

Computer Architecture

Lecture 19: Multiprocessors, Consistency, Coherence

Prof. Onur Mutlu

ETH Zürich

Fall 2017

29 November 2017

Summary of Last Week's Lectures

- Memory Latency Tolerance
- Runahead Execution
- Wrong Path Effects
- Prefetching

Today

- Multiprocessors
- Memory Consistency
- Cache Coherence

Multiprocessors and Issues in Multiprocessing

Readings: Multiprocessing

■ Required

- Amdahl, “[Validity of the single processor approach to achieving large scale computing capabilities](#),” AFIPS 1967.

■ Recommended

- Mike Flynn, “[Very High-Speed Computing Systems](#),” Proc. of IEEE, 1966
- Hill, Jouppi, Sohi, “[Multiprocessors and Multicomputers](#),” pp. 551-560 in Readings in Computer Architecture.
- Hill, Jouppi, Sohi, “[Dataflow and Multithreading](#),” pp. 309-314 in Readings in Computer Architecture.

Memory Consistency

- Required

- Lamport, "How to Make a Multiprocessor Computer That Correctly Executes Multiprocess Programs," IEEE Transactions on Computers, 1979

Readings: Cache Coherence

■ Required

- Papamarcos and Patel, “A low-overhead coherence solution for multiprocessors with private cache memories,” ISCA 1984.

■ Recommended:

- Culler and Singh, *Parallel Computer Architecture*
 - Chapter 5.1 (pp 269 – 283), Chapter 5.3 (pp 291 – 305)
- P&H, *Computer Organization and Design*
 - Chapter 5.8 (pp 534 – 538 in 4th and 4th revised eds.)

D-INFK Distinguished Lecture (Monday)

- Monday, 4 December 2017
- Prof. Michael Scott (Univ of Rochester)
- 16:15-17:15 @ CAB G61

Data Structures for Persistent Memory

4 December 2017 | Colloquium

ABSTRACT: Several forthcoming successors to DRAM have the intriguing property that they retain their content when power is lost. Wide availability of such memory raises the possibility that long-lived, pointer-rich data might be maintained directly in main memory, across program runs and even system crashes. Special care must be taken, however, to assure that the values in memory will be mutually consistent in the wake of a crash, given that caches write values back in arbitrary order. This talk will survey research currently under way at the University of Rochester that aims to - formally specify what it means for a persistent data structure to be implemented correctly, - implement failure atomicity at reasonable cost on current and future architectures, - protect memory-mapped persistent structures from buggy applications.

CHALLENGE: CRASH CONSISTENCY

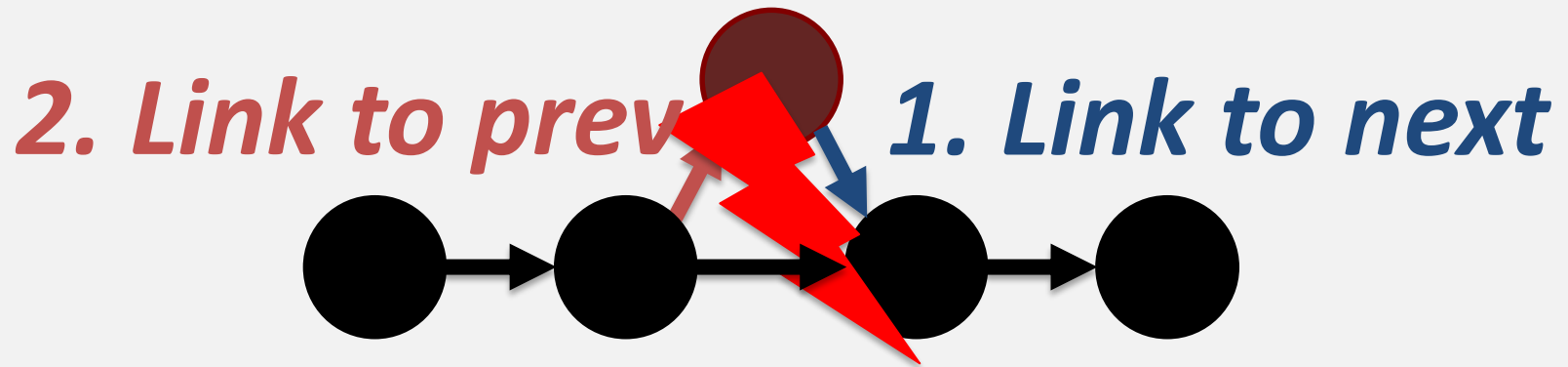


Persistent Memory System

**System crash can result in
permanent data corruption in NVM**

CRASH CONSISTENCY PROBLEM

Example: Add a node to a linked list



**System crash can result in
inconsistent memory state**

One Key Challenge in Persistent Memory

- How to ensure consistency of system/data if all memory is persistent?
- Two extremes
 - Programmer transparent: Let the system handle it
 - Programmer only: Let the programmer handle it
- Many alternatives in-between...

CURRENT SOLUTIONS

Explicit interfaces to manage consistency

– NV-Heaps [ASPLOS'11], BPFS [SOSP'09], Mnemosyne [ASPLOS'11]

```
AtomicBegin {  
    Insert a new node;  
} AtomicEnd;
```

Limits adoption of NVM

Have to rewrite code with clear partition
between volatile and non-volatile data

Burden on the programmers

How do we make legacy code work?

A NEW APPROACH: ThyNVM

Goal:
**Software transparent consistency in
persistent memory systems**

ThyNVM: Summary

A new hardware-based checkpointing mechanism

- **Checkpoints** at *multiple granularities* to reduce both checkpointing latency and metadata overhead
- **Overlaps** *checkpointing* and *execution* to reduce checkpointing latency
- **Adapts** to *DRAM and NVM* characteristics

Performs within **4.9%** of an *idealized DRAM* with zero cost consistency

More About ThyNVM

- Jinglei Ren, Jishen Zhao, Samira Khan, Jongmoo Choi, Yongwei Wu, and Onur Mutlu,
"ThyNVM: Enabling Software-Transparent Crash Consistency in Persistent Memory Systems"
*Proceedings of the 48th International Symposium on Microarchitecture (**MICRO**), Waikiki, Hawaii, USA, December 2015.*
[[Slides \(pptx\)](#)] [[pdf](#)] [[Lightning Session Slides \(pptx\)](#)] [[pdf](#)] [[Poster \(pptx\)](#)] [[pdf](#)]
[[Source Code](#)]

ThyNVM: Enabling Software-Transparent Crash Consistency in Persistent Memory Systems

Jinglei Ren^{*†} Jishen Zhao[‡] Samira Khan^{†'} Jongmoo Choi^{+†} Yongwei Wu^{*} Onur Mutlu[†]

[†]Carnegie Mellon University ^{*}Tsinghua University

[‡]University of California, Santa Cruz [']University of Virginia ⁺Dankook University

Programming Ease to Exploit Persistence

Tools/Libraries to Help Programmers

- Himanshu Chauhan, Irina Calciu, Vijay Chidambaram, Eric Schkufza, Onur Mutlu, and Pratap Subrahmanyam,
"NVMove: Helping Programmers Move to Byte-Based Persistence"

*Proceedings of the 4th Workshop on Interactions of NVM/Flash with Operating Systems and Workloads (**INFLOW**), Savannah, GA, USA, November 2016.*

[[Slides \(pptx\)](#) ([pdf](#))]

NVMOVE: Helping Programmers Move to Byte-Based Persistence

Himanshu Chauhan *

UT Austin

Irina Calciu

VMware Research Group

Vijay Chidambaram

UT Austin

Eric Schkufza

VMware Research Group

Onur Mutlu

ETH Zürich

Pratap Subrahmanyam

VMware

Multiprocessors and Issues in Multiprocessing

Remember: Flynn's Taxonomy of Computers

- Mike Flynn, “**Very High-Speed Computing Systems**,” Proc. of IEEE, 1966
- **SISD**: Single instruction operates on single data element
- **SIMD**: Single instruction operates on multiple data elements
 - Array processor
 - Vector processor
- **MISD**: Multiple instructions operate on single data element
 - Closest form: systolic array processor, streaming processor
- **MIMD**: Multiple instructions operate on multiple data elements (multiple instruction streams)
 - Multiprocessor
 - Multithreaded processor

Why Parallel Computers?

- **Parallelism: Doing multiple things at a time**
- **Things: instructions, operations, tasks**
- **Main (or Original) Goal**
 - **Improve performance (Execution time or task throughput)**
 - Execution time of a program governed by Amdahl's Law
- **Other Goals**
 - **Reduce power consumption**
 - (4N units at freq F/4) consume less power than (N units at freq F)
 - Why?
 - **Improve cost efficiency and scalability, reduce complexity**
 - Harder to design a single unit that performs as well as N simpler units
 - **Improve dependability: Redundant execution in space**

Types of Parallelism and How to Exploit Them

■ Instruction Level Parallelism

- Different instructions within a stream can be executed in parallel
- Pipelining, out-of-order execution, speculative execution, VLIW
- Dataflow

■ Data Parallelism

- Different pieces of data can be operated on in parallel
- SIMD: Vector processing, array processing
- Systolic arrays, streaming processors

■ Task Level Parallelism

- Different “tasks/threads” can be executed in parallel
- Multithreading
- Multiprocessing (multi-core)

Task-Level Parallelism: Creating Tasks

- Partition a single problem into multiple related tasks (threads)
 - Explicitly: Parallel programming
 - Easy when tasks are natural in the problem
 - Web/database queries
 - Difficult when natural task boundaries are unclear
 - Transparently/implicitly: Thread level speculation
 - Partition a single thread speculatively
- Run many independent tasks (processes) together
 - Easy when there are many processes
 - Batch simulations, different users, cloud computing workloads
 - Does not improve the performance of a single task

Multiprocessing Fundamentals

Multiprocessor Types

- Loosely coupled multiprocessors
 - No shared global memory address space
 - Multicomputer network
 - Network-based multiprocessors
 - Usually programmed via message passing
 - Explicit calls (send, receive) for communication

- Tightly coupled multiprocessors
 - Shared global memory address space
 - Traditional multiprocessing: symmetric multiprocessing (SMP)
 - Existing multi-core processors, multithreaded processors
 - Programming model similar to uniprocessors (i.e., multitasking uniprocessor) except
 - Operations on shared data require synchronization

Main Design Issues in Tightly-Coupled MP

- Shared memory synchronization
 - How to handle locks, atomic operations
- Cache coherence
 - How to ensure correct operation in the presence of private caches
- Memory consistency: Ordering of memory operations
 - What should the programmer expect the hardware to provide?
- Shared resource management
- Communication: Interconnects

Main Programming Issues in Tightly-Coupled MP

- Load imbalance

- How to partition a single task into multiple tasks

- Synchronization

- How to synchronize (efficiently) between tasks
- How to communicate between tasks
- Locks, barriers, pipeline stages, condition variables, semaphores, atomic operations, ...

- Ensuring correct operation while optimizing for performance

Aside: Hardware-based Multithreading

- Coarse grained
 - Quantum based
 - Event based (switch-on-event multithreading), e.g., switch on L3 miss

- Fine grained
 - Cycle by cycle
 - Thornton, “[CDC 6600: Design of a Computer](#),” 1970.
 - Burton Smith, “[A pipelined, shared resource MIMD computer](#),” ICPP 1978.

- Simultaneous
 - Can dispatch instructions from multiple threads at the same time
 - Good for improving execution unit utilization

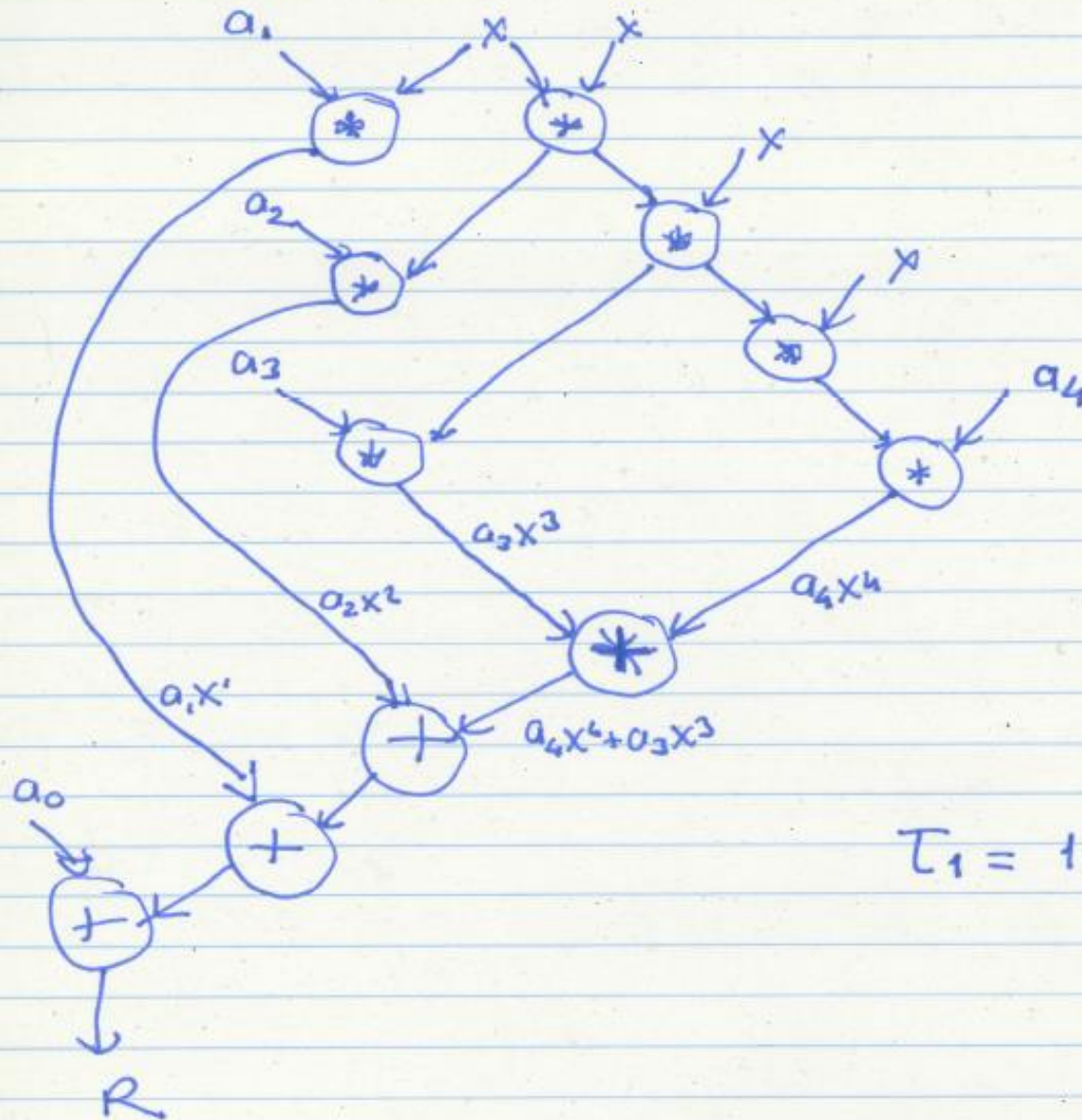
Limits of Parallel Speedup

Parallel Speedup Example

- $a_4x^4 + a_3x^3 + a_2x^2 + a_1x + a_0$
- Assume given inputs: x and each a_i
- Assume each operation 1 cycle, no communication cost, each op can be executed in a different processor
- How fast is this with a single processor?
 - Assume no pipelining or concurrent execution of instructions
- How fast is this with 3 processors?

$$R = a_4x^4 + a_3x^3 + a_2x^2 + a_1x + a_0$$

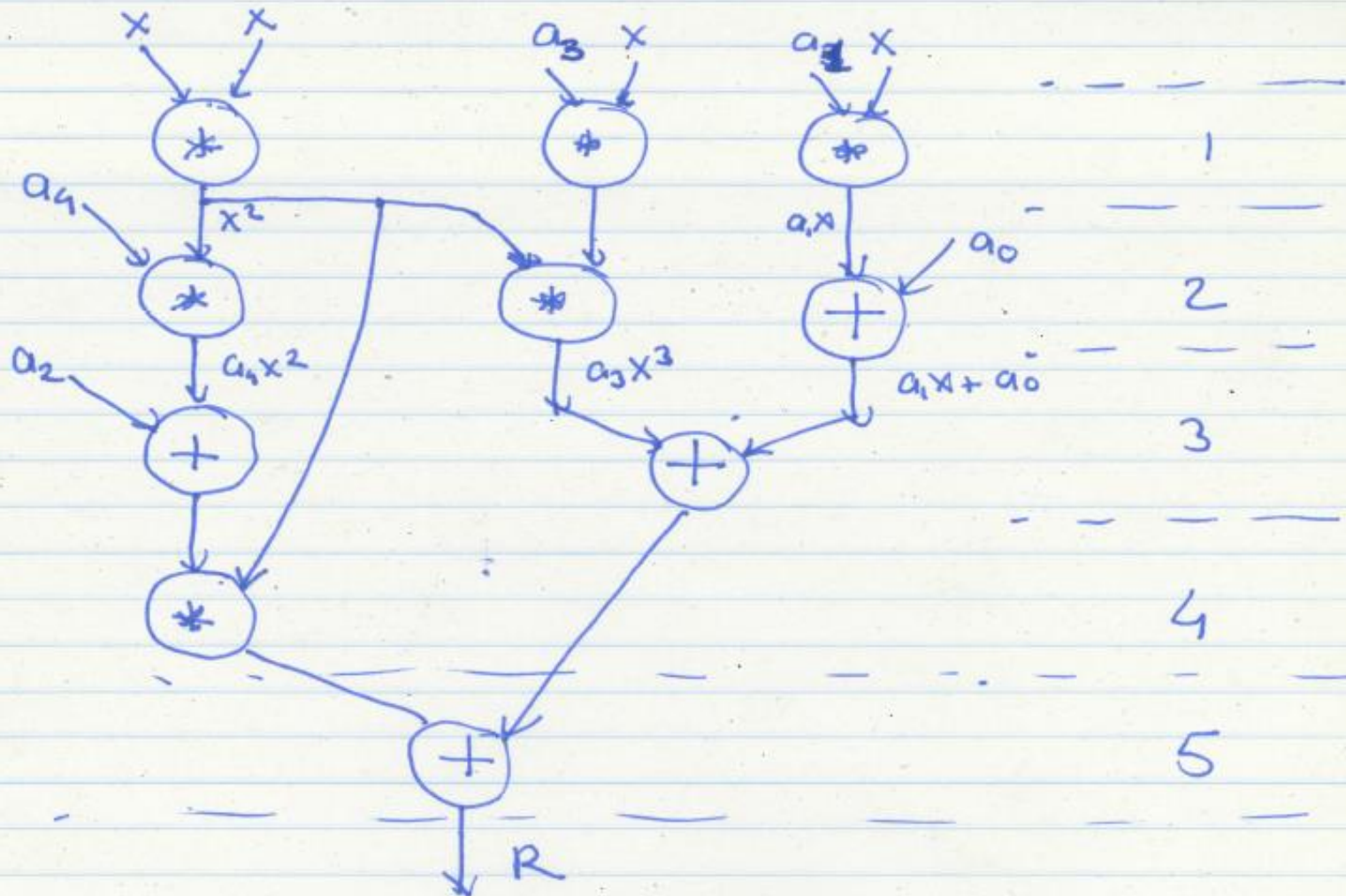
Single processor : 11 operations (data flow graph) ^{draw the}



$T_1 = 11$ cycles

$$R = a_4x^4 + a_3x^3 + a_2x^2 + a_1x + a_0$$

Three processors : T_3 (excc.time with 3 proc.)



$$T_3 = \underline{5 \text{ cycles}}$$

Speedup with 3 Processors

$$T_3 = \underline{5 \text{ cycles}}$$

$$\text{Speedup with 3 processors} = \frac{11}{5} = 2.2$$

$$\left(\frac{T_1}{T_3} \right)$$

Is this a fair comparison?

Revisiting the Single-Processor Algorithm

Revisit T_1

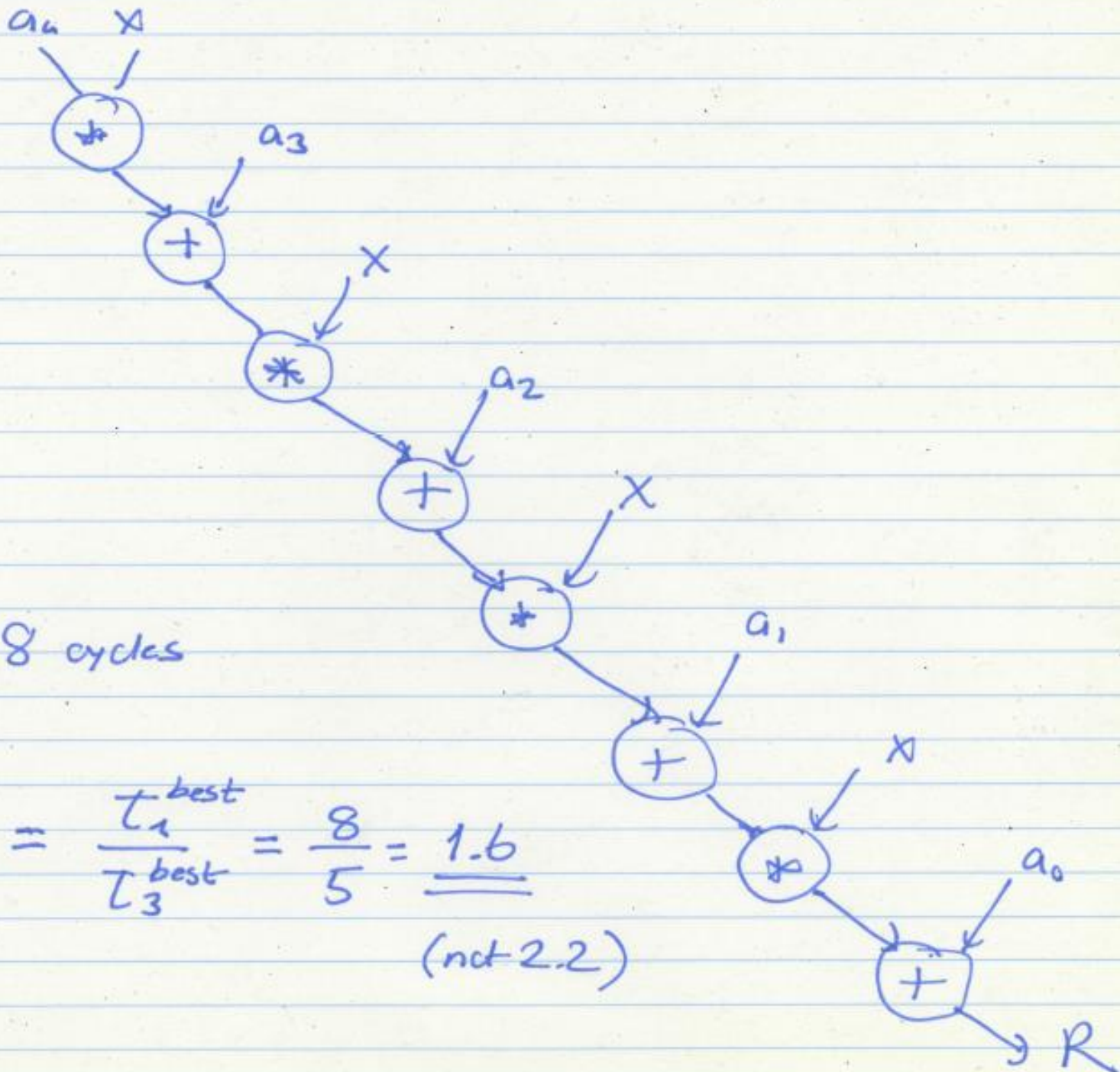
Better single-processor algorithm:

$$R = a_4 x^4 + a_3 x^3 + a_2 x^2 + a_1 x + a_0$$

$$R = (((a_4 x + a_3) x + a_2) x + a_1) x + a_0$$

(Horner's method)

Horner, "A new method of solving numerical equations of all orders, by continuous approximation," Philosophical Transactions of the Royal Society, 1819.

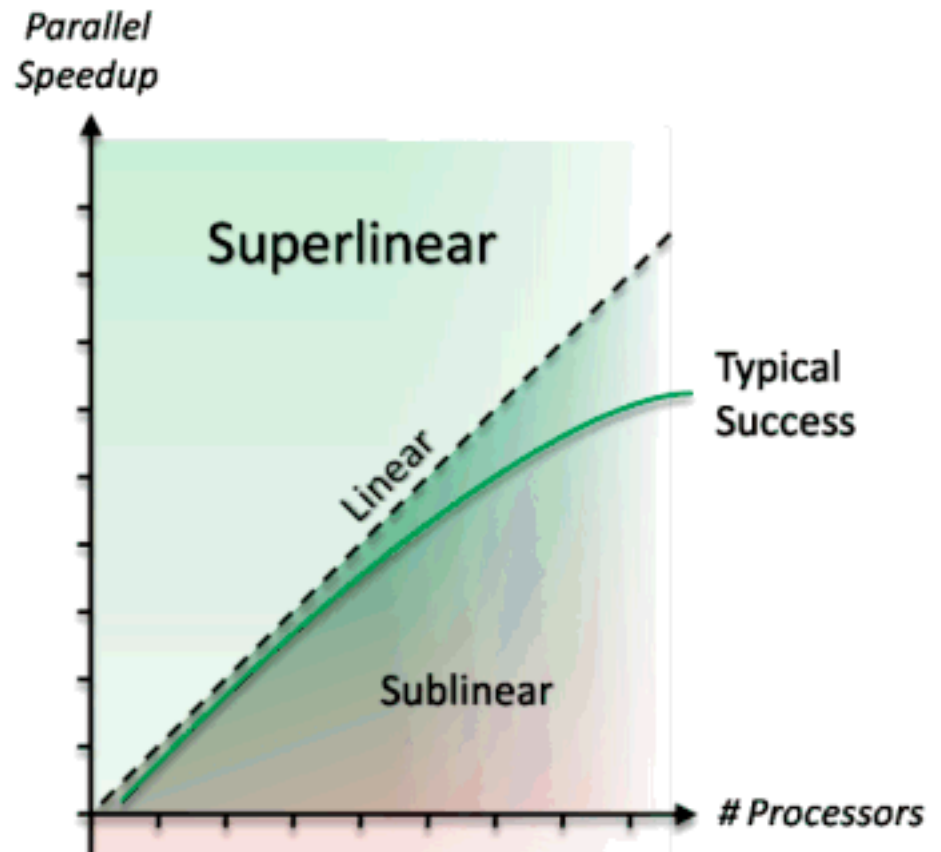


$T_1 = 8$ cycles

Speedup with 3 procs. $= \frac{T_1^{\text{best}}}{T_3^{\text{best}}} = \frac{8}{5} = \underline{\underline{1.6}}$
(not 2.2)

Superlinear Speedup

- Can speedup be greater than P with P processing elements?
- **Unfair comparisons**
Compare best parallel algorithm to wimpy serial algorithm \rightarrow unfair
- **Cache/memory effects**
More processors \rightarrow
more cache or memory \rightarrow
fewer misses in cache/mem



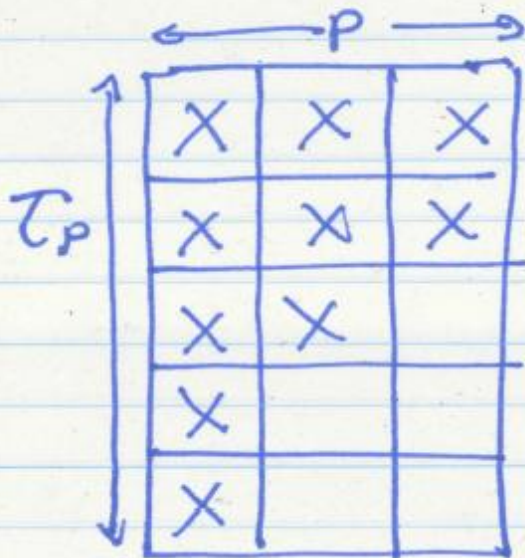
Utilization, Redundancy, Efficiency

- Traditional metrics
 - Assume all P processors are tied up for parallel computation
- Utilization: How much processing capability is used
 - $U = (\# \text{ Operations in parallel version}) / (\text{processors} \times \text{Time})$
- Redundancy: how much extra work is done with parallel processing
 - $R = (\# \text{ of operations in parallel version}) / (\# \text{ operations in best single processor algorithm version})$
- Efficiency
 - $E = (\text{Time with 1 processor}) / (\text{processors} \times \text{Time with } P \text{ processors})$
 - $E = U/R$

Utilization of a Multiprocessor

Multiprocessor metrics

Utilization : How much processing capability we use



$$U = \frac{10 \text{ operations (in parallel version)}}{3 \text{ processors} \times 5 \text{ time units}} = \frac{10}{15}$$

$$U = \frac{\text{Ops with } p \text{ proc.}}{p \times T_p}$$

Redundancy: How much extra work due to multiprocessing

$$R = \frac{\text{Ops with } p \text{ proc.}^{\text{best}}}{\text{Ops with 1 proc.}^{\text{best}}} = \frac{10}{8}$$

R is always ≥ 1

Efficiency: How much resource we use compared to how much resource we can get away with

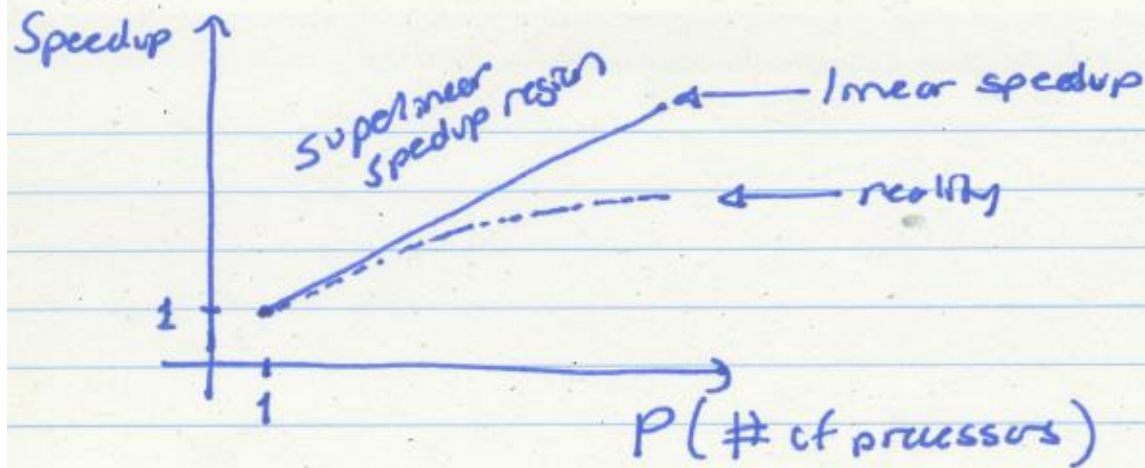
$$E = \frac{1 \cdot T_1^{\text{best}}}{p \cdot T_p^{\text{best}}} \quad \begin{array}{l} \text{(tying up 1 proc for } T_p \text{ time units)} \\ \text{(tying up } p \text{ proc. for } T_p \text{ time units)} \end{array}$$

$$= \frac{8}{15}$$

$$\left(E = \frac{U}{R} \right)$$

Amdahl's Law and Caveats of Parallelism

Caveats of Parallelism (I)



Why the reality? (diminishing returns)

$$T_p = \alpha \cdot \frac{T_1}{p} + (1-\alpha) \cdot T_1$$

┌
└
↓
parallelizable part/fraction
of the single-processor
program

┌
└
└
non-parallelizable part

Amdahl's Law

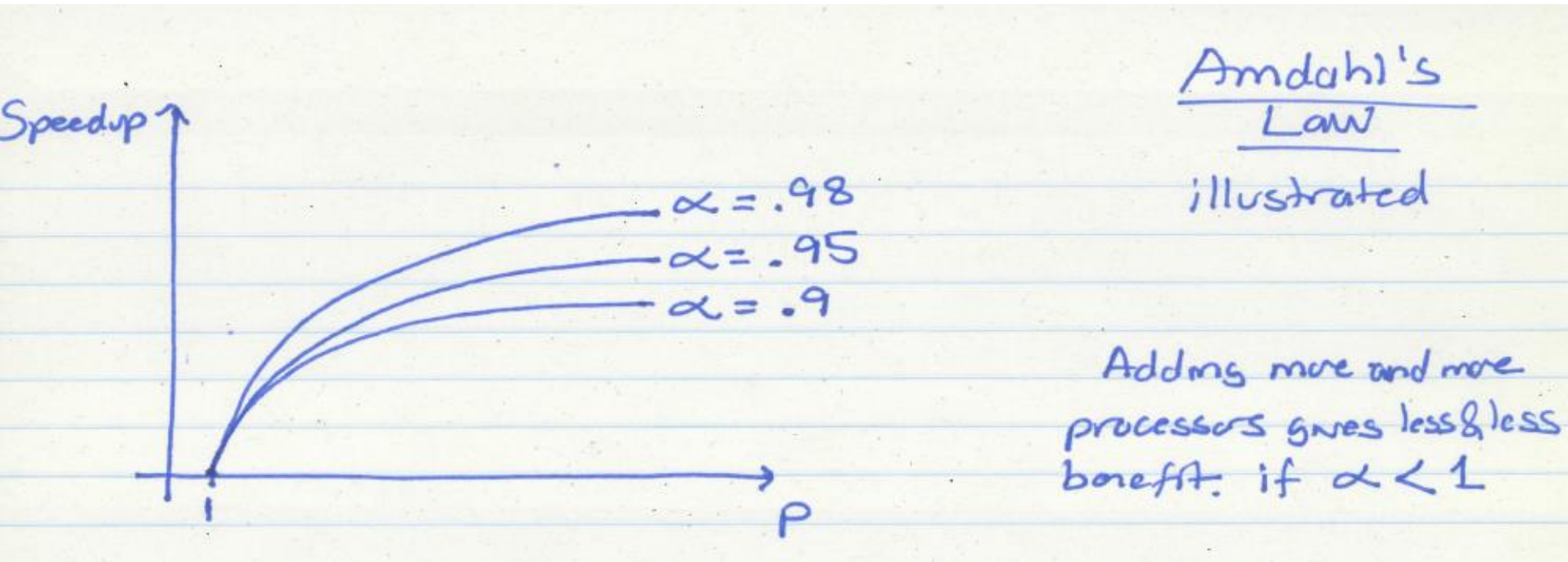
$$\text{Speedup}_{\text{with } p \text{ proc.}} = \frac{T_1}{T_p} = \frac{1}{\frac{\alpha}{p} + (1-\alpha)}$$

$$\text{Speedup}_{\text{as } p \rightarrow \infty} = \frac{1}{1 - \alpha}$$

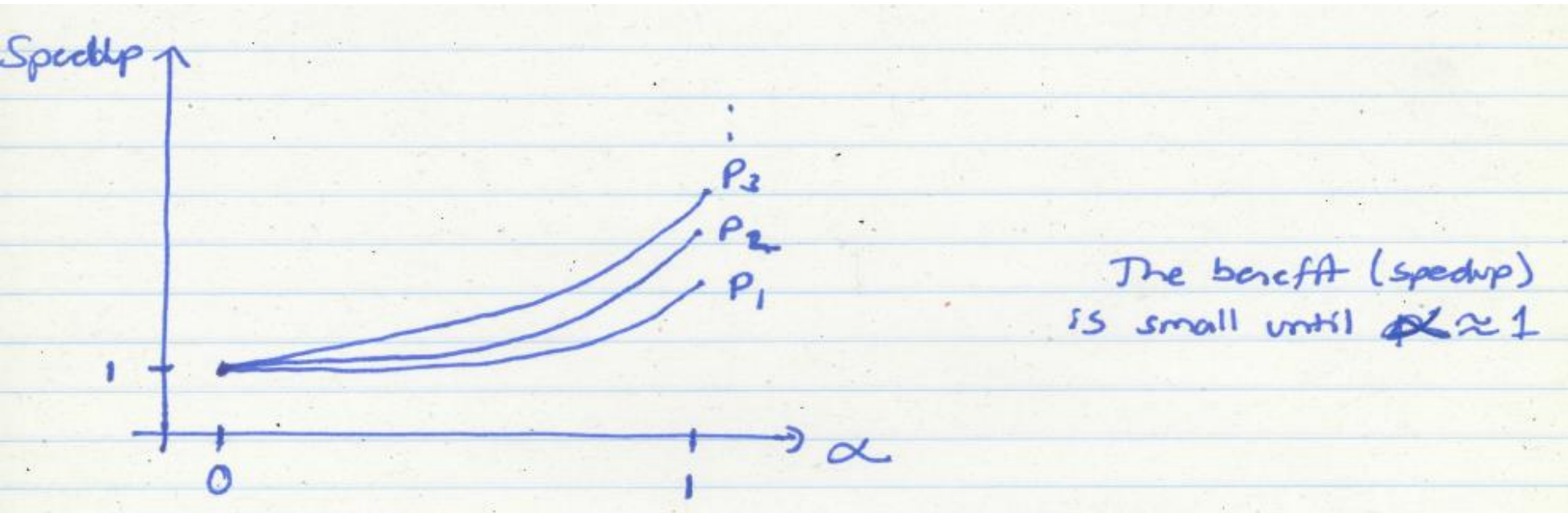
α → bottleneck for parallel Speedup

Amdahl, “Validity of the single processor approach to achieving large scale computing capabilities,” AFIPS 1967.

Amdahl's Law Implication 1



Amdahl's Law Implication 2



Caveats of Parallelism (II)

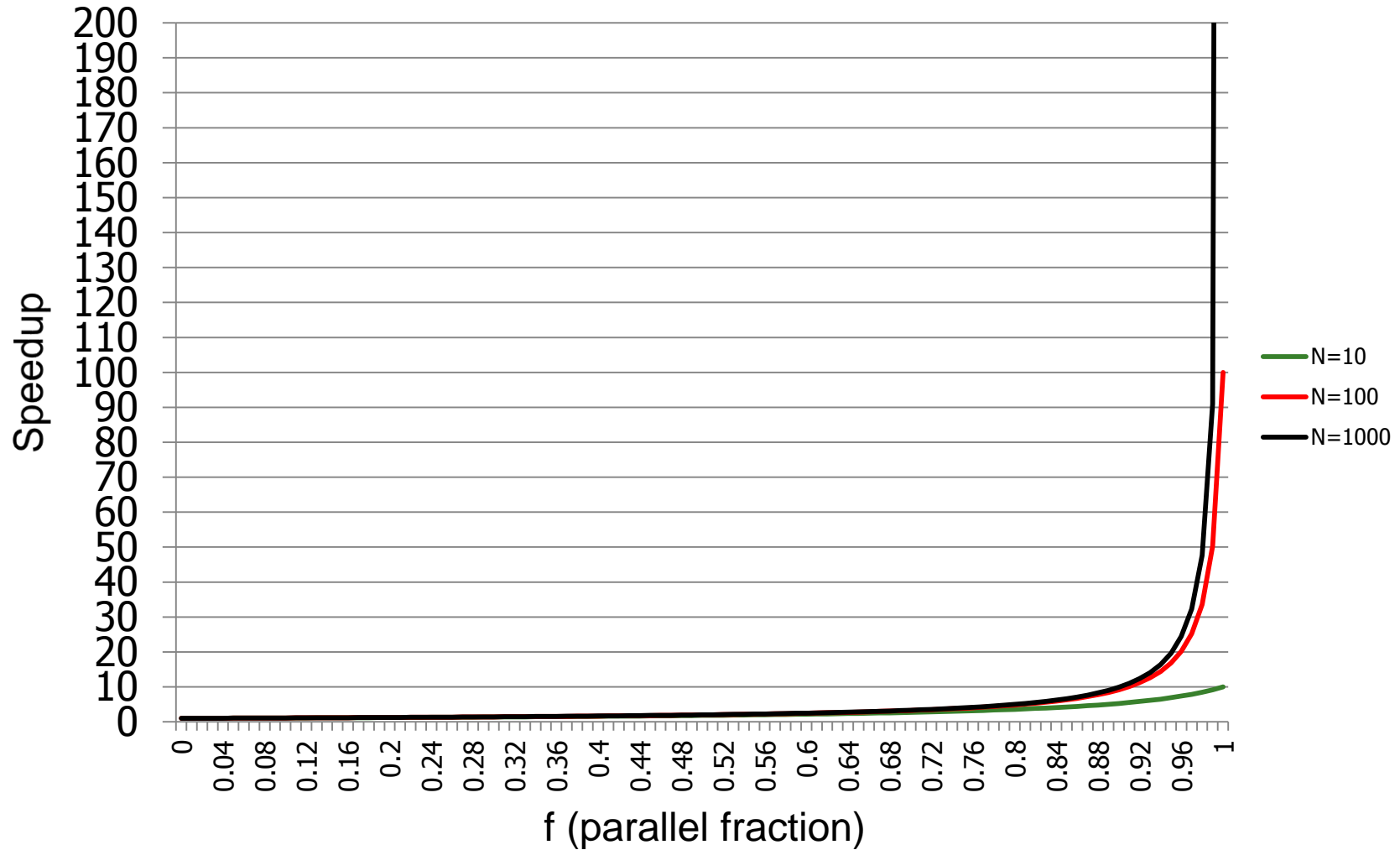
■ Amdahl's Law

- f: Parallelizable fraction of a program
- N: Number of processors

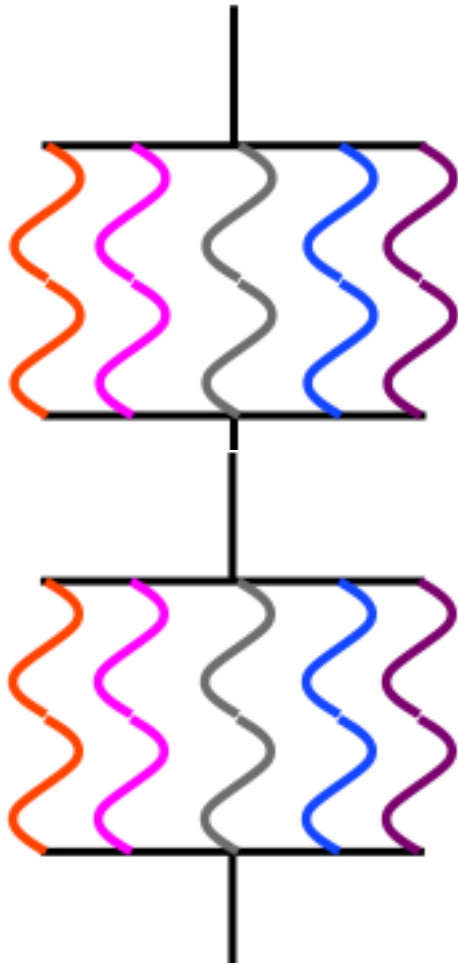
$$\text{Speedup} = \frac{1}{1 - f + \frac{f}{N}}$$

- Amdahl, “Validity of the single processor approach to achieving large scale computing capabilities,” AFIPS 1967.
- **Maximum speedup limited by serial portion: Serial bottleneck**
- **Parallel portion is usually not perfectly parallel**
 - Synchronization overhead (e.g., updates to shared data)
 - Load imbalance overhead (imperfect parallelization)
 - Resource sharing overhead (contention among N processors)

Sequential Bottleneck



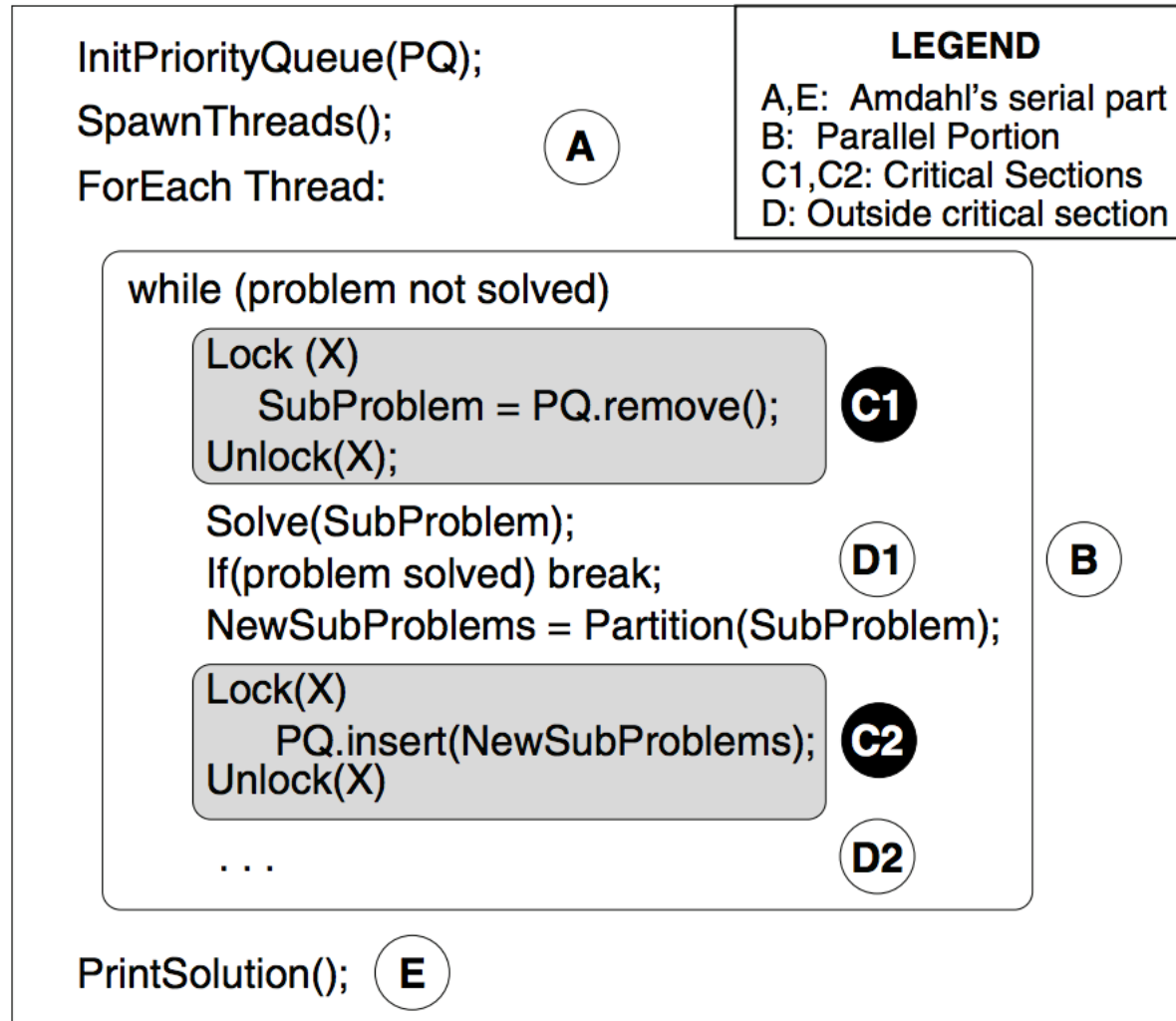
Why the Sequential Bottleneck?



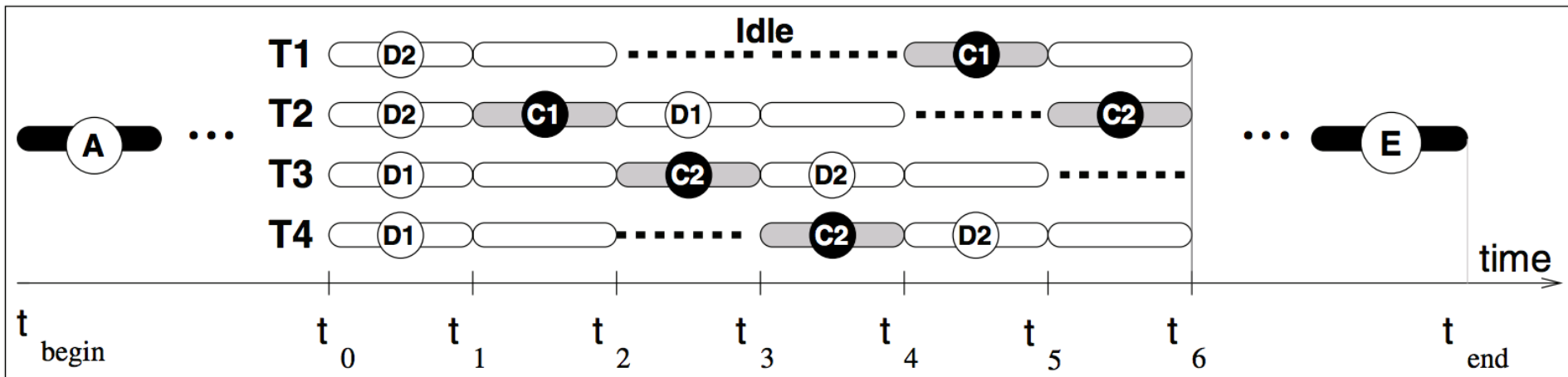
- Parallel machines have the sequential bottleneck
- Main cause: **Non-parallelizable operations on data** (e.g. non-parallelizable loops)

```
for ( i = 0 ; i < N; i++)  
    A[i] = (A[i] + A[i-1]) / 2
```
- There are other causes as well:
 - Single thread prepares data and spawns parallel tasks (usually sequential)

Another Example of Sequential Bottleneck (I)



Another Example of Sequential Bottleneck (II)



Bottlenecks in Parallel Portion

- **Synchronization:** Operations manipulating shared data cannot be parallelized
 - Locks, mutual exclusion, barrier synchronization
 - **Communication:** Tasks may need values from each other
 - Causes thread serialization when shared data is contended
- **Load Imbalance:** Parallel tasks may have different lengths
 - Due to imperfect parallelization or microarchitectural effects
 - Reduces speedup in parallel portion
- **Resource Contention:** Parallel tasks can share hardware resources, delaying each other
 - Replicating all resources (e.g., memory) expensive
 - Additional latency not present when each task runs alone

Bottlenecks in Parallel Portion: Another View

- Threads in a multi-threaded application can be inter-dependent
 - As opposed to threads from different applications
- Such threads can synchronize with each other
 - Locks, barriers, pipeline stages, condition variables, semaphores, ...
- Some threads can be on the critical path of execution due to synchronization; some threads are not
- Even within a thread, some “code segments” may be on the critical path of execution; some are not

Remember: Critical Sections

- Enforce mutually exclusive access to shared data
- Only one thread can be executing it at a time
- Contended critical sections make threads wait → threads causing serialization can be on the critical path

Each thread:

```
loop {
```

```
  Compute
```

N

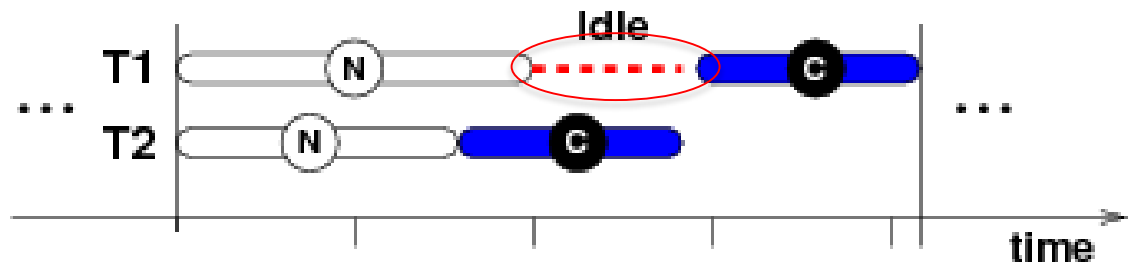
```
  lock(A)
```

```
    Update shared data
```

```
  unlock(A)
```

C

```
}
```

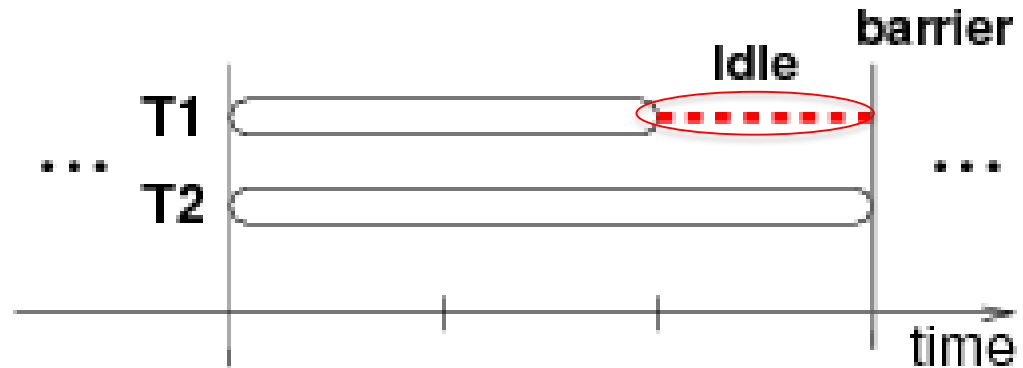


Remember: Barriers

- Synchronization point
- Threads have to wait until all threads reach the barrier
- Last thread arriving to the barrier is on the critical path

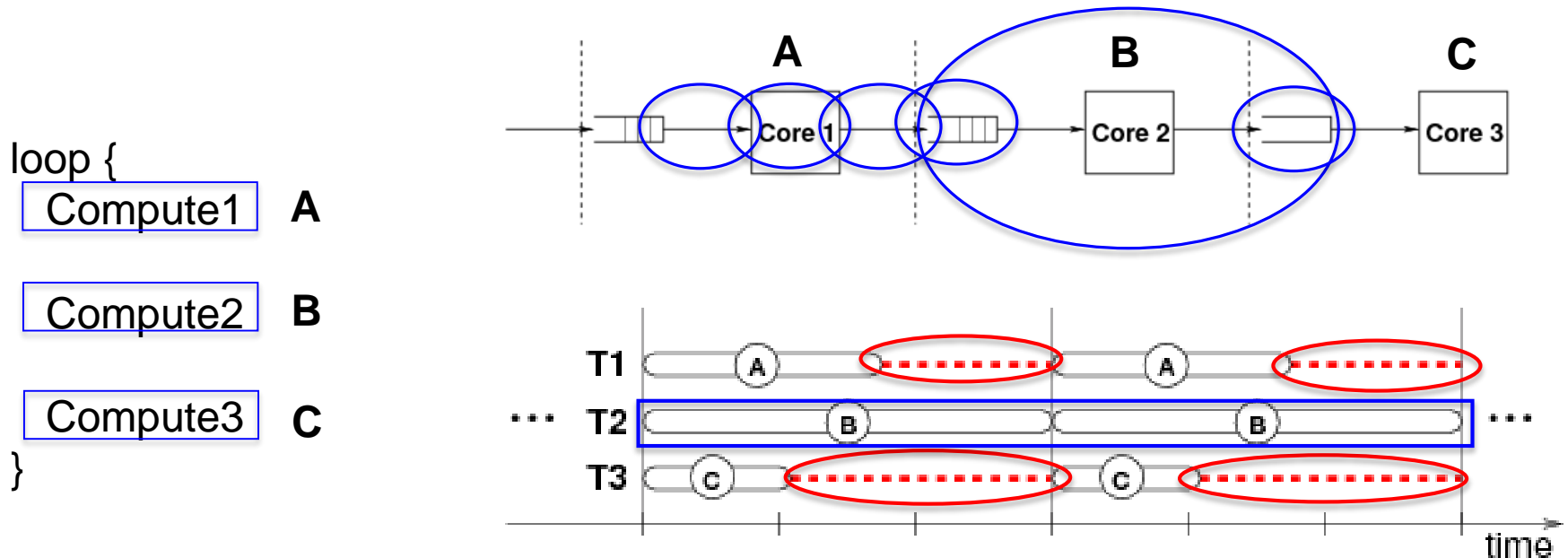
Each thread:

```
loop1 {  
    Compute  
}  
barrier  
loop2 {  
    Compute  
}
```



Remember: Stages of Pipelined Programs

- Loop iterations are statically divided into code segments called *stages*
- Threads execute stages on different cores
- Thread executing the slowest stage is on the critical path



Difficulty in Parallel Programming

- Little difficulty if parallelism is natural
 - “Embarrassingly parallel” applications
 - Multimedia, physical simulation, graphics
 - Large web servers, databases?
- Difficulty is in
 - Getting parallel programs to work correctly
 - Optimizing performance in the presence of bottlenecks
- Much of **parallel computer architecture** is about
 - Designing machines that overcome the sequential and parallel bottlenecks to achieve higher performance and efficiency
 - Making programmer’s job easier in writing correct and high-performance parallel programs

Memory Ordering in Multiprocessors

Readings: Memory Consistency

■ Required

- ❑ Lamport, "How to Make a Multiprocessor Computer That Correctly Executes Multiprocess Programs," IEEE Transactions on Computers, 1979

■ Recommended

- ❑ Gharachorloo et al., "Memory Consistency and Event Ordering in Scalable Shared-Memory Multiprocessors," ISCA 1990.
- ❑ Gharachorloo et al., "Two Techniques to Enhance the Performance of Memory Consistency Models," ICPP 1991.
- ❑ Ceze et al., "BulkSC: bulk enforcement of sequential consistency," ISCA 2007.

Memory Consistency vs. Cache Coherence

- **Consistency** is about ordering of **all memory operations** from different processors (i.e., to different memory locations)
 - **Global ordering** of accesses to *all* memory *locations*
- **Coherence** is about ordering of **operations** from different processors **to the same memory location**
 - **Local ordering** of accesses to *each* cache *block*

Difficulties of Multiprocessing

- Much of **parallel computer architecture** is about
 - ❑ Designing machines that overcome the sequential and parallel bottlenecks to achieve higher performance and efficiency
 - ❑ Making programmer's job easier in writing correct and high-performance parallel programs

Ordering of Operations

- Operations: A, B, C, D
 - In what order should the hardware execute (and report the results of) these operations?
- A contract between programmer and microarchitect
 - Specified by the ISA
- Preserving an “expected” (more accurately, “agreed upon”) order simplifies programmer’s life
 - Ease of debugging; ease of state recovery, exception handling
- Preserving an “expected” order usually makes the hardware designer’s life difficult
 - Especially if the goal is to design a high performance processor: Recall load-store queues in out of order execution and their complexity

Memory Ordering in a Single Processor

- Specified by the von Neumann model
- Sequential order
 - Hardware **executes** the load and store operations **in the order specified by the sequential program**
- Out-of-order execution does not change the semantics
 - Hardware **retires (reports to software the results of)** the load and store operations **in the order specified by the sequential program**
- Advantages: 1) Architectural state is precise within an execution.
2) Architectural state is consistent across different runs of the program
→ Easier to debug programs
- Disadvantage: Preserving order adds overhead, reduces performance, increases complexity, reduces scalability

Memory Ordering in a Dataflow Processor

- A memory operation executes when its operands are ready
- Ordering specified only by data dependencies
- Two operations can be executed and retired in any order if they have no dependency
- Advantage: Lots of parallelism → high performance
- Disadvantages:
 - Precise state is very hard to maintain (No specified order)
→ Very hard to debug
 - Order can change across runs of the same program
→ Very hard to debug

Memory Ordering in a MIMD Processor

- Each processor's memory operations are in sequential order with respect to the "thread" running on that processor (assume each processor obeys the von Neumann model)
- Multiple processors execute memory operations concurrently
- How does the memory see the order of operations from all processors?
 - In other words, what is the ordering of operations across different processors?

Why Does This Even Matter?

- Ease of debugging

- It is nice to have the same execution done at different times to have the same order of execution → Repeatability

- Correctness

- Can we have incorrect execution if the order of memory operations is different from the point of view of different processors?

- Performance and overhead

- Enforcing a strict “sequential ordering” can make life harder for the hardware designer in implementing performance enhancement techniques (e.g., OoO execution, caches)

When Could Order Affect Correctness?

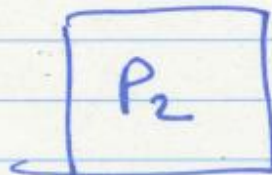
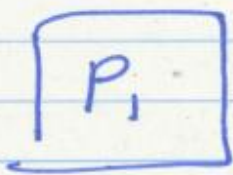
- When protecting shared data

Protecting Shared Data

- Threads are not allowed to update shared data concurrently
 - For correctness purposes
- Accesses to shared data are encapsulated inside *critical sections* or protected via *synchronization constructs* (locks, semaphores, condition variables)
- Only one thread can execute a critical section at a given time
 - Mutual exclusion principle
- A multiprocessor should provide the *correct* execution of synchronization primitives to enable the programmer to protect shared data

Supporting Mutual Exclusion

- Programmer needs to make sure mutual exclusion (synchronization) is correctly implemented
 - We will assume this
 - But, correct parallel programming is an important topic
 - Reading: Dijkstra, “[Cooperating Sequential Processes](http://www.cs.utexas.edu/users/EWD/transcriptions/EWD01xx/EWD123.html),” 1965.
 - <http://www.cs.utexas.edu/users/EWD/transcriptions/EWD01xx/EWD123.html>
 - See Dekker’s algorithm for mutual exclusion
- Programmer relies on hardware primitives to support correct synchronization
- If hardware primitives are not correct (or unpredictable), programmer’s life is tough
- If hardware primitives are correct but not easy to reason about or use, programmer’s life is still tough



Protecting Shared Data

$F_1 = \emptyset$



A $F_1 = 1$

B IF ($F_2 = \emptyset$) THEN
 {Critical section}
 $F_1 \neq \emptyset$

ELSE
 {...}

$F_2 = \emptyset$



X $F_2 = 1$

Y IF ($F_1 = \emptyset$) THEN
 {Critical section}

ELSE
 {...}

Only P_1 or
 P_2 should
be in this
section
at any given
time,
not both

Assume P_1 is in critical section.

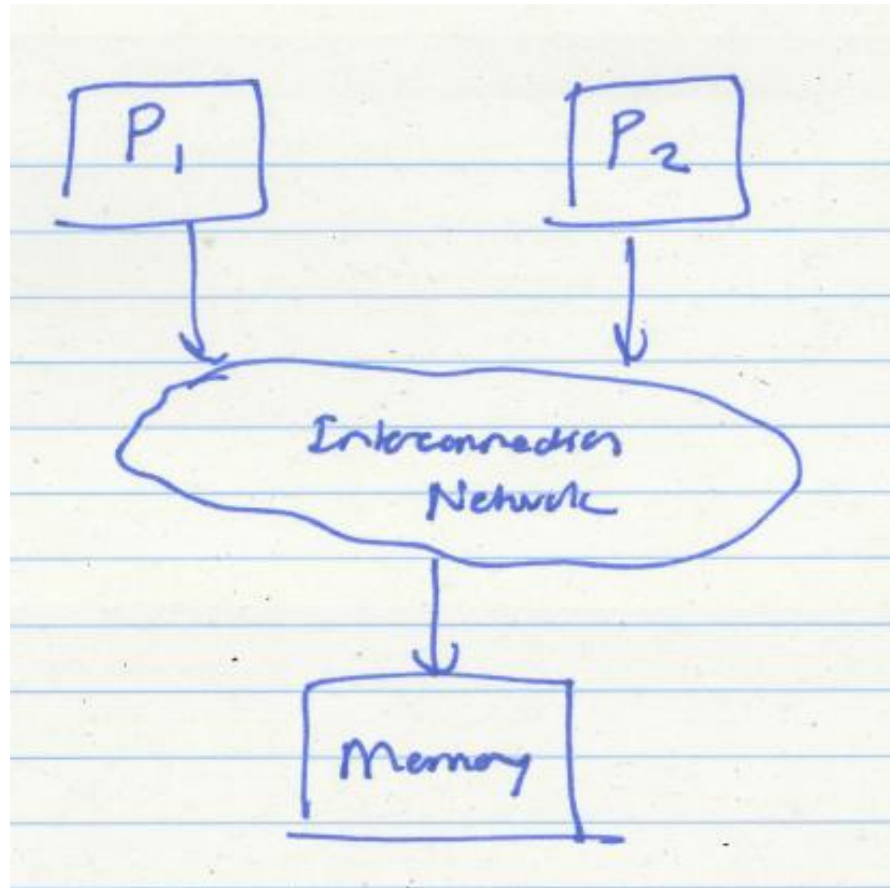
Intuitively, it must have executed A,

which means F_1 must be 1 (as A happens before B),

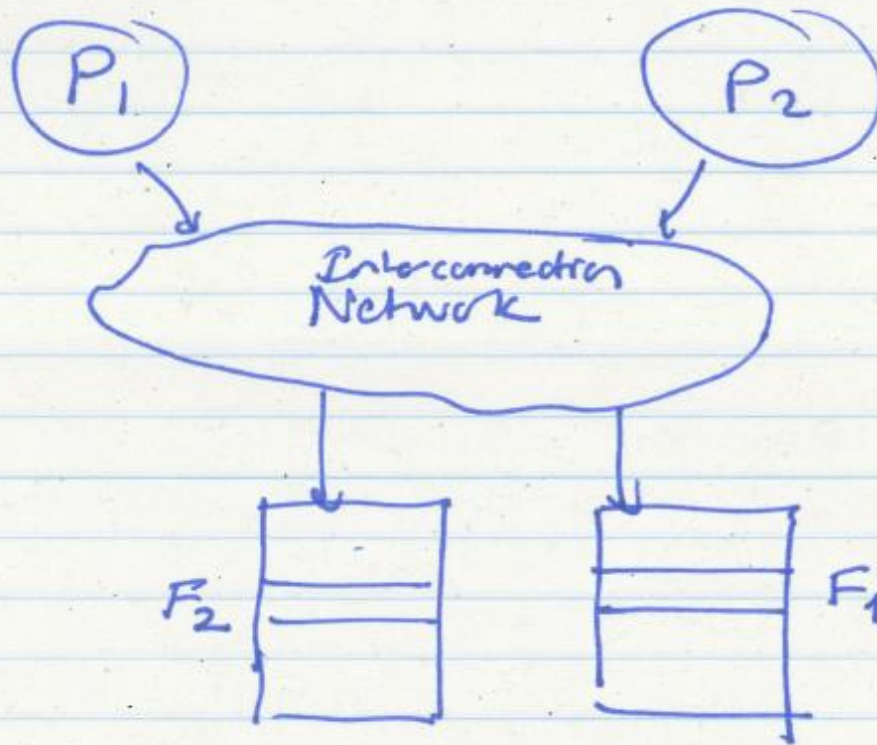
which means P_2 should not enter the critical section.

A Question

- Can the two processors be in the critical section at the same time given that they both obey the von Neumann model?
- Answer: yes



An Incorrect Result (due to an implementation that does not provide sequential consistency)



time 0: P_1 executes A
(set $F_1 = 1$) $S + F_1$ complete
A is sent to memory (from P_1 's view)

P_2 executes X
(set $F_2 = 1$) $S + F_2$ complete
X is sent to memory (from P_2 's view)

Both Processors in Critical Section

time 0: P_1 executes A
(set $F_1 = 1$) st F_1 complete (from P_1 's view)
A is sent to memory

P_2 executes X
(set $F_2 = 1$) st F_2 complete (from P_2 's view)
X is sent to memory

time 1: P_1 executes B
(test $F_2 == 0$) ld F_2 started
B is sent to memory

P_2 executes Y
(test $F_1 == 0$) ld F_1 started
Y is sent to memory

time 50: Memory sends back to P_1
 $F_2 (0)$ ld F_2 complete

Memory sends back to P_2
 $(F_1 \neq 0)$ ld F_1 complete

time 51: P_1 is in critical section
~~execute~~

P_2 is in critical section

time 100: Memory completes A
 $F_1 = 1$ in memory
(too late!)

Memory completes ~~X~~
 $F_2 = 1$ in memory
(too late!)

What happened?

P₁'s view of mem. ops

A	(F ₁ =1)
B	(test F ₂ =0)
X	(F ₂ =1)

P₂'s view

X	(F ₂ =1)
Y	(test ₂ F ₁ =0)
A	(F ₁ =1)

A appeared to happen
before X

X appeared to happen
before A



Problem!

These two processors did
not see the same order
of operations in memory

The Problem

- The two processors did **NOT** see the same order of operations to memory
- The “happened before” relationship between multiple updates to memory was inconsistent between the two processors’ points of view
- As a result, each processor thought the other was **not** in the critical section

How Can We Solve The Problem?

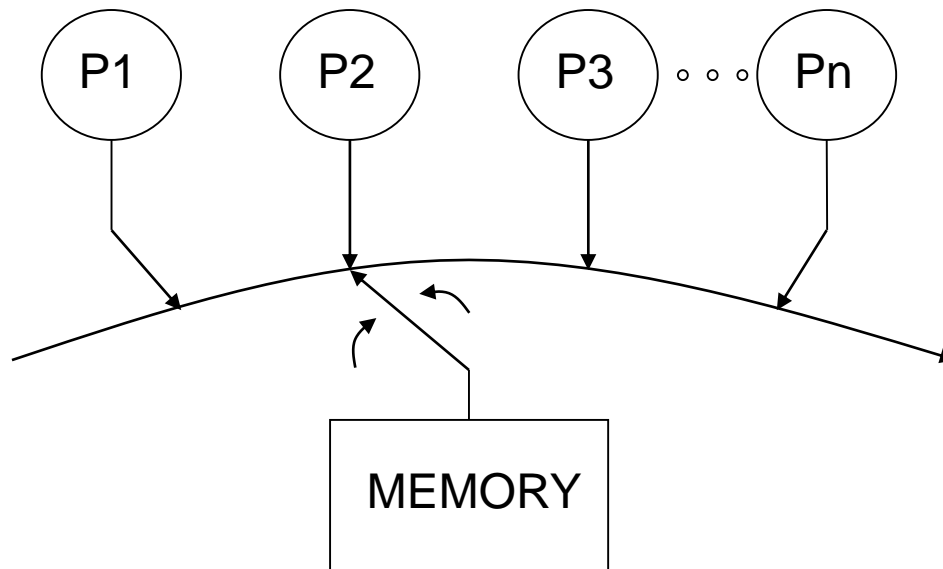
- Idea: Sequential consistency
- All processors see the same order of operations to memory
- i.e., all memory operations happen in an order (called the global total order) that is consistent across all processors
- Assumption: within this global order, each processor's operations appear in sequential order with respect to its own operations.

Sequential Consistency

- Lamport, "How to Make a Multiprocessor Computer That Correctly Executes Multiprocess Programs," IEEE Transactions on Computers, 1979
 - A multiprocessor system is sequentially consistent if:
 - the result of any execution is the same as if the operations of all the processors were executed in some sequential order
- AND
- the operations of each individual processor appear in this sequence in the order specified by its program
 - This is a memory ordering model, or memory model
 - Specified by the ISA

Programmer's Abstraction

- Memory is a switch that services one load or store at a time from any processor
- All processors see the currently serviced load or store at the same time
- Each processor's operations are serviced in program order



Sequentially Consistent Operation Orders

- Potential correct global orders (all are correct):
 - A B X Y
 - A X B Y
 - A X Y B
 - X A B Y
 - X A Y B
 - X Y A B
- Which order (interleaving) is observed depends on implementation and dynamic latencies

Consequences of Sequential Consistency

■ Corollaries

1. Within the same execution, all processors see the same global order of operations to memory
 - No correctness issue
 - Satisfies the “happened before” intuition
2. Across different executions, different global orders can be observed (each of which is sequentially consistent)
 - Debugging is still difficult (as order changes across runs)

Issues with Sequential Consistency?

- Nice abstraction for programming, but two issues:
 - Too conservative ordering requirements
 - Limits the aggressiveness of performance enhancement techniques

- Is the total global order requirement too strong?
 - Do we need a global order across all operations and all processors?
 - How about a global order only across all stores?
 - Total store order memory model; unique store order model
 - How about enforcing a global order only at the boundaries of synchronization?
 - Relaxed memory models
 - Acquire-release consistency model

Issues with Sequential Consistency?

- Performance enhancement techniques that could make SC implementation difficult
- Out-of-order execution
 - Loads happen out-of-order with respect to each other and with respect to independent stores → makes it difficult for all processors to see the same global order of all memory operations
- Caching
 - A memory location is now present in multiple places
 - Prevents the effect of a store to be seen by other processors → makes it difficult for all processors to see the same global order of all memory operations

Weaker Memory Consistency

- The ordering of operations is important when the order affects operations on shared data → i.e., when processors need to synchronize to execute a “program region”
- Weak consistency
 - Idea: Programmer specifies regions in which memory operations do not need to be ordered
 - “Memory fence” instructions delineate those regions
 - All memory operations before a fence must complete before fence is executed
 - All memory operations after the fence must wait for the fence to complete
 - Fences complete in program order
 - All synchronization operations act like a fence

Tradeoffs: Weaker Consistency

■ Advantage

- No need to guarantee a very strict order of memory operations
 - Enables the hardware implementation of performance enhancement techniques to be **simpler**
 - Can be **higher performance** than stricter ordering

■ Disadvantage

- More **burden on the programmer** or software (need to get the “fences” correct)

■ Another example of the programmer-microarchitect tradeoff

Example Question (I)

■ Question 4 in

□ <http://www.ece.cmu.edu/~ece447/s13/lib/exe/fetch.php?media=final.pdf>

4. Sequential Consistency [30 points]

Two threads (A and B) are concurrently running on a dual-core processor that implements a *sequentially consistent* memory model. Assume that the value at address 0x1000 is initialized to 0.

Thread A

X1: st 0x1, (0x1000)
X2: ld \$r1, (0x1000)
X3: st 0x2, (0x1000)
X4: ld \$r2, (0x1000)

Thread B

Y1: st 0x3, (0x1000)
Y2: ld \$r3, (0x1000)
Y3: st 0x4, (0x1000)
Y4: ld \$r4, (0x1000)

(a) List all possible values that can be stored in \$r3 after both threads have finished executing.

Example Question (II)

- (b) After both threads have finished executing, you find that $(\$r1, \$r2, \$r3, \$r4) = (1, 2, 3, 4)$. How many different *instruction interleavings* of the two threads produce this result?

- (c) What is the total number of all possible instruction interleavings? You need not expand factorials.

- (d) On a *non-sequentially consistent* processor, is the total number of all possible instruction interleavings less than, equal to, or greater than your answer to question (c)?

Caching in Multiprocessors

- Caching not only complicates ordering of **all operations**...
 - A memory location can be present in multiple caches
 - Prevents the effect of a store or load to be seen by other processors → **makes it difficult for all processors to see the same global order of (all) *memory operations***
- ... but it also complicates ordering of **operations on a single memory location**
 - A single memory location can be present in multiple caches
 - **Makes it difficult for processors that have cached the same location to have the correct value of that location (in the presence of updates to that location)**

Cache Coherence

Readings: Cache Coherence

■ Required

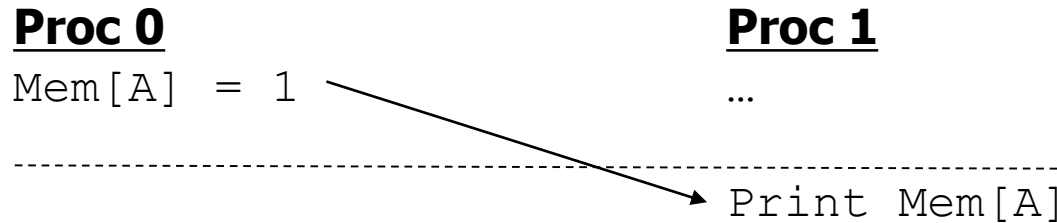
- Culler and Singh, *Parallel Computer Architecture*
 - Chapter 5.1 (pp 269 – 283), Chapter 5.3 (pp 291 – 305)
- P&H, *Computer Organization and Design*
 - Chapter 5.8 (pp 534 – 538 in 4th and 4th revised eds.)
- Papamarcos and Patel, “A low-overhead coherence solution for multiprocessors with private cache memories,” ISCA 1984.

■ Recommended

- Censier and Feautrier, “A new solution to coherence problems in multicache systems,” IEEE Trans. Computers, 1978.
- Goodman, “Using cache memory to reduce processor-memory traffic,” ISCA 1983.
- Laudon and Lenoski, “The SGI Origin: a ccNUMA highly scalable server,” ISCA 1997.
- Martin et al, “Token coherence: decoupling performance and correctness,” ISCA 2003.
- Baer and Wang, “On the inclusion properties for multi-level cache hierarchies,” ISCA 1988.

Shared Memory Model

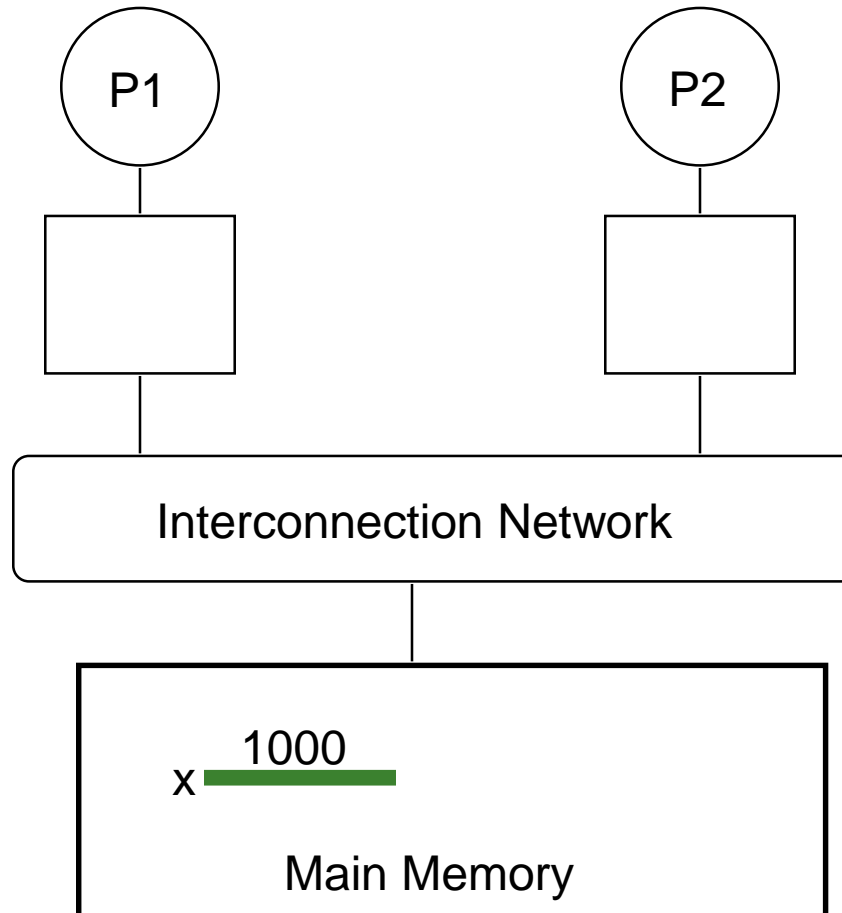
- Many parallel programs communicate through *shared memory*
- Proc 0 writes to an address, followed by Proc 1 reading
 - This implies communication between the two



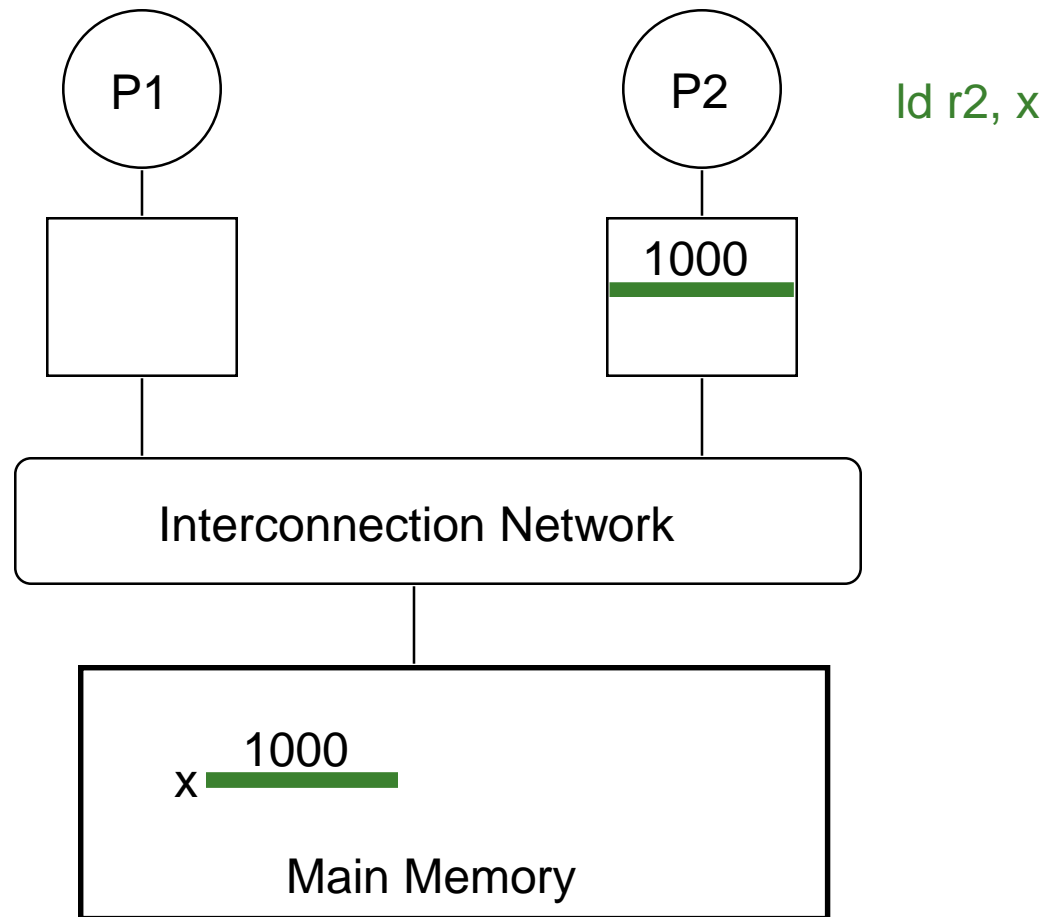
- Each read should receive the value last written by anyone
 - This requires synchronization (what does last written mean?)
- What if Mem[A] is cached (at either end)?

Cache Coherence

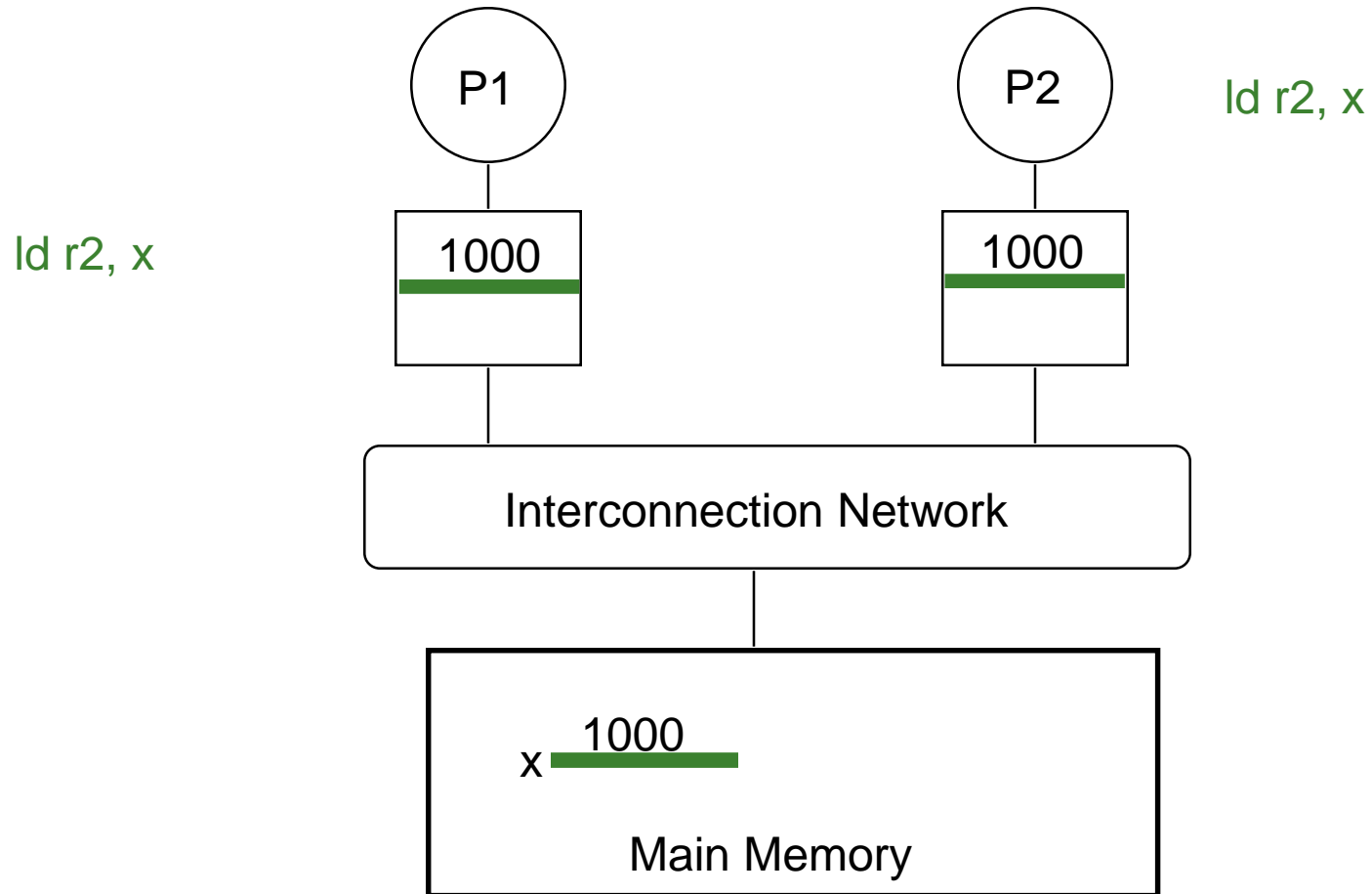
- Basic question: If multiple processors cache the same block, how do they ensure they all see a consistent state?



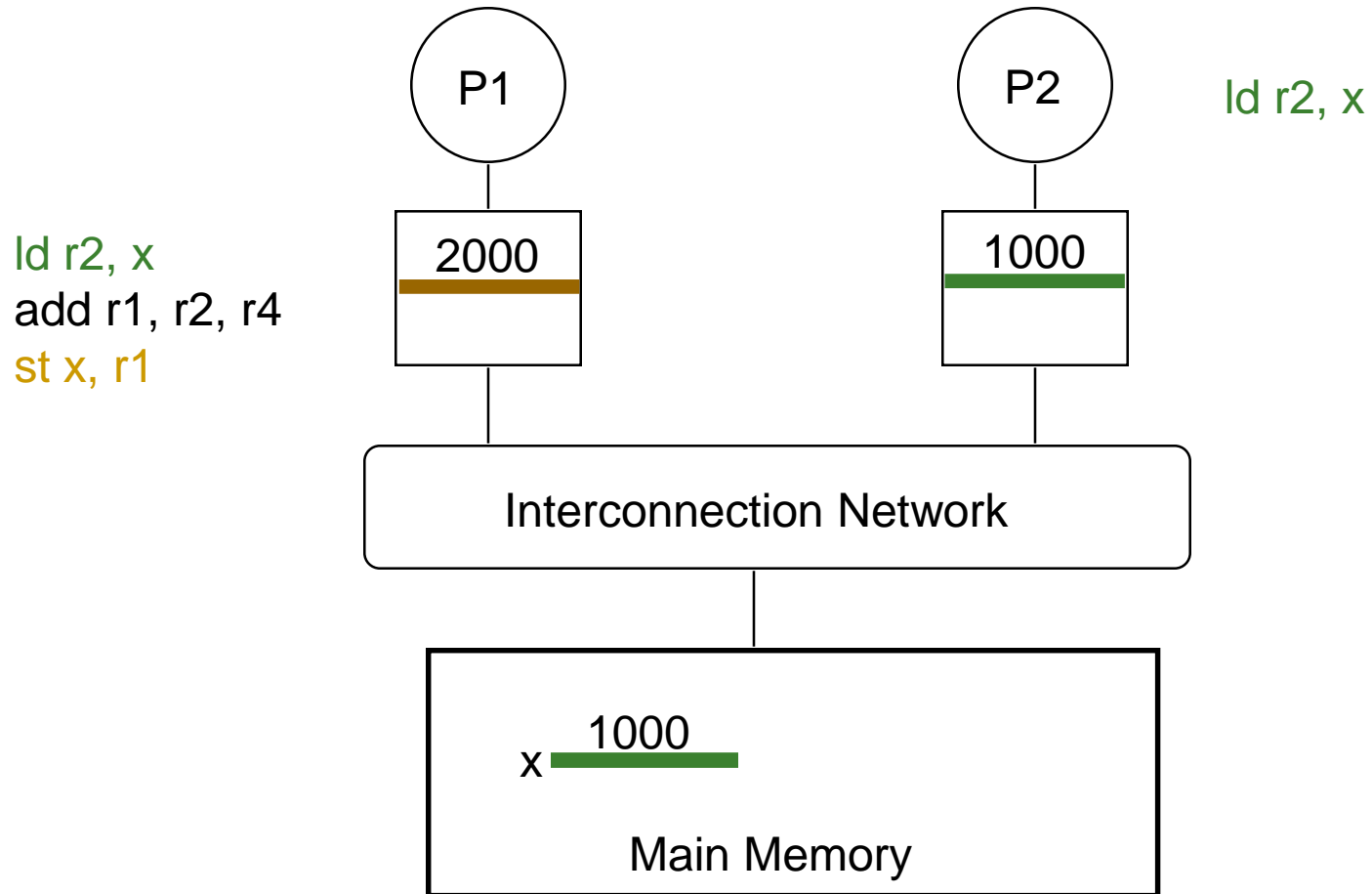
The Cache Coherence Problem



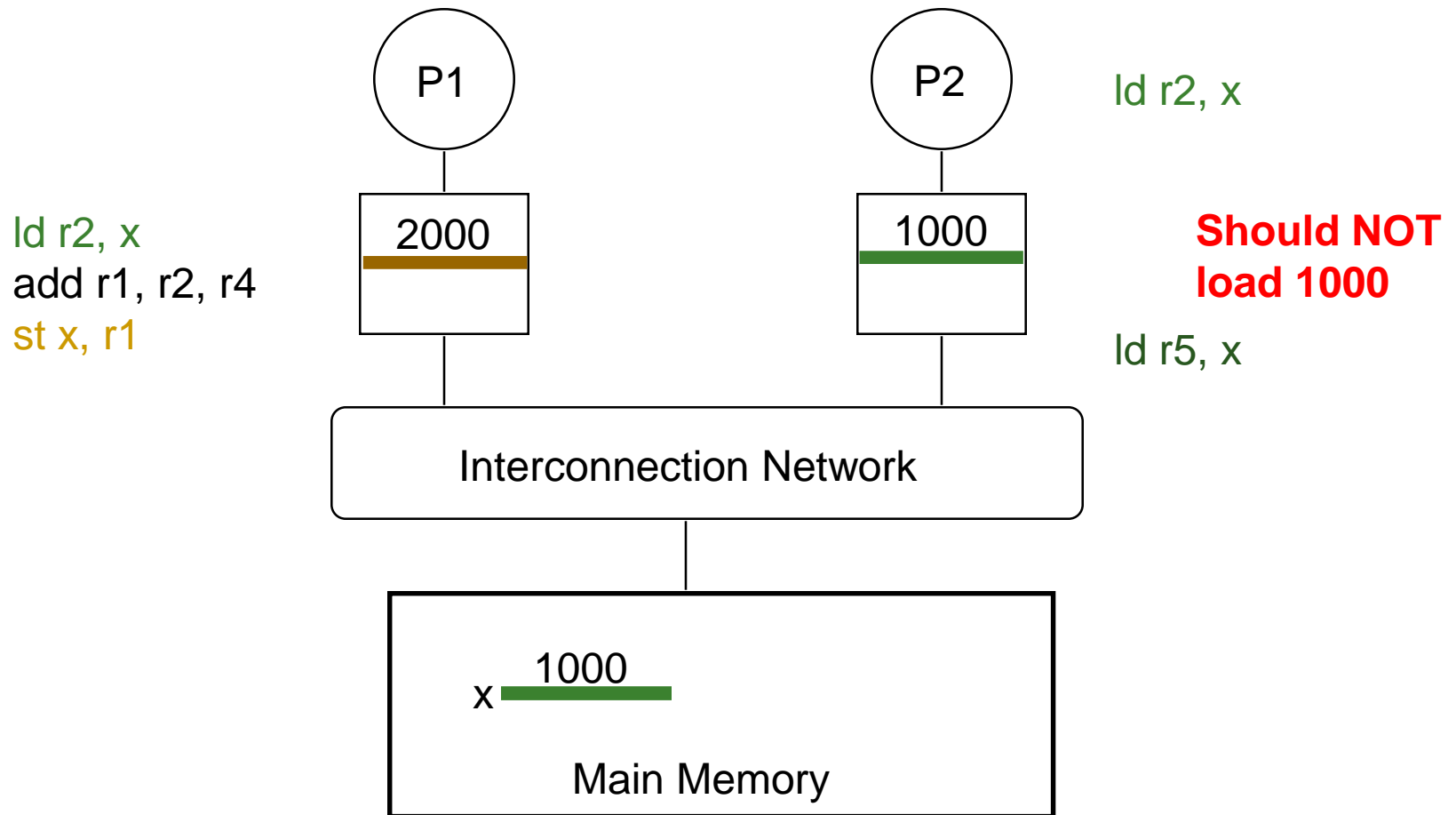
The Cache Coherence Problem



The Cache Coherence Problem



The Cache Coherence Problem



Cache Coherence: Whose Responsibility?

■ Software

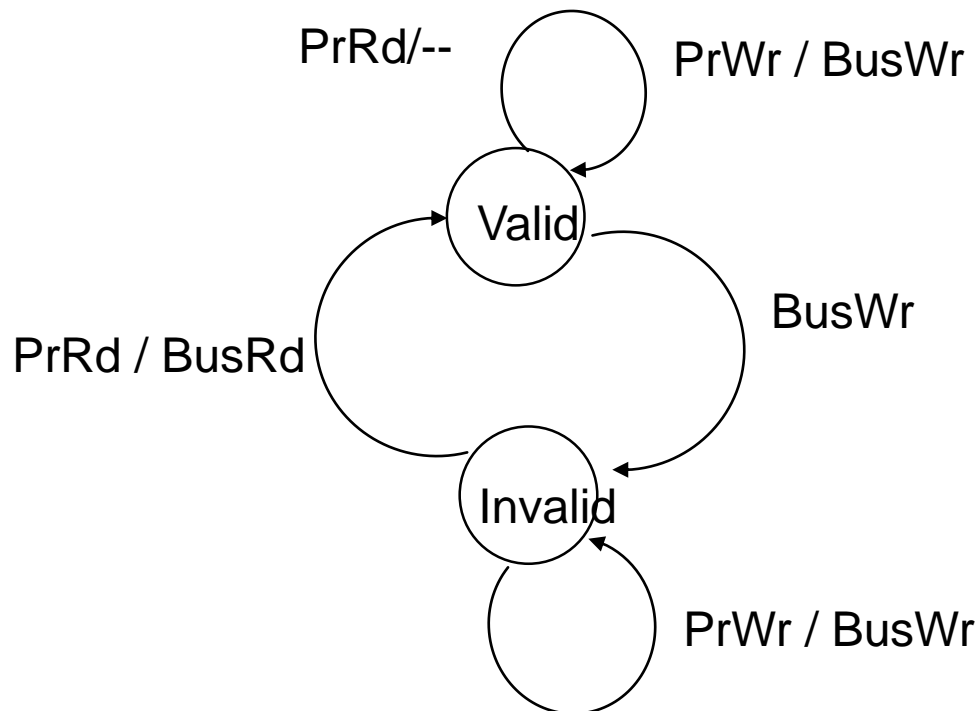
- ❑ Can the programmer ensure coherence if caches are invisible to software?
- ❑ What if the ISA provided a cache flush instruction?
 - FLUSH-LOCAL A: Flushes/invalidates the cache block containing address A from a processor's local cache.
 - FLUSH-GLOBAL A: Flushes/invalidates the cache block containing address A from all other processors' caches.
 - FLUSH-CACHE X: Flushes/invalidates all blocks in cache X.

■ Hardware

- ❑ Simplifies software's job
- ❑ One idea: Invalidate all other copies of block A when a processor writes to it

A Very Simple Coherence Scheme (VI)

- Caches “snoop” (observe) each other’s write/read operations. If a processor writes to a block, all others invalidate the block.
- A simple protocol:



- Write-through, no-write-allocate cache
- Actions of the local processor on the cache block: **PrRd**, **PrWr**,
- Actions that are broadcast on the bus for the block: **BusRd**, **BusWr**

(Non-)Solutions to Cache Coherence

■ No hardware based coherence

- Keeping caches coherent is software's responsibility
- + Makes microarchitect's life easier
- Makes average programmer's life much harder
 - need to worry about hardware caches to maintain program correctness?
- Overhead in ensuring coherence in software (e.g., page protection and page-based software coherence)

■ All caches are shared between all processors

- + No need for coherence
- Shared cache becomes the bandwidth bottleneck
- Very hard to design a scalable system with low-latency cache access this way

Maintaining Coherence

- Need to guarantee that all processors see a consistent value (i.e., consistent updates) for the same memory location
- Writes to location A by P0 should be seen by P1 (eventually), and all writes to A should appear in some order
- Coherence needs to provide:
 - **Write propagation:** guarantee that updates will propagate
 - **Write serialization:** provide a consistent order seen by all processors for the same memory location
- Need a global point of serialization for this store ordering

Hardware Cache Coherence

- Basic idea:
 - A processor/cache broadcasts its write/update to a memory location to all other processors
 - Another cache that has the location either updates or invalidates its local copy

Coherence: Update vs. Invalidate

- How can we *safely update replicated data*?
 - Option 1 (Update protocol): push an update to all copies
 - Option 2 (Invalidate protocol): ensure there is only one copy (local), update it
- **On a Read:**
 - If local copy is Invalid, put out request
 - (If another node has a copy, it returns it, otherwise memory does)

Coherence: Update vs. Invalidate (II)

■ **On a Write:**

- ❑ Read block into cache as before

Update Protocol:

- ❑ Write to block, and simultaneously broadcast written data and address to sharers
- ❑ (Other nodes update the data in their caches if block is present)

Invalidate Protocol:

- ❑ Write to block, and simultaneously broadcast invalidation of address to sharers
- ❑ (Other nodes invalidate block in their caches if block is present)

Update vs. Invalidate Tradeoffs

- Which do we want?
 - Write frequency and sharing behavior are critical
- **Update**
 - + If sharer set is constant and updates are infrequent, avoids the cost of invalidate-reacquire (broadcast update pattern)
 - If data is rewritten without intervening reads by other cores, updates were useless
 - Write-through cache policy → bus becomes bottleneck
- **Invalidate**
 - + After invalidation broadcast, core has exclusive access rights
 - + Only cores that keep reading after each write retain a copy
 - If write contention is high, leads to ping-ponging (rapid invalidation-reacquire traffic from different processors)

Two Cache Coherence Methods

- ❑ How do we ensure that the proper caches are updated?
- ❑ **Snoopy Bus** [Goodman ISCA 1983, Papamarcos+ ISCA 1984]
 - Bus-based, *single point of serialization for all memory requests*
 - Processors observe other processors' actions
 - ❑ E.g.: P1 makes “read-exclusive” request for A on bus, P0 sees this and invalidates its own copy of A
- ❑ **Directory** [Censier and Feautrier, IEEE ToC 1978]
 - *Single point of serialization per block*, distributed among nodes
 - Processors make explicit requests for blocks
 - Directory tracks which caches have each block
 - Directory coordinates invalidation and updates
 - ❑ E.g.: P1 asks directory for exclusive copy, directory asks P0 to invalidate, waits for ACK, then responds to P1

Directory Based Cache Coherence

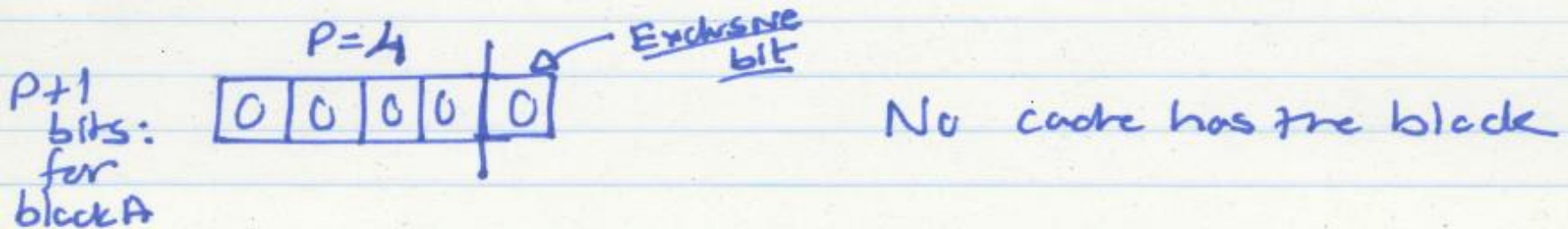
Directory Based Coherence

- Idea: A logically-central directory keeps track of where the copies of each cache block reside. Caches consult this directory to ensure coherence.

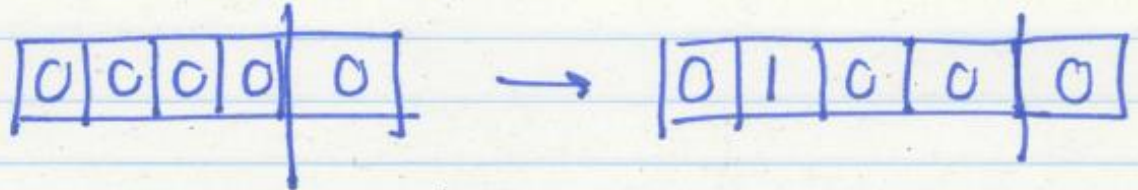
- An example mechanism:
 - For each cache block in memory, store $P+1$ bits in directory
 - One bit for each cache, indicating whether the block is in cache
 - Exclusive bit: indicates that a cache has the only copy of the block and can update it without notifying others
 - On a read: set the cache's bit and arrange the supply of data
 - On a write: invalidate all caches that have the block and reset their bits
 - Have an "exclusive bit" associated with each block in each cache (so that the cache can update the exclusive block silently)

Directory Based Coherence Example (I)

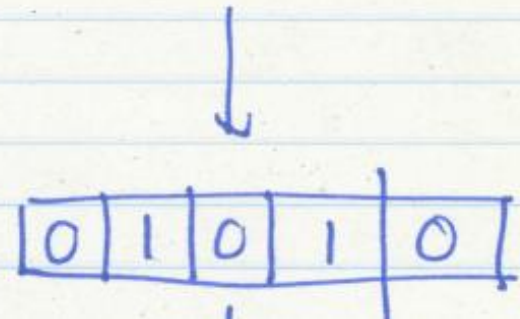
Example directory based scheme



① P_1 takes a read miss to block A



② P_3 takes a read miss



③ P₂ takes a write miss

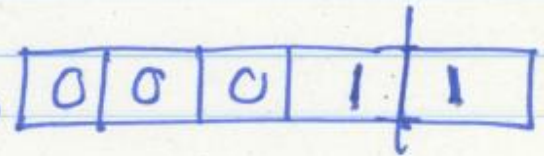
- invalidate P₁ & P₃'s caches
- write request → P₂ has the exclusive copy of the block now. Set the Exclusive bit



- P₂ can now update the block without notifying any other processor or the directory
- P₂ needs to have a bit in its cache indicating it can perform exclusive updates to that block
 - private/exclusive bit per cache block

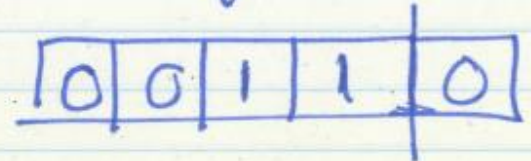
④ P₃ takes a write miss

- Mem ~~Controller~~ requests ~~the~~ block from P₂
- Mem Controller gives block to P₃
- P₂ invalidates its copy



⑤ P₂ takes a read miss

- P₃ supplies it



Directory Optimizations

- Directory is the coordinator for all actions to be performed on the block by any processor
 - Guarantees correctness, ordering
- Yet, there are many opportunities for optimization
 - Enabled by bypassing the directory and directly communicating between caches
 - We will see this later

We did not cover the following slides in lecture. They are for your benefit.

Computer Architecture

Lecture 19: Multiprocessors, Consistency, Coherence

Prof. Onur Mutlu

ETH Zürich

Fall 2017

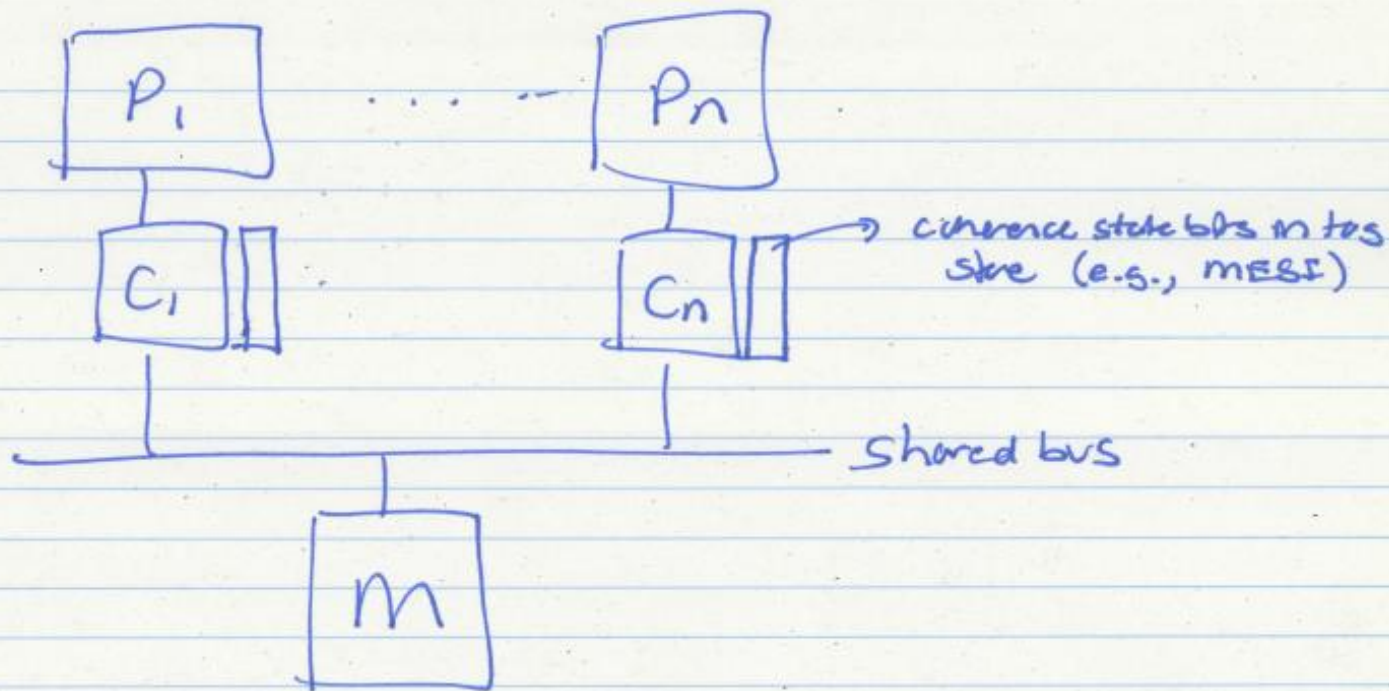
29 November 2017

Snoopy Cache Coherence

Snoopy Cache Coherence

- Idea:
 - ❑ All caches “snoop” all other caches’ read/write requests and keep the cache block coherent
 - ❑ Each cache block has “coherence metadata” associated with it in the tag store of each cache

- Easy to implement if all caches share a common bus
 - ❑ Each cache broadcasts its read/write operations on the bus
 - ❑ Good for small-scale multiprocessors
 - ❑ What if you would like to have a 1000-node multiprocessor?



SNOOPY CACHE

Each Cache observes its own processor & the bus
 - Changes the state of the cached block based on observed actions by processor & the bus

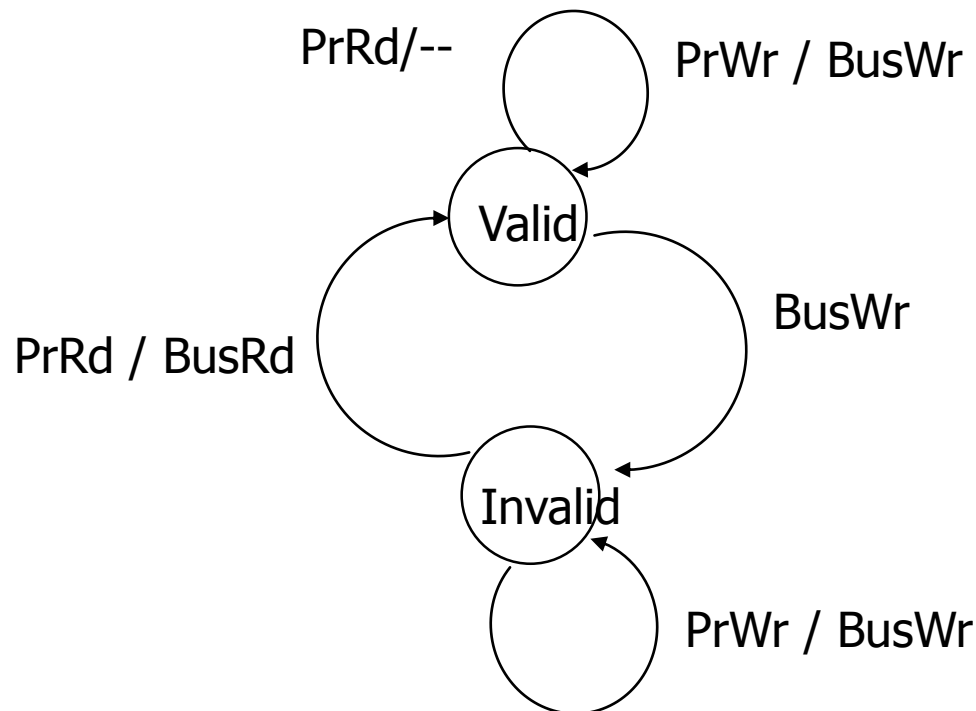
Processor actions to a block: PR (Proc. Read)
 RW (Proc. Write)

Bus actions to a block : BR (Bus Read)
 (comming from another processor) BW (Bus Write)

or BRx (Bus Read Exclusive)

A Simple Snoopy Cache Coherence Protocol

- Caches “snoop” (observe) each others’ write/read operations
- A simple protocol (VI protocol):



- **Write-through**, no-write-allocate cache
- Actions of the local processor on the cache block: **PrRd**, **PrWr**,
- Actions that are broadcast on the bus for the block: **BusRd**, **BusWr**

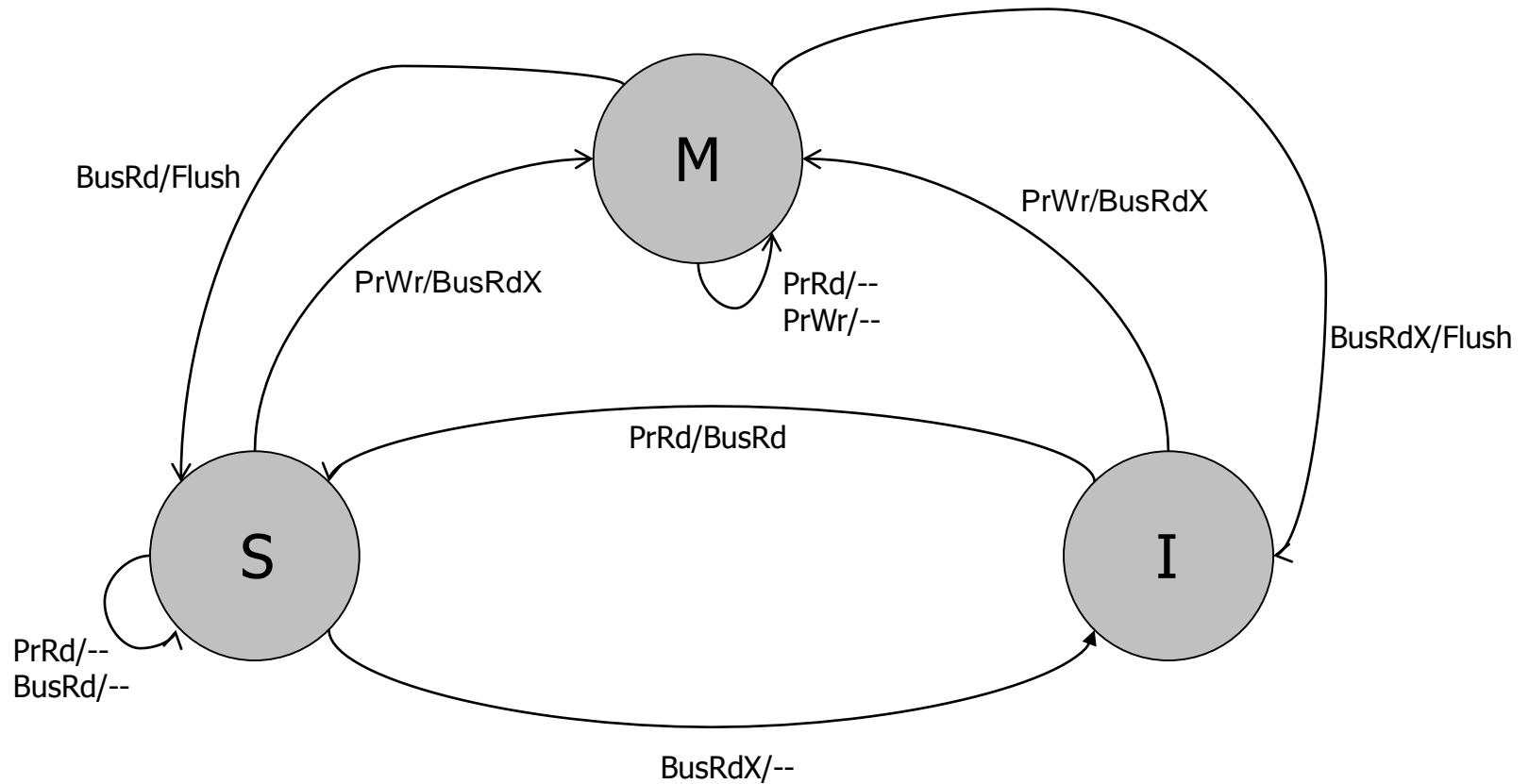
Extending the Protocol

- What if you want write-back caches?
 - We want a “modified” state

A More Sophisticated Protocol: MSI

- Extend metadata per block to encode three states:
 - **M**(odified): cache line is the only cached copy and is dirty
 - **S**(hared): cache line is potentially one of several cached copies
 - **I**(nvalid): cache line is not present in this cache
- Read miss makes a *Read* request on bus, transitions to **S**
- Write miss makes a *ReadEx* request, transitions to **M** state
- When a processor snoops *ReadEx* from another writer, it must invalidate its own copy (if any)
- $S \rightarrow M$ *upgrade* can be made without re-reading data from memory (via *Invalidations*)

MSI State Machine



ObservedEvent/Action

[Culler/Singh96]

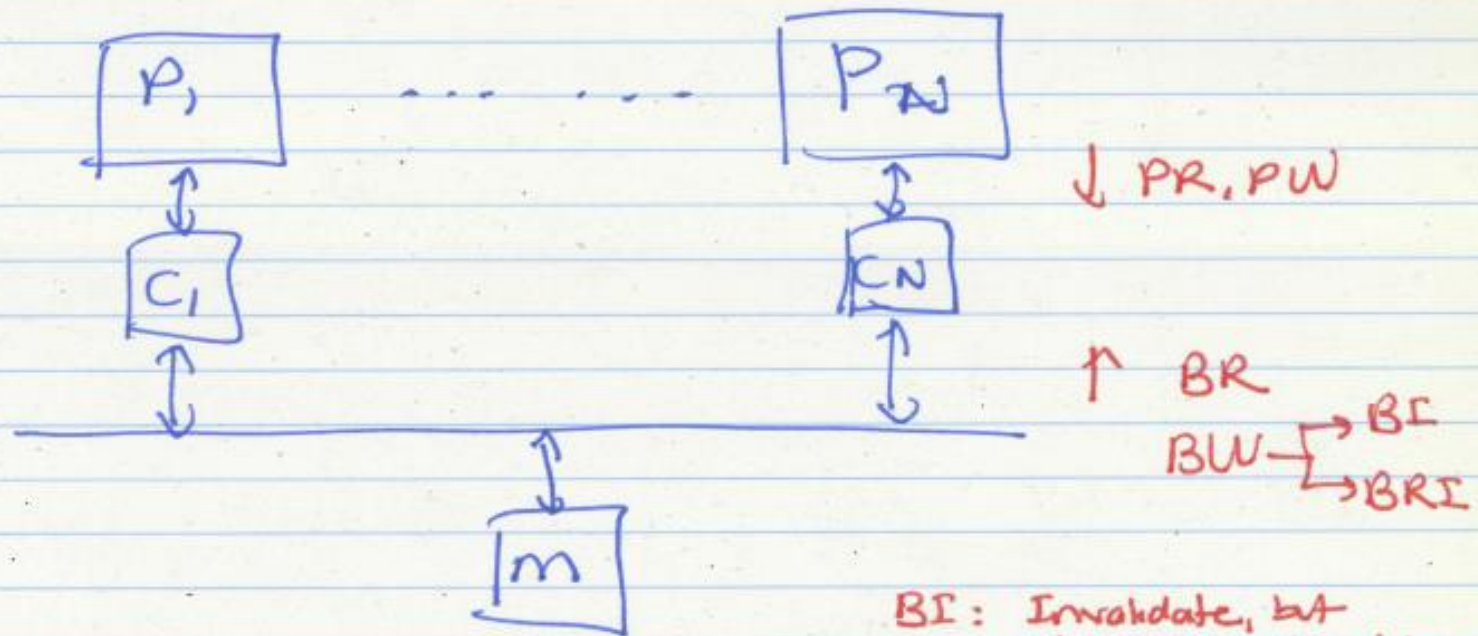
The Problem with MSI

- A block is in no cache to begin with
- Problem: On a read, the block immediately goes to “Shared” state although it may be the only copy to be cached (i.e., no other processor will cache it)
- Why is this a problem?
 - Suppose the cache that read the block wants to write to it at some point
 - It needs to broadcast “invalidate” even though it has the only cached copy!
 - *If the cache knew it had the only cached copy in the system, it could have written to the block without notifying any other cache → saves unnecessary broadcasts of invalidations*

The Solution: MESI

- Idea: Add another state indicating that this is the only cached copy and it is clean.
 - *Exclusive* state
- Block is placed into the *exclusive* state if, during *BusRd*, no other cache had it
 - Wired-OR “shared” signal on bus can determine this: snooping caches assert the signal if they also have a copy
- Silent transition *Exclusive* → *Modified* is possible on write!
- MESI is also called the *Illinois protocol*
 - Papamarcos and Patel, “A low-overhead coherence solution for multiprocessors with private cache memories,” ISCA 1984.

Illinois Protocol



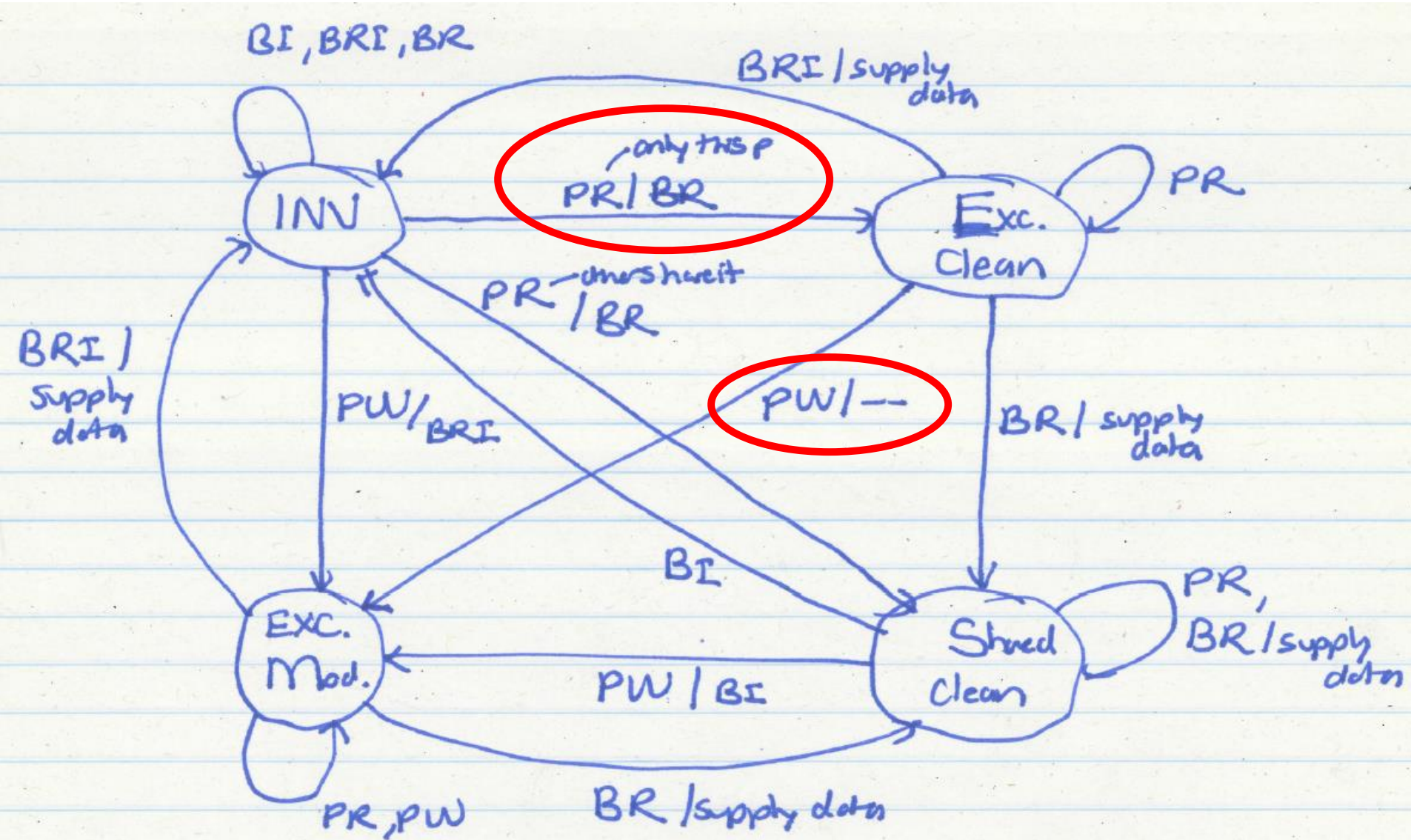
BI: Invalidate, but
already have the data
(do not supply it)

BRI: Invalidate, but
also need the data
(supply it)

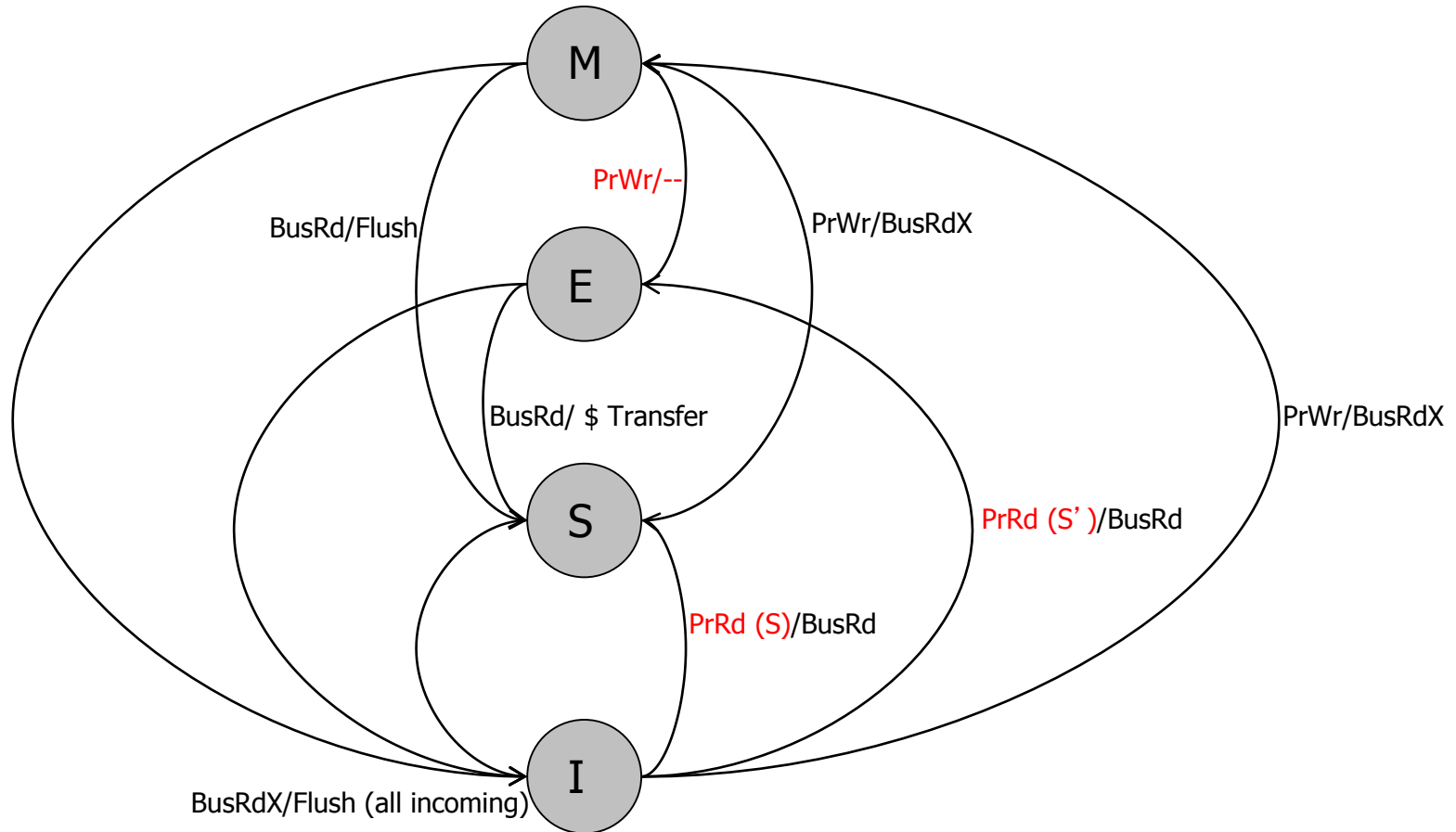
4 States

- M: Modified (Exclusive copy, modified)
- E: Exclusive (" " , clean)
- S: Shared (Shared copy, clean)
- I: Invalid

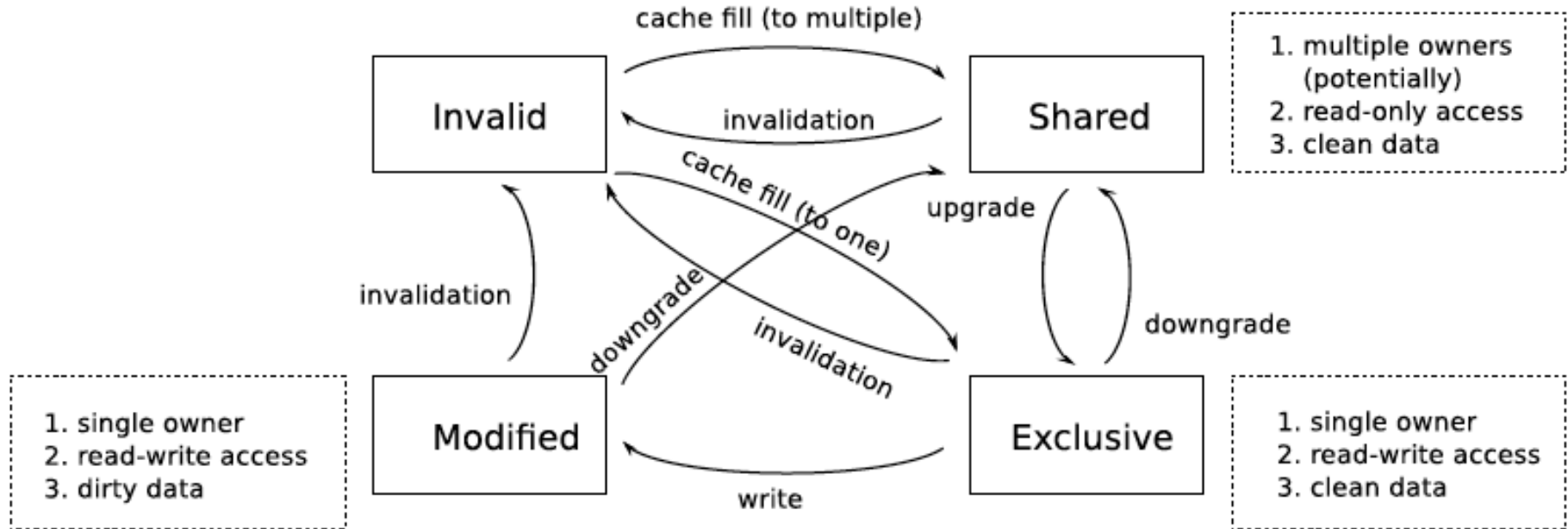
MESI State Machine



MESI State Machine



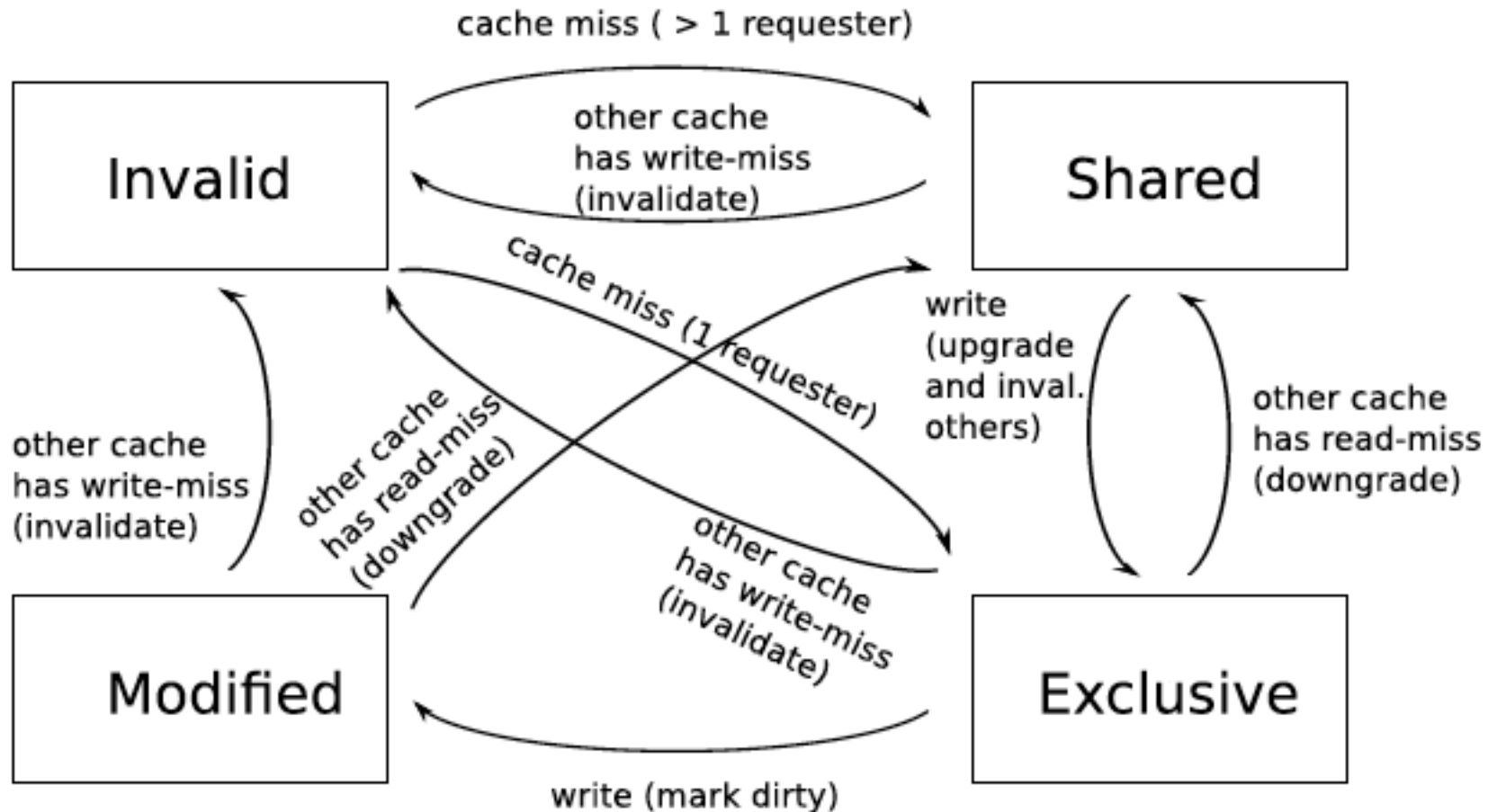
MESI State Machine from Optional Lab 5



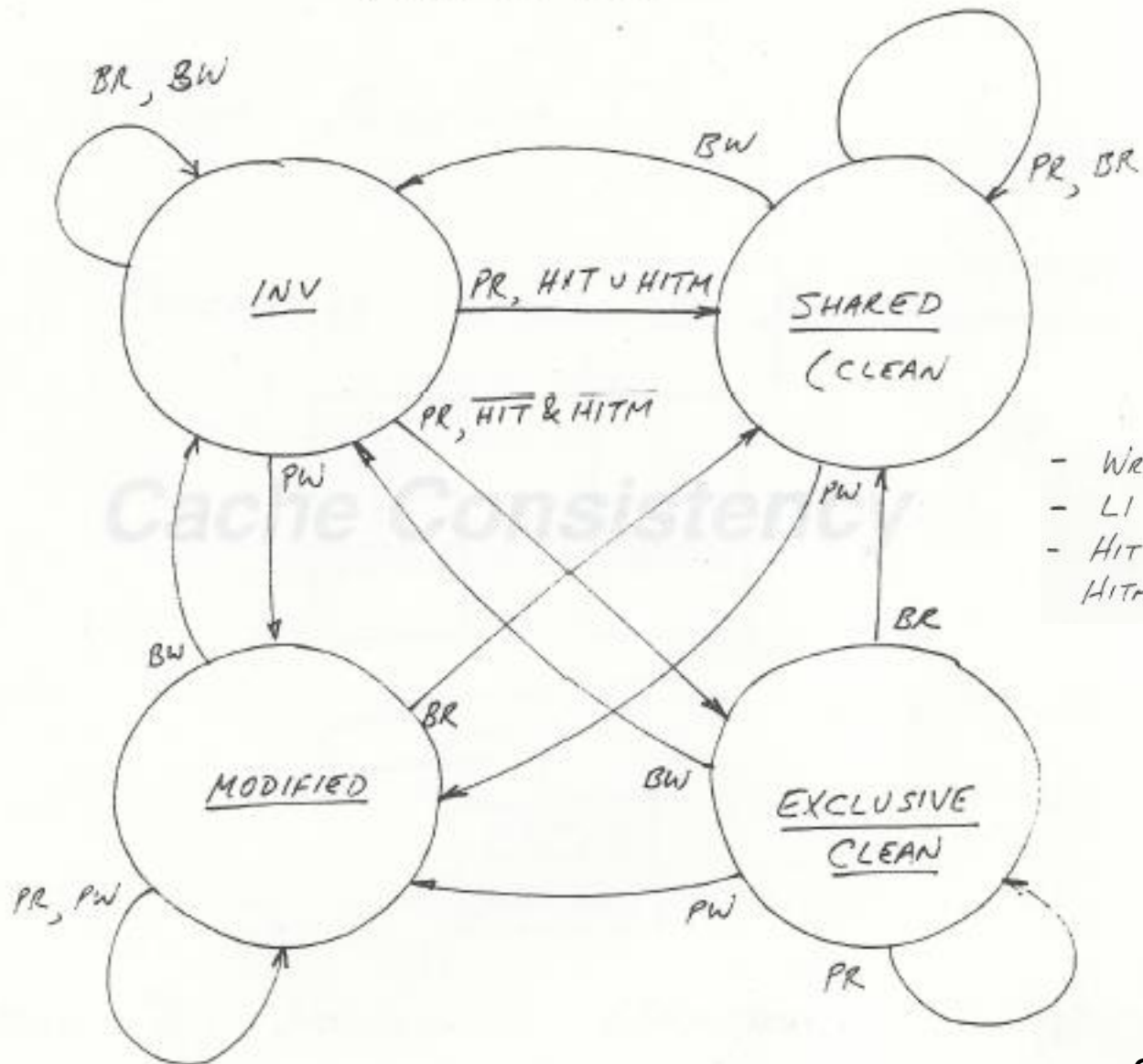
A transition from a single-owner state (Exclusive or Modified) to Shared is called a **downgrade**, because the transition takes away the owner's right to modify the data

A transition from Shared to a single-owner state (Exclusive or Modified) is called an **upgrade**, because the transition grants the ability to the owner (the cache which contains the respective block) to write to the block.

MESI State Machine from Optional Lab 5



Intel Pentium Pro



- WRITE ALLOCATE
- L1 CAN HAVE DATA NOT IN L2
- HIT : SOMEONE HAS IT CLEAN
HITM : SOMEONE HAS IT DIRTY

Snoopy Invalidation Tradeoffs

- Should a downgrade from M go to S or I?
 - S: if data is likely to be reused (before it is written to by another processor)
 - I: if data is likely to be not reused (before it is written to by another)
- Cache-to-cache transfer
 - On a BusRd, should data come from another cache or memory?
 - Another cache
 - May be faster, if memory is slow or highly contended
 - Memory
 - Simpler: no need to wait to see if another cache has the data first
 - Less contention at the other caches
 - Requires writeback on M downgrade
- Writeback on Modified->Shared: necessary?
 - One possibility: **Owner** (O) state (MOESI protocol)
 - One cache owns the latest data (memory is not updated)
 - Memory writeback happens when all caches evict copies

The Problem with MESI

- Observation: Shared state requires the data to be clean
 - i.e., all caches that have the block have the up-to-date copy and so does the memory
- Problem: Need to write the block to memory when BusRd happens when the block is in Modified state
- Why is this a problem?
 - Memory can be updated unnecessarily → some other processor may want to write to the block again

Improving on MESI

- Idea 1: Do not transition from $M \rightarrow S$ on a BusRd. Invalidate the copy and supply the modified block to the requesting processor directly without updating memory
- Idea 2: Transition from $M \rightarrow S$, but designate one cache as the owner (O), who will write the block back when it is evicted
 - Now “Shared” means “Shared and potentially dirty”
 - This is a version of the MOESI protocol

Tradeoffs in Sophisticated Cache Coherence Protocols

- The protocol can be optimized with more states and prediction mechanisms to
 - + Reduce unnecessary invalidates and transfers of blocks
- However, more states and optimizations
 - Are more difficult to design and verify (lead to more cases to take care of, race conditions)
 - Provide diminishing returns

Revisiting Two Cache Coherence Methods

- ❑ How do we ensure that the proper caches are updated?
- ❑ **Snoopy Bus** [Goodman ISCA 1983, Papamarcos+ ISCA 1984]
 - Bus-based, *single point of serialization for all memory requests*
 - Processors observe other processors' actions
 - ❑ E.g.: P1 makes “read-exclusive” request for A on bus, P0 sees this and invalidates its own copy of A
- ❑ **Directory** [Censier and Feautrier, IEEE ToC 1978]
 - *Single point of serialization per block*, distributed among nodes
 - Processors make explicit requests for blocks
 - Directory tracks which caches have each block
 - Directory coordinates invalidation and updates
 - ❑ E.g.: P1 asks directory for exclusive copy, directory asks P0 to invalidate, waits for ACK, then responds to P1

Snoopy Cache vs. Directory Coherence

■ Snoopy Cache

- + Miss latency (critical path) is short: request → bus transaction to mem.
- + Global serialization is easy: bus provides this already (arbitration)
- + Simple: can adapt bus-based uniprocessors easily
- Relies on broadcast messages to be seen by all caches (in same order):
 - single point of serialization (bus): *not scalable*
 - *need a virtual bus (or a totally-ordered interconnect)*

■ Directory

- Adds indirection to miss latency (critical path): request → dir. → mem.
- Requires extra storage space to track sharer sets
 - Can be approximate (false positives are OK for correctness)
- Protocols and race conditions are more complex (for high-performance)
- + Does not require broadcast to all caches
- + Exactly as scalable as interconnect and directory storage
(much more scalable than bus)

Revisiting Directory-Based Cache Coherence

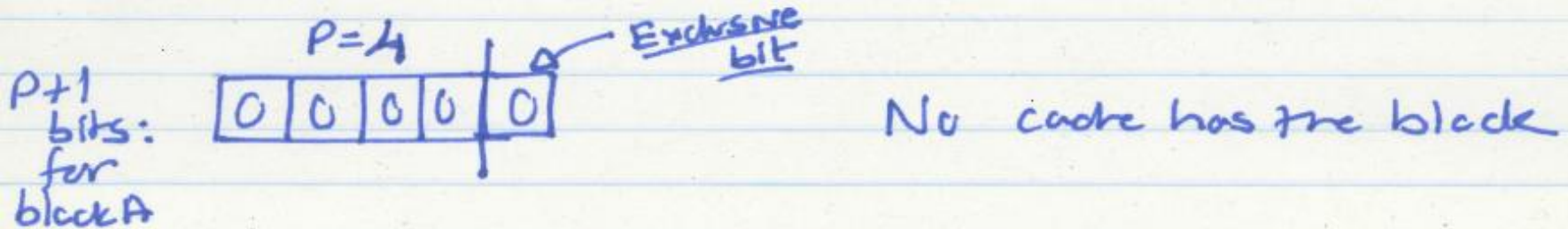
Remember: Directory Based Coherence

- Idea: A logically-central directory keeps track of where the copies of each cache block reside. Caches consult this directory to ensure coherence.

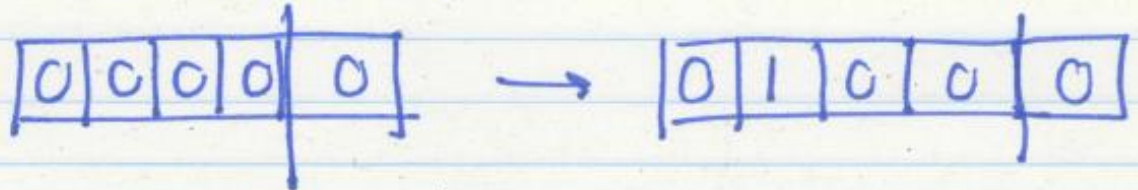
- An example mechanism:
 - For each cache block in memory, store $P+1$ bits in directory
 - One bit for each cache, indicating whether the block is in cache
 - Exclusive bit: indicates that the cache that has the only copy of the block and can update it without notifying others
 - On a read: set the cache's bit and arrange the supply of data
 - On a write: invalidate all caches that have the block and reset their bits
 - Have an "exclusive bit" associated with each block in each cache

Remember: Directory Based Coherence

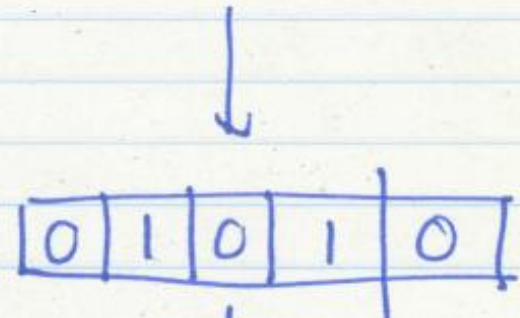
Example directory based scheme



① P_1 takes a read miss to block A



② P_3 takes a read miss



Directory-Based Protocols

- Required when scaling past the capacity of a single bus
- Distributed, *but*:
 - Coherence still requires single point of serialization (for write serialization)
 - Serialization location can be different for every block (striped across nodes/memory-controllers)
- We can reason about the protocol for a single block: one *server* (directory node), many *clients* (private caches)
- Directory receives *Read* and *ReadEx* requests, and sends *Inv!* requests: invalidation is explicit (as opposed to snoopy buses)

Directory: Data Structures

0x00	Shared: {P0, P1, P2}
0x04	---
0x08	Exclusive: P2
0x0C	---
...	---

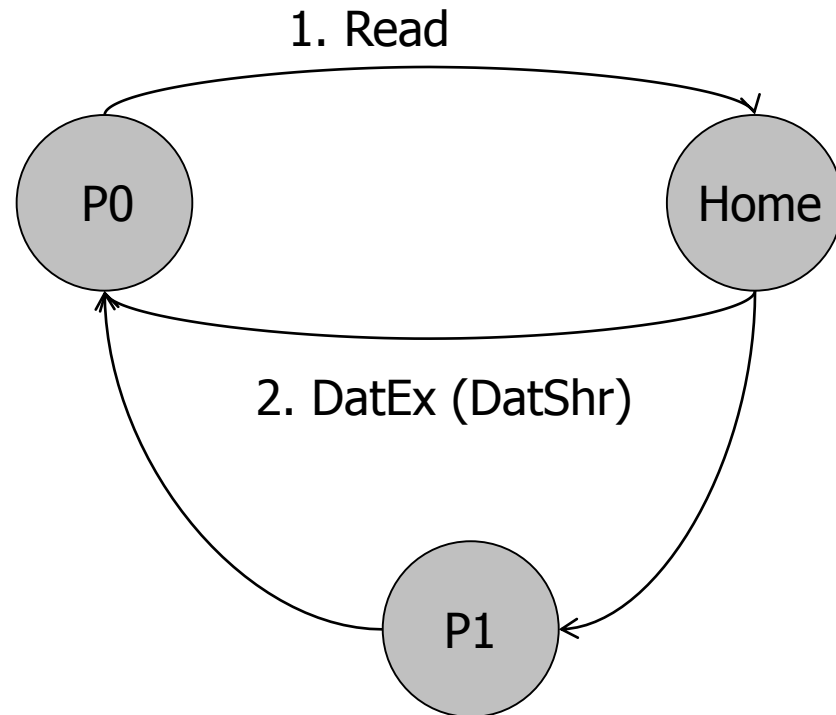
- Required to support invalidation and cache block requests
- Key operation to support is *set inclusion test*
 - False positives are OK: want to know which caches *may* contain a copy of a block, and spurious invalidations are ignored
 - False positive rate determines *performance*
- Most accurate (and expensive): full bit-vector
- Compressed representation, linked list, Bloom filters are all possible

Directory: Basic Operations

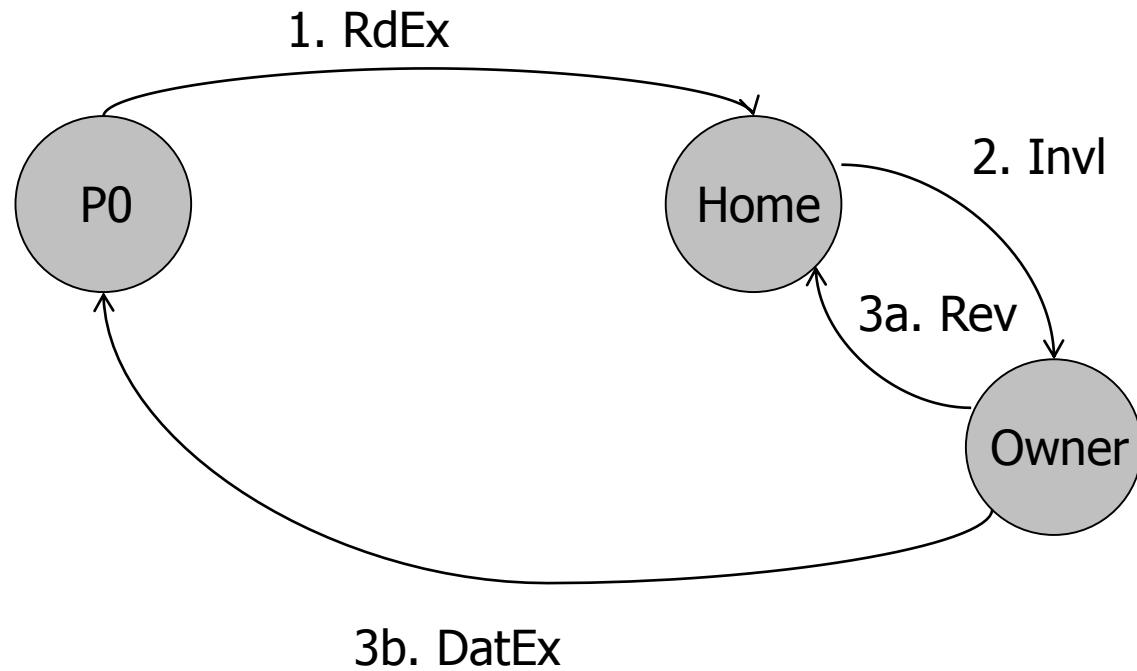
- Follow *semantics* of snoop-based system
 - but with explicit request, reply messages
- Directory:
 - Receives *Read, ReadEx, Upgrade* requests from nodes
 - Sends *Inval/Downgrade* messages to sharers if needed
 - Forwards request to memory if needed
 - Replies to requestor and updates sharing state
- Protocol design is flexible
 - Exact forwarding paths depend on implementation
 - For example, do cache-to-cache transfer?

MESI Directory Transaction: Read

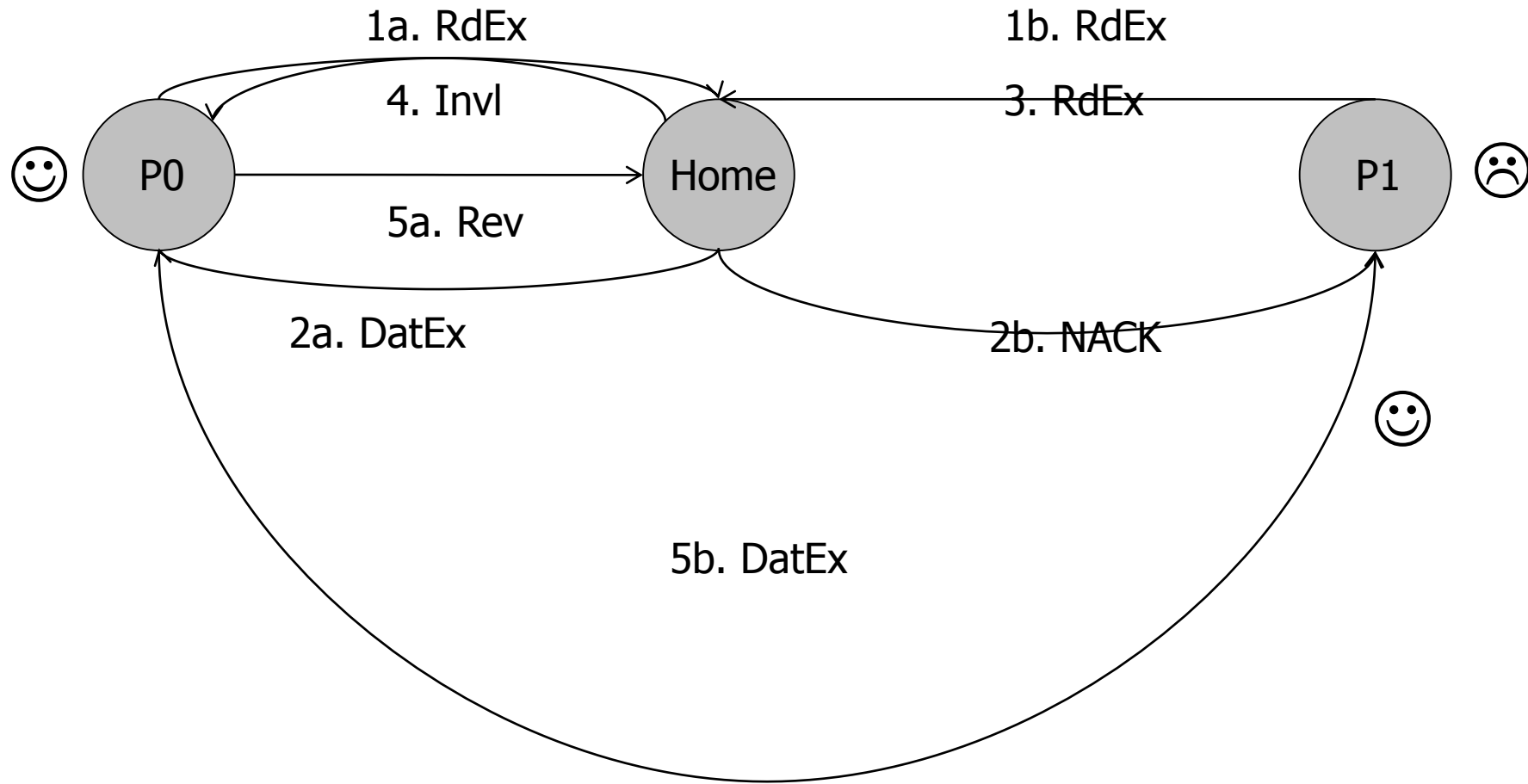
P0 acquires an address for reading:



RdEx with Former Owner



Contention Resolution (for Write)



Issues with Contention Resolution

- Need to escape race conditions by:
 - NACKing requests to busy (pending invalidate) entries
 - Original requestor retries
 - OR, queuing requests and granting in sequence
 - (Or some combination thereof)
- Fairness
 - Which requestor should be preferred in a conflict?
 - Interconnect delivery order, and distance, both matter
- Ping-ponging is a higher-level issue
 - With solutions like combining trees (for locks/barriers) and better shared-data-structure design

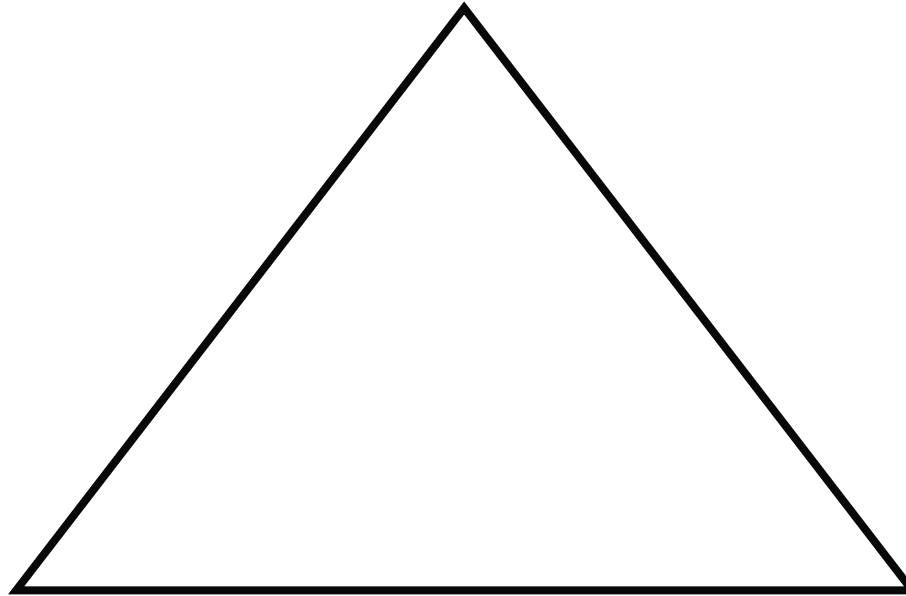
Scaling the Directory: Some Questions

- How large is the directory?
- How can we reduce the access latency to the directory?
- How can we scale the system to thousands of nodes?
- Can we get the best of snooping and directory protocols?
 - Heterogeneity
 - E.g., token coherence [Martin+, ISCA 2003]

Advancing Coherence

Motivation: Three Desirable Attributes

Low-latency cache-to-cache misses



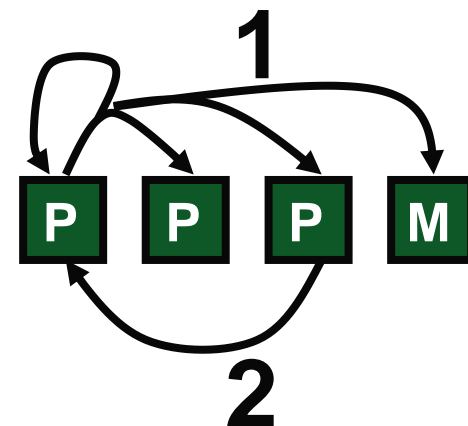
No bus-like interconnect

Bandwidth efficient

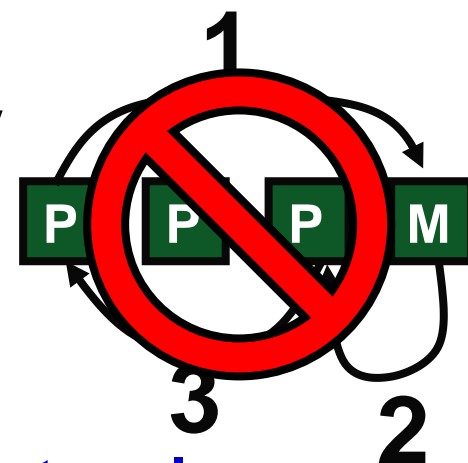
Dictated by workload and technology trends

Workload Trends

- **Commercial workloads**
 - Many cache-to-cache misses
 - Clusters of small multiprocessors
- **Goals:**
 - **Direct cache-to-cache misses (2 hops, not 3 hops)**
 - Moderate scalability



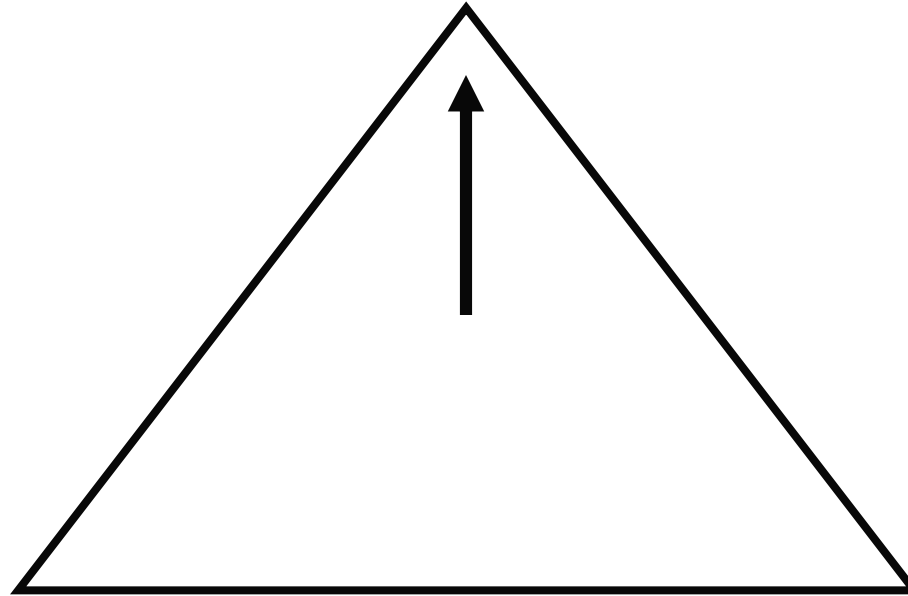
Directory
Protocol



Workload trends → snooping protocols

Workload Trends

Low-latency cache-to-cache misses



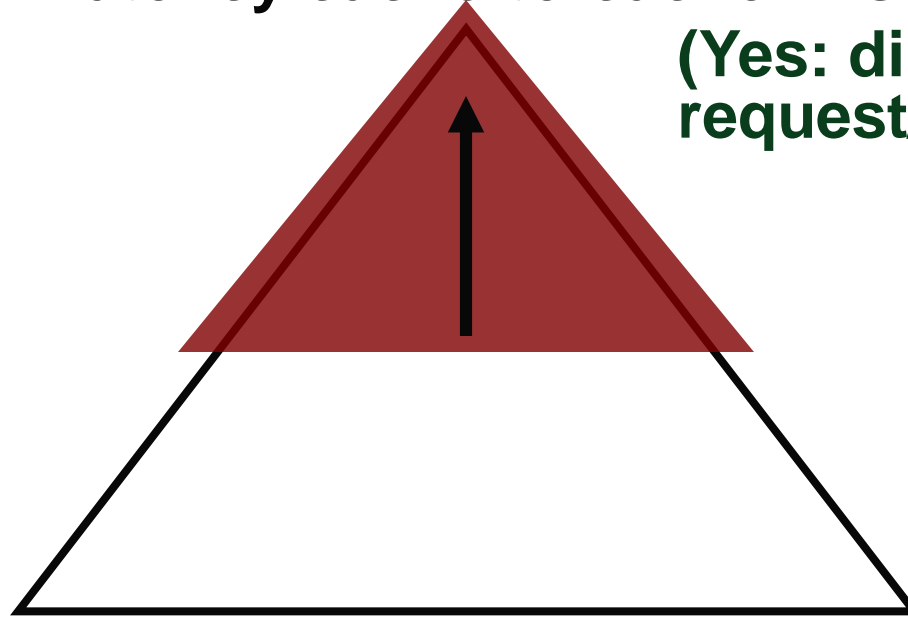
No bus-like interconnect

Bandwidth efficient

Workload Trends □ Snooping Protocols

Low-latency cache-to-cache misses

**(Yes: direct
request/response)**



No bus-like interconnect

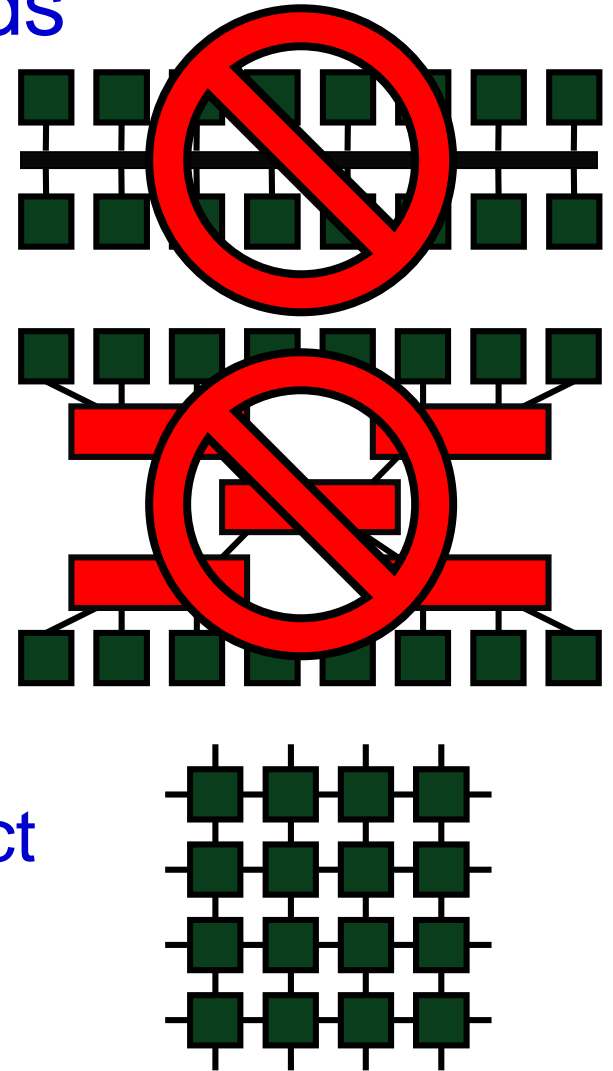
(No: requires a “virtual bus”)

Bandwidth efficient

(No: broadcast always)

Technology Trends

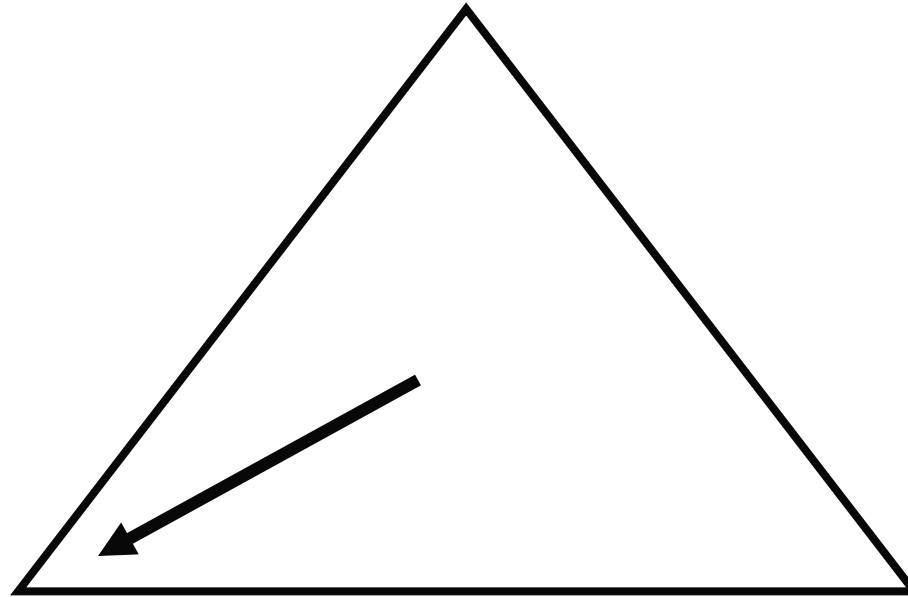
- High-speed point-to-point links
 - No (multi-drop) busses
- Increasing design integration
 - “Glueless” multiprocessors
 - Improve cost & latency
- Desire: low-latency interconnect
 - Avoid “virtual bus” ordering
 - Enabled by directory protocols



Technology trends → unordered interconnects

Technology Trends

Low-latency cache-to-cache misses



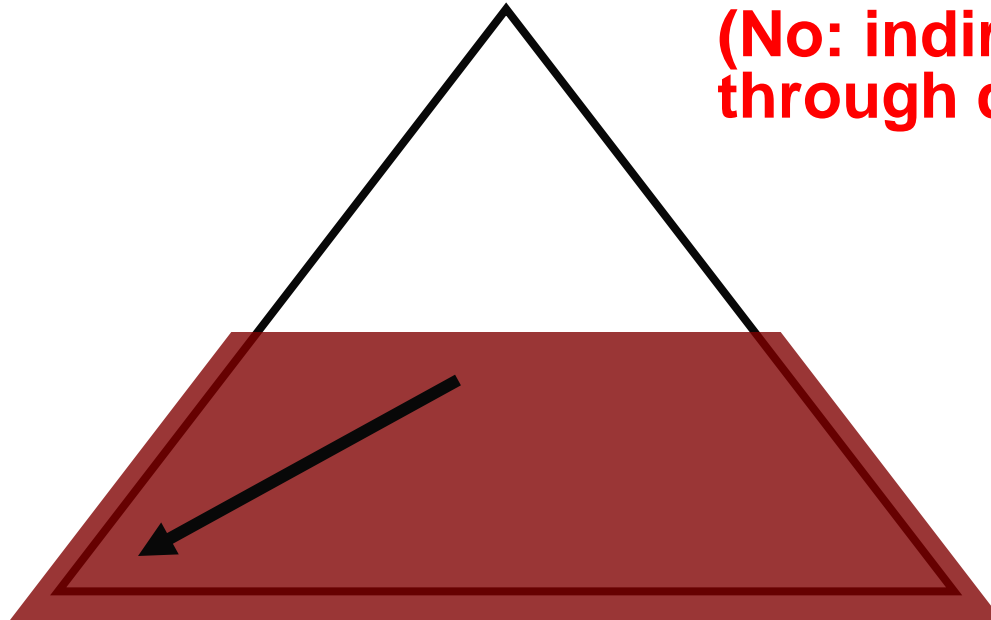
No bus-like interconnect

Bandwidth efficient

Technology Trends □ Directory Protocols

Low-latency cache-to-cache misses

**(No: indirection
through directory)**



No bus-like interconnect
(Yes: no ordering required)

Bandwidth efficient
(Yes: avoids broadcast)

Goal: All Three Attributes

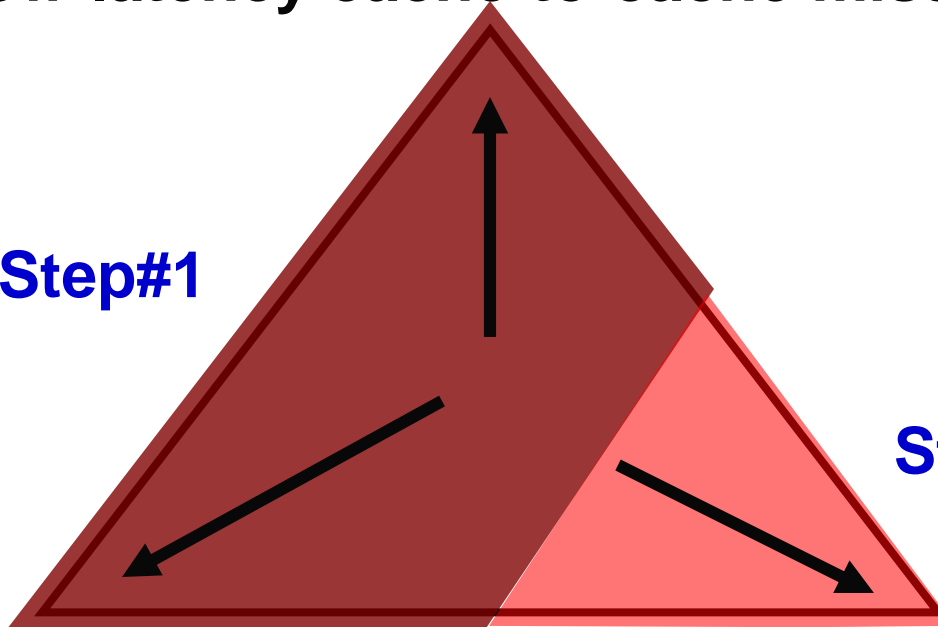
Low-latency cache-to-cache misses

Step#1

Step#2

No bus-like interconnect

Bandwidth efficient



Token Coherence: Key Insight

- Goal of invalidation-based coherence
 - Invariant: **many readers -or- single writer**
 - Enforced by **globally** coordinated actions

Key insight

- Enforce this invariant directly using **tokens**
 - **Fixed number of tokens** per block
 - **One token to read, all tokens to write**
- Guarantees **safety** in all cases
 - Global invariant enforced with only **local** rules
 - Independent of races, request ordering, etc.