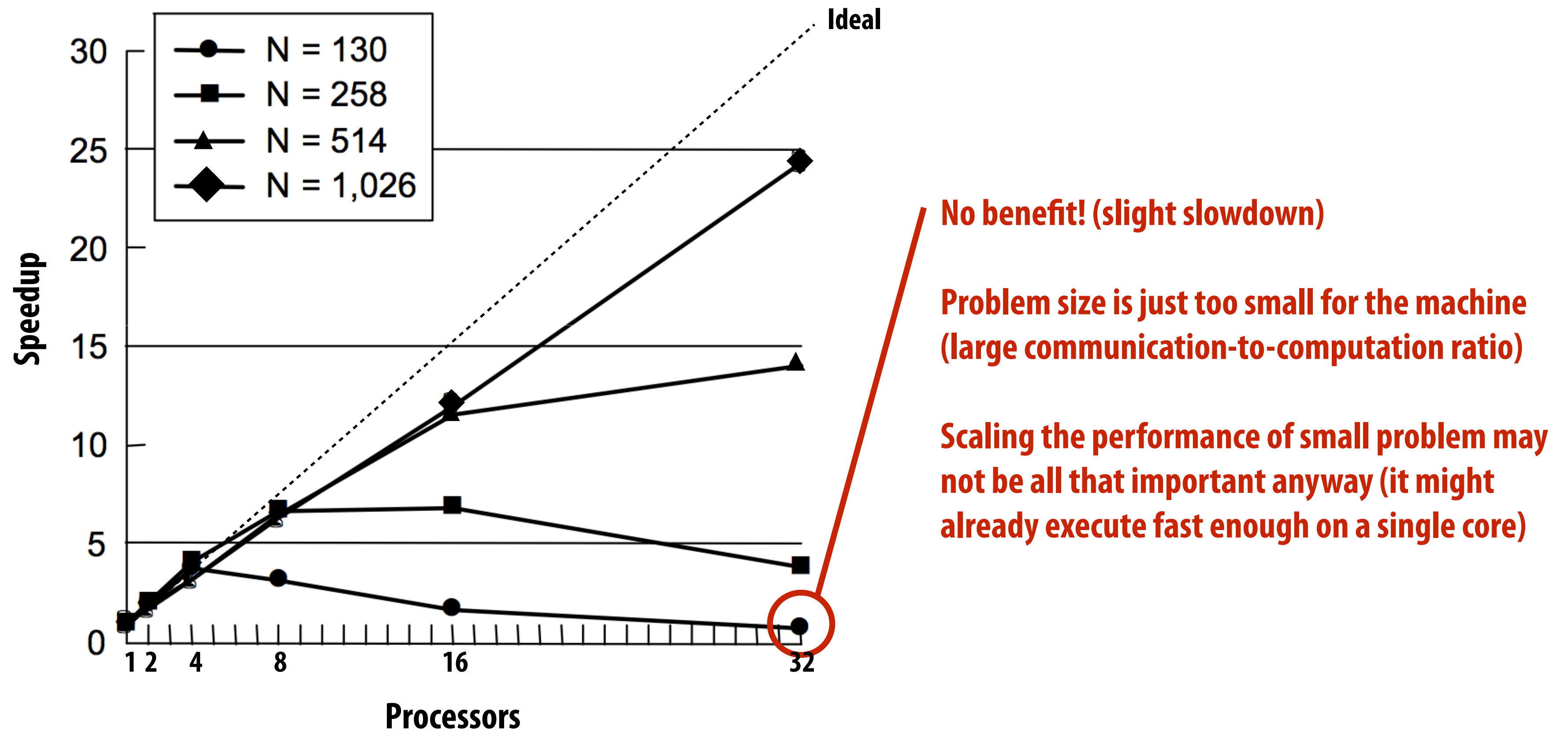
arithmetic intensity:

$$\frac{N}{\sqrt{P}}$$

Small N (or large P) yields low arithmetic intensity!

Pitfalls of fixed problem size speedup analysis

Solver execution on 32 processor SGI Origin 2000

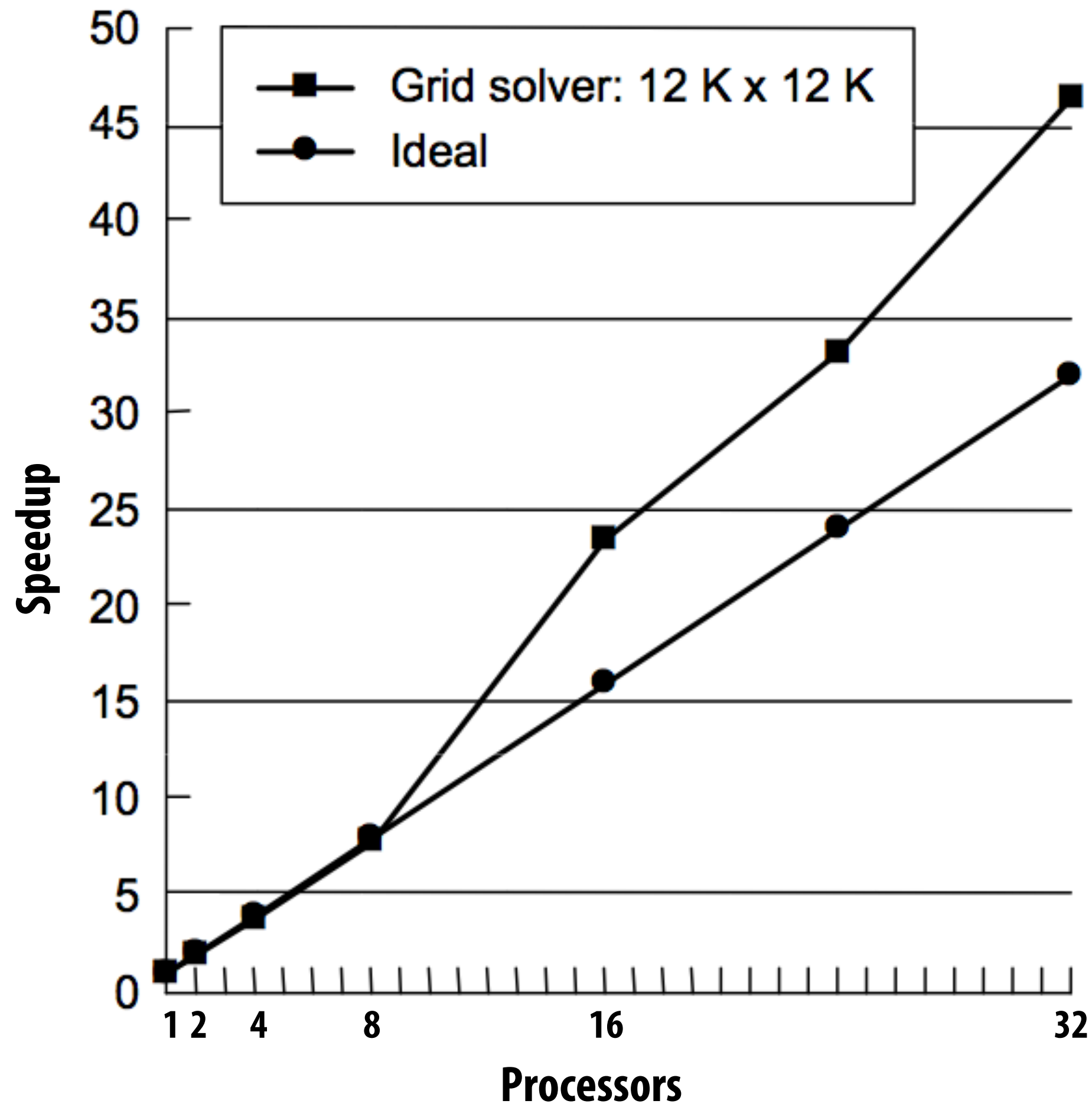


258 x 258 grid on 32 processors: ~ 310 grid cells per processor

1K x 1K grid on 32 processors: ~ 32K grid cells per processor

Pitfalls of fixed problem size speedup analysis

Execution on 32 processor SGI Origin 2000



Here: super-linear speedup! with enough processors, chunk of grid assigned to each processor begins to fit in cache (key working set fits in per-processor cache)

Another example: if problem size is too large for a single machine, working set may not fit in memory: causing thrashing to disk

(this would make speedup on a bigger parallel machine with more memory look amazing!)

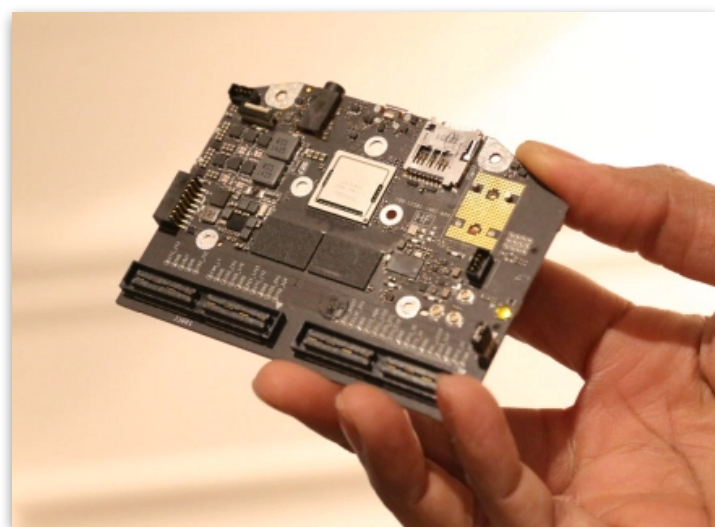
Understanding scaling

- **There can be complex interactions between the size of the problem to solve and the size of the parallel computer**
 - Can impact load balance, overhead, arithmetic intensity, locality of data access
 - Effects can be dramatic and application dependent
- **Evaluating a machine with a fixed problem size can be problematic**
 - **Too small a problem:**
 - Parallelism overheads dominate parallelism benefits (may even result in slow downs)
 - Problem size may be appropriate for small machines, but inappropriate for large ones (does not reflect realistic usage of large machine!)
 - **Too large a problem: (problem size chosen to be appropriate for large machine)**
 - Key working set may not “fit” in small machine (causing thrashing to disk, or key working set exceeds cache capacity, or can’t run at all)
 - When problem working set “fits” in a large machine but not small one, super-linear speedups can occur
- **Can be desirable to scale problem size as machine sizes grow**
(buy a bigger machine to compute more, rather than just compute the same problem faster)

Architects also think about scaling

A common question: “Does an architecture scale?”

- **Scaling up**: how does architecture’s performance scale with increasing core count?
 - Will design scale to the high end?
- **Scaling down**: how does architecture’s performance scale with decreasing core count?
 - Will design scale to the low end?
- **Parallel architectures are designed to work in a range of contexts**
 - Same architecture used for low-end, medium-scale, and high-end systems
 - GPUs are a great example
 - Same SMM core architecture, different numbers of SMM cores per chip



**Tegra X1: 2 SMM cores
(mobile SoC)**



**GTX 950: 6 SMM cores
(90 watts)**



**GTX 980: 16 SMM cores
(165 watts)**



**Titan X: 24 SMM cores
(250 watts)**

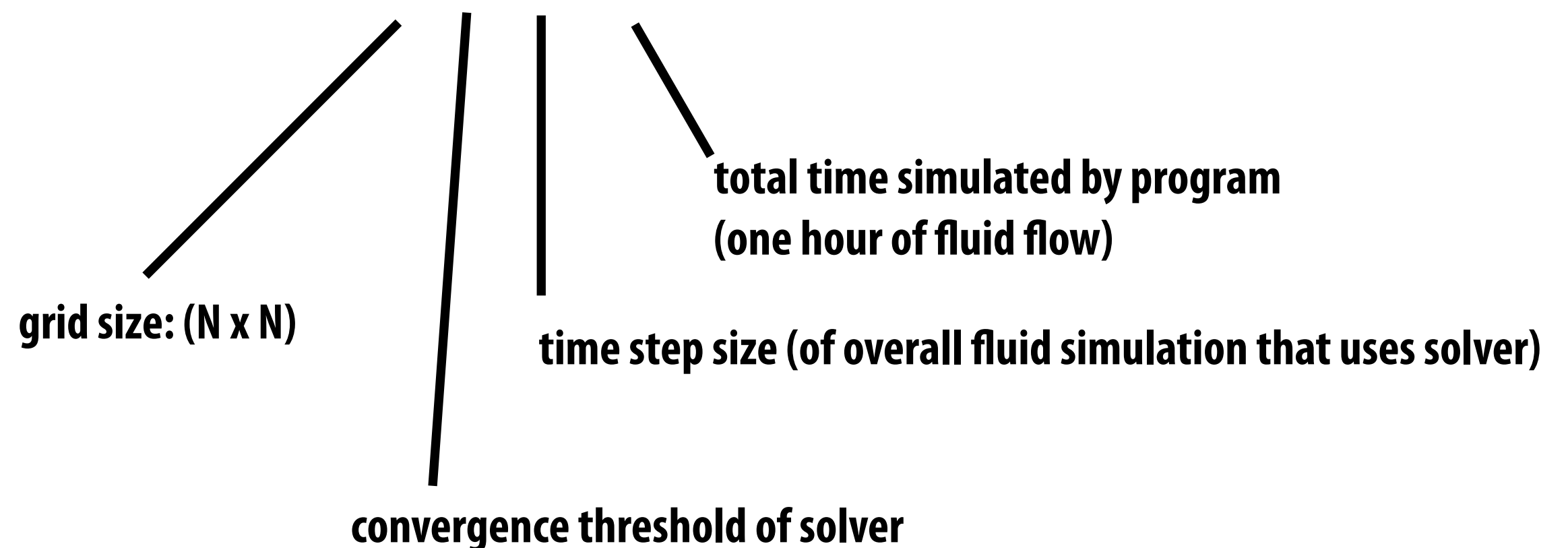
Questions to ask when scaling a problem

■ Under what constraints should the problem be scaled?

- “Work done by program” may no longer be the quantity that is fixed
- Fixed data set size, fixed memory usage per processor, fixed execution time, etc.?

■ How should the problem be scaled?

- Problem size is often determined by more than one parameter
- Solver example: problem defined by $(N, \epsilon, \Delta t, T)$



Problem-constrained scaling *

- Focus: use a parallel computer to solve the same problem faster

$$\text{Speedup} = \frac{\text{time 1 processor}}{\text{time P processors}}$$

- Recall pitfalls from earlier in lecture (small problems may not be realistic workloads for large machines, big problems may not fit on small machines)
- Examples of problem-constrained scaling:
 - Almost everything we've considered parallelizing in class so far

* Problem-constrained scaling is often called "hard scaling".

Time-constrained scaling

- **Focus: completing more work in a fixed amount of time**
 - Execution time kept fixed as the machine (and problem) scales

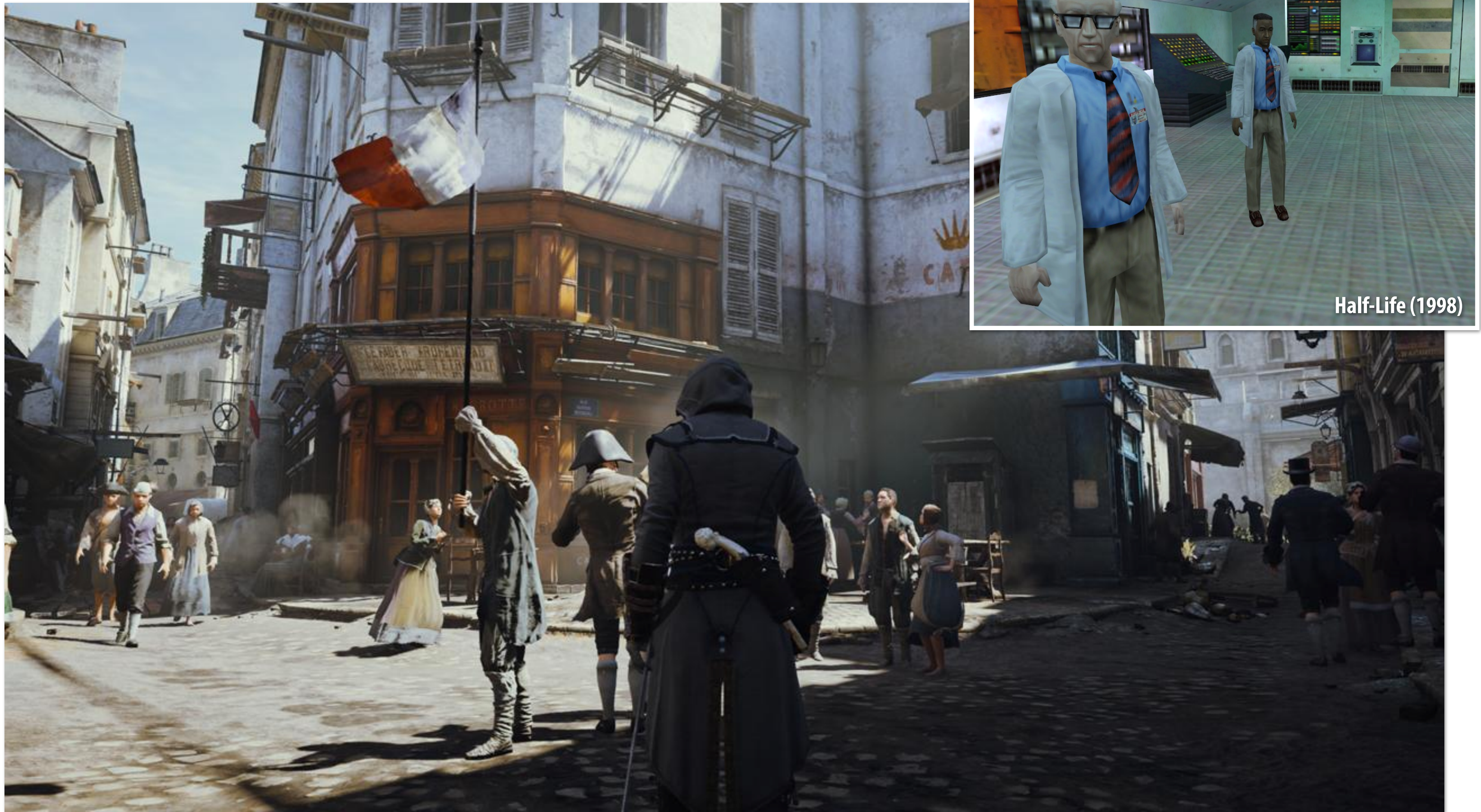
$$\text{Speedup} = \frac{\text{work done by } P \text{ processors}}{\text{work done by 1 processor}}$$

- **How to measure “work”?**
 - Challenge: “work done” may not be linear function of problem inputs (e.g. matrix multiplication is $O(N^3)$ work for $O(N^2)$ sized inputs)
 - One approach: “work done” is defined by execution time of same computation on a single processor (but consider effects of thrashing if problem too big)
 - Ideally, a measure of work is:
 - Simple to understand
 - Scales linearly with sequential run time (so ideal speedup remains linear in P)

Time-constrained scaling example

Real-time 3D graphics: more compute power allows for rendering of much more complex scene

Problem size metrics: number of polygons, texels sampled, shader length, etc.



Assassin's Creed Unity (2014)

Image credits:

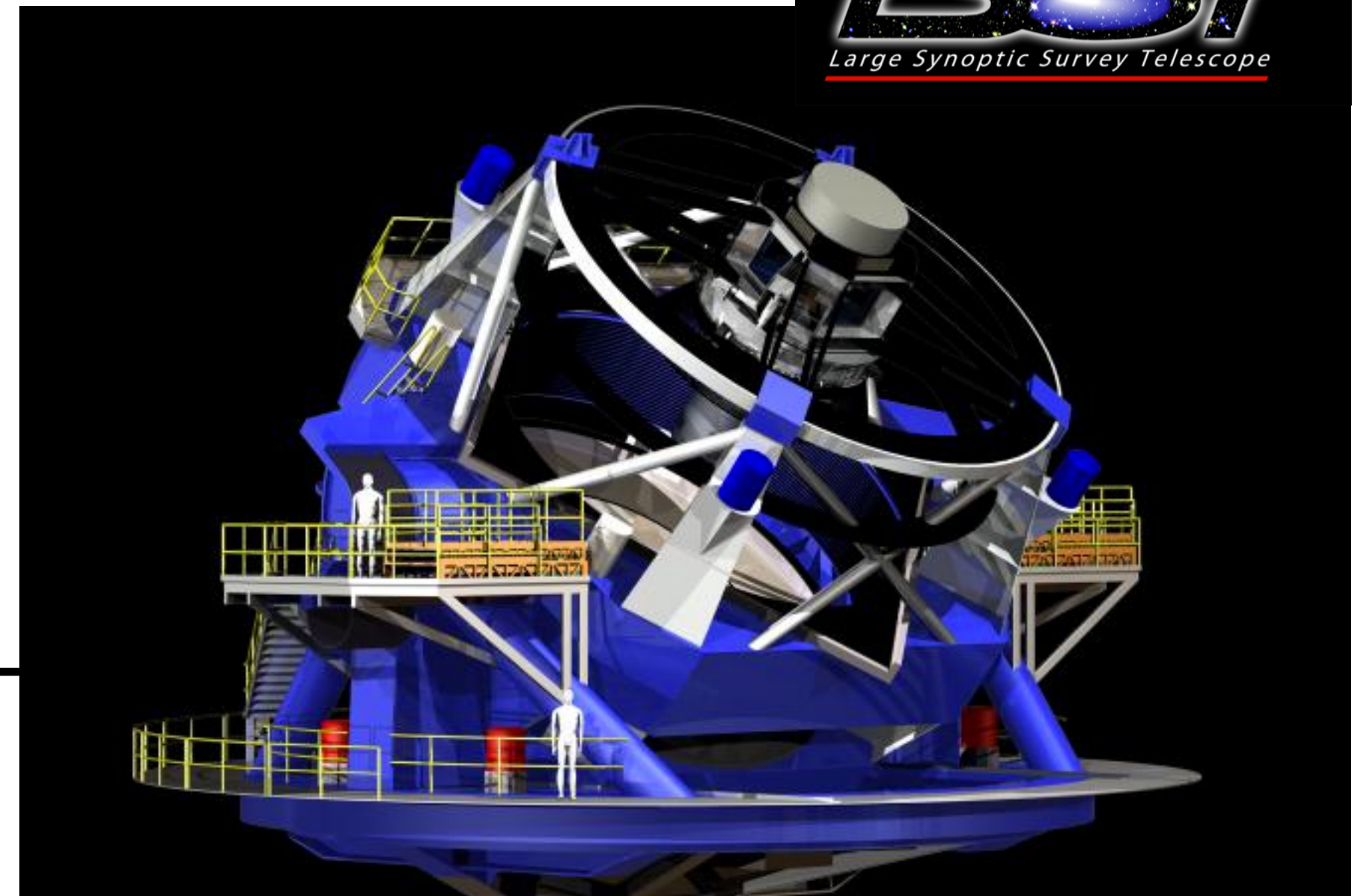
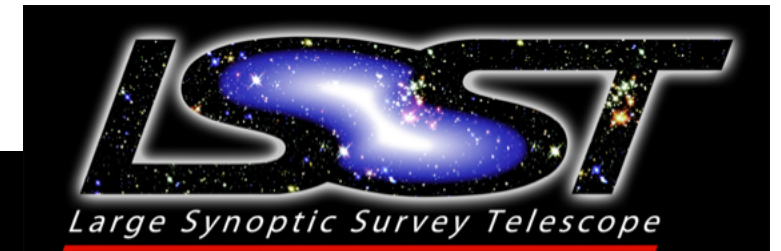
<http://www.gamespot.com/forums/system-wars-314159282/assassin-s-creed-unity-best-graphics-of-2014-31696528/>

http://www.game-weavers.com/?page_id=490

Time-constrained scaling example

Large Synoptic Survey Telescope (LSST)

- Estimated completion in 2019
- Acquire high-resolution survey of sky (3-gigapixel image every 15 seconds, every night for many years)



LSST will be located on top of Cerro Pachón Mountain, Chile

Rapid Image analysis compute platform
(detect "potentially"
interesting events)

Notify other observatories if
potential event detected.



Increasing compute capability allows
for more sophisticated detection
algorithms (fewer false positives,
detect broader class of events)

More time-constrained scaling examples

- **Computational finance**

- **Run most sophisticated model possible in: 1 ms, 1 minute, overnight, etc.**

- **Modern web sites**

- **Want to generate complex page, respond to user in X milliseconds (studies show site usage directly corresponds to page load latency)**

- **Real-time computer vision for robotics**

- **Consider self-driving car: want best-quality obstacle detection in 5 ms**

Memory-constrained scaling *

- **Focus: run the largest problem possible without overflowing main memory ****
- **Memory per processor is held fixed (e.g., add more machines to cluster)**
- **Neither work or execution time are held constant**

$$\text{Speedup} = \frac{\text{work (P processors)} \times \text{time (1 processor)}}{\text{time (P processors)} \times \text{work (1 processor)}}$$
$$= \frac{\text{work per unit time on P processors}}{\text{work per unit time on 1 processor}}$$

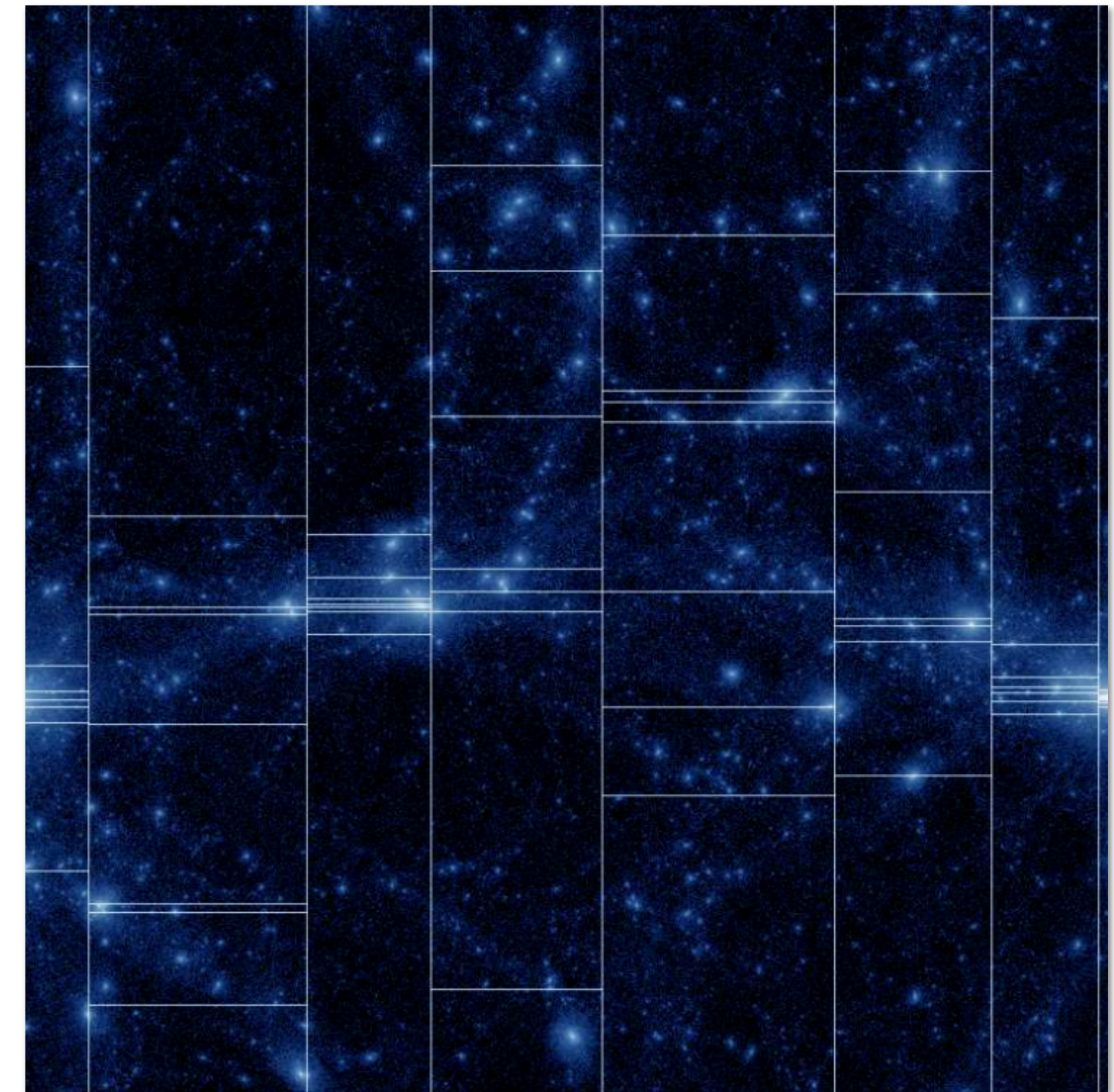
- **Note: scaling problem size can make runtimes very large**
 - **Consider $O(N^3)$ matrix multiplication on $O(N^2)$ matrices**

* Memory-constrained scaling is often called “weak scaling”

** Assumptions: (1) memory resources scale with processor count (2) spilling to disk is infeasible behavior (too slow)

Memory-constrained scaling examples

- **One motivation to use supercomputers and large clusters is simply to be able to fit large problems in memory**
- **Large N-body problems**
 - **2012 Supercomputing Gordon Bell Prize Winner: 1,073,741,824,000 particle N-body simulation on Japan's K-Computer**
- **Large-scale machine learning**
 - **Billions of clicks, documents, etc.**
- **Memcached (in memory caching system for web apps)**
 - **More servers = more available cache**



2D domain decomposition of N-body simulation

Scaling examples at PIXAR

- **Rendering a “shot” (a sequence of frames) in a movie**
 - **Goal: minimize time to completion (problem constrained)**
 - **Assign each frame to a different machine in the cluster**
- **Artists working to design lighting for a scene**
 - **Provide interactive frame rate to artist (time constrained)**
 - **More performance = higher fidelity representation shown to artist in allotted time**
- **Physical simulation: e.g., fluid simulation**
 - **Parallelize simulation across multiple machines to fit simulation grid in aggregate memory of processors (memory constrained)**
- **Final render of images for movie**
 - **Scene complexity is typically bounded by memory available on farm machines**
 - **One barrier to exploiting additional parallelism within a machine is that required footprint often increases with number of processors**



Summary of tips

- **Measure, measure, measure...**
- **Establish high watermarks for your program**
 - **Are you compute, synchronization, or bandwidth bound?**
- **Be aware of scaling issues. Is the problem well matched for the machine?**