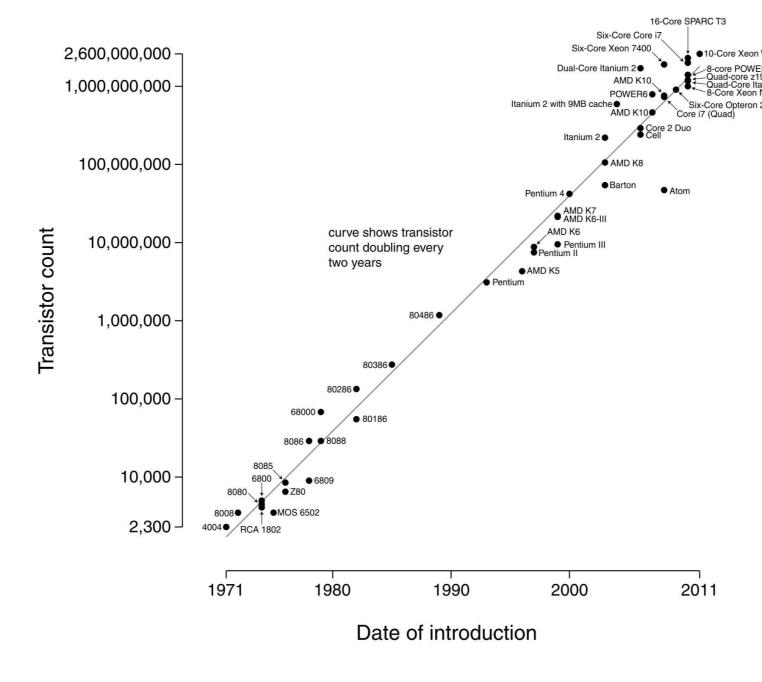
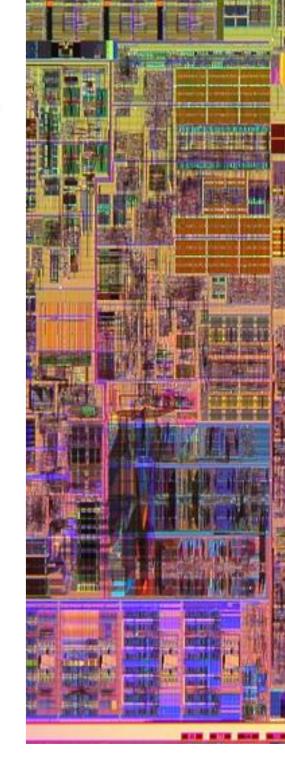




Multicore Computer Architecture Issues: From hardware to software 18-613, Spring 2019

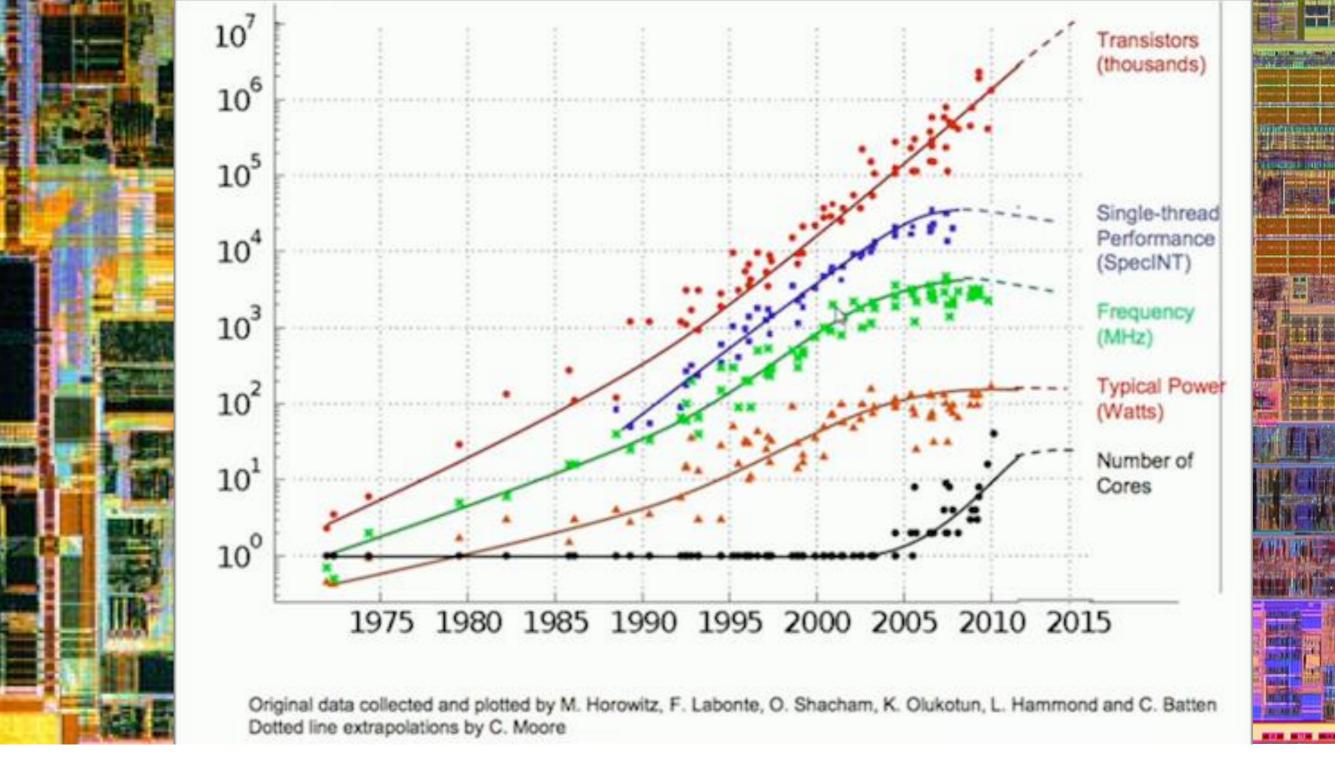
#### Microprocessor Transistor Counts 1971-2011 & Moore's Law





# Technology Improves

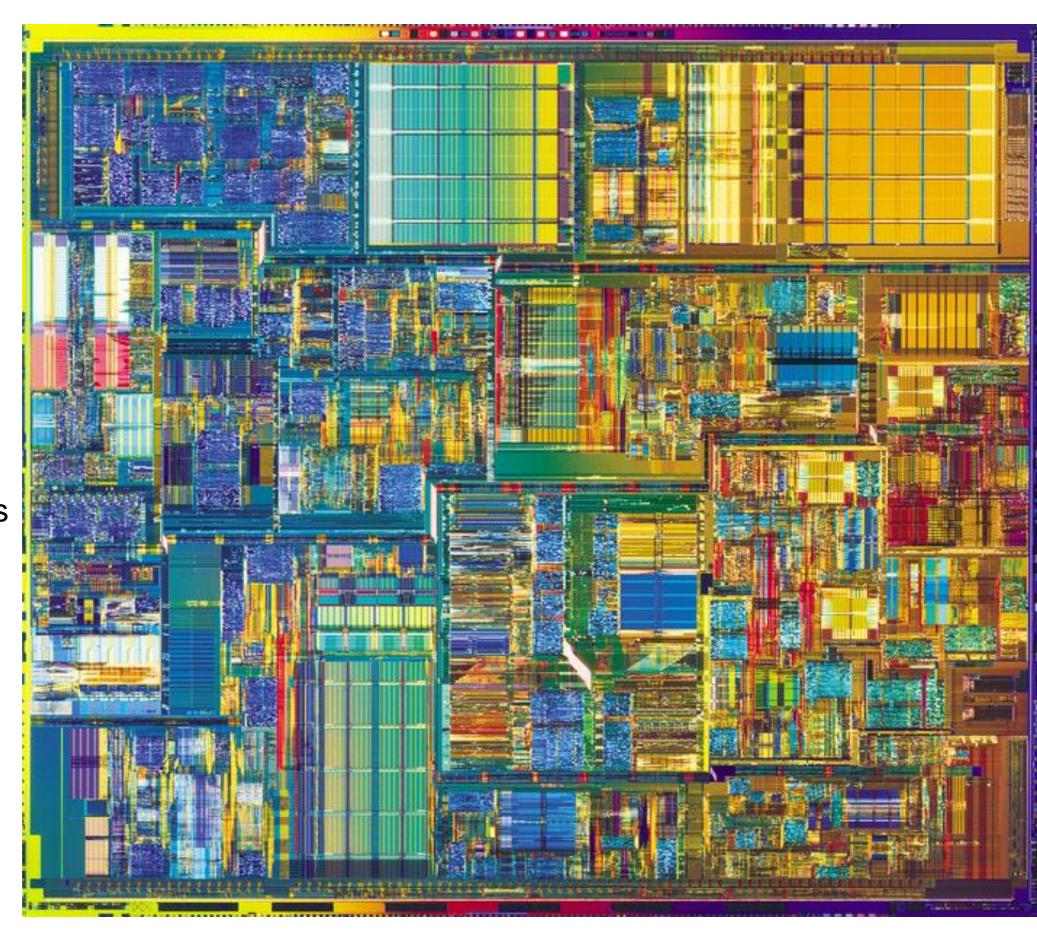
Transistors get smaller, clock speeds go up, power stays roughly constant. Party!



Performance Scaling Hit a Wall!

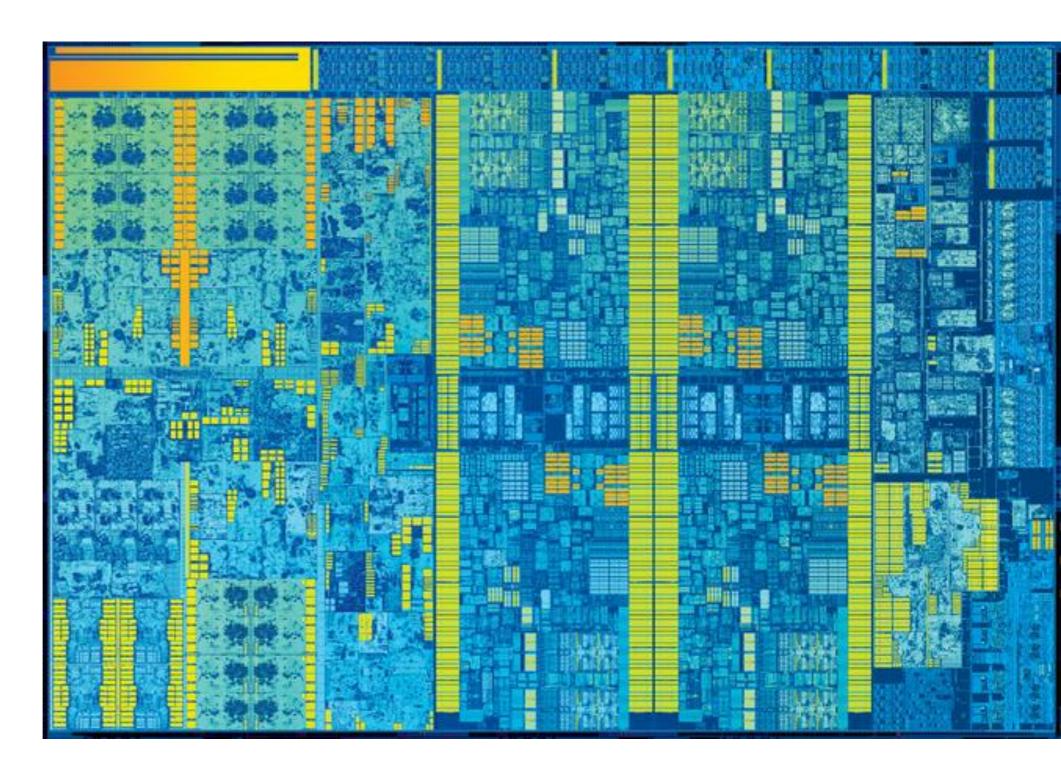
Transistors still doubling, performance tapers off.
Architects need to be creative!

- Willamette core
- •180 nm process
- •217 mm² die size
- •42,000,000 transistors



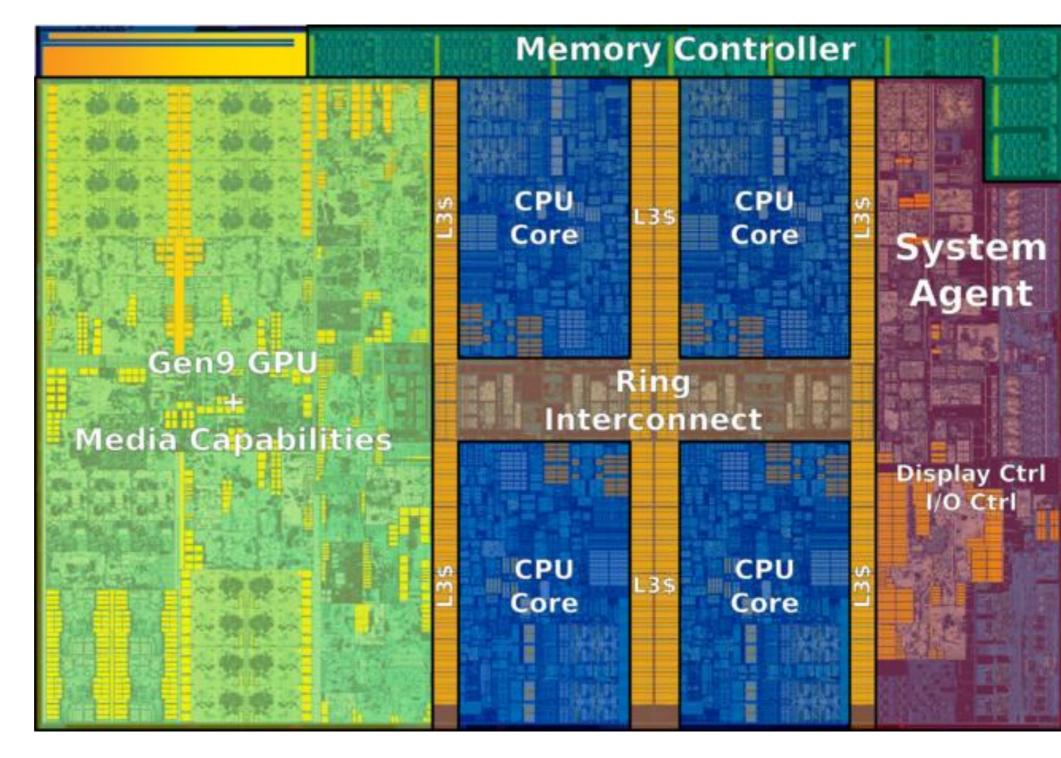
#### •14 nm process

- •11 metal layers
- •~1,750,000,000 transistors
- •~9.19 mm x ~11.08 mm
- •~101.83 mm² die size
- •4 CPU cores + 24 GPU EUs



#### •14 nm process

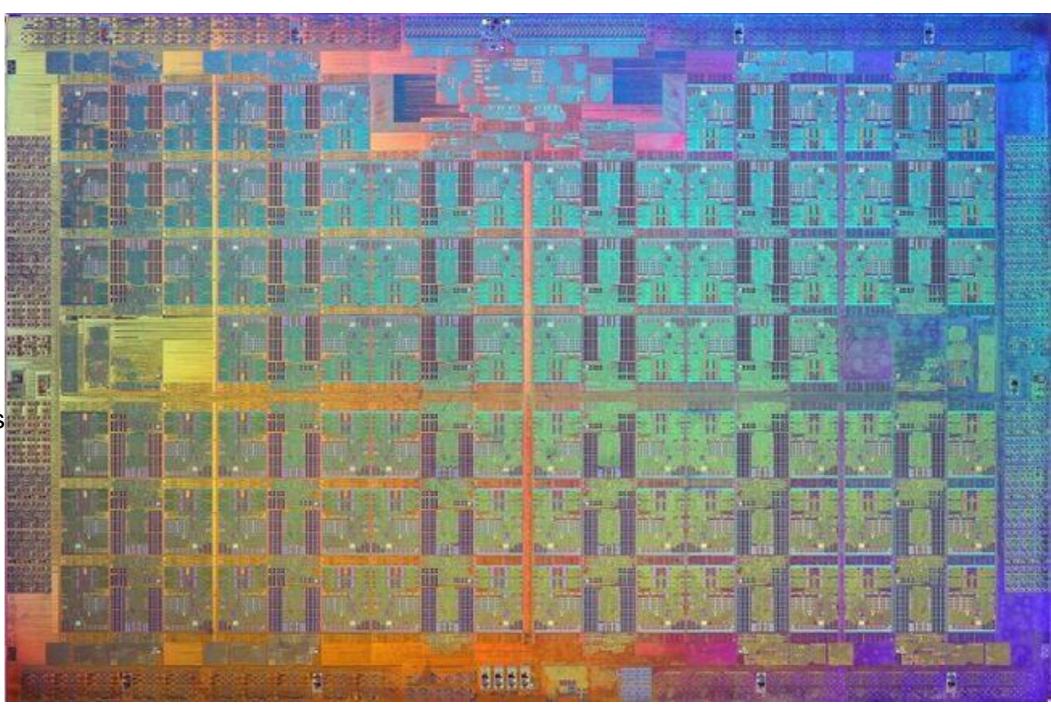
- •11 metal layers
- •~1,750,000,000 transistors
- •~9.19 mm x ~11.08 mm
- •~101.83 mm² die size
- •4 CPU cores + 24 GPU EUs



# Shared memory multi-threading

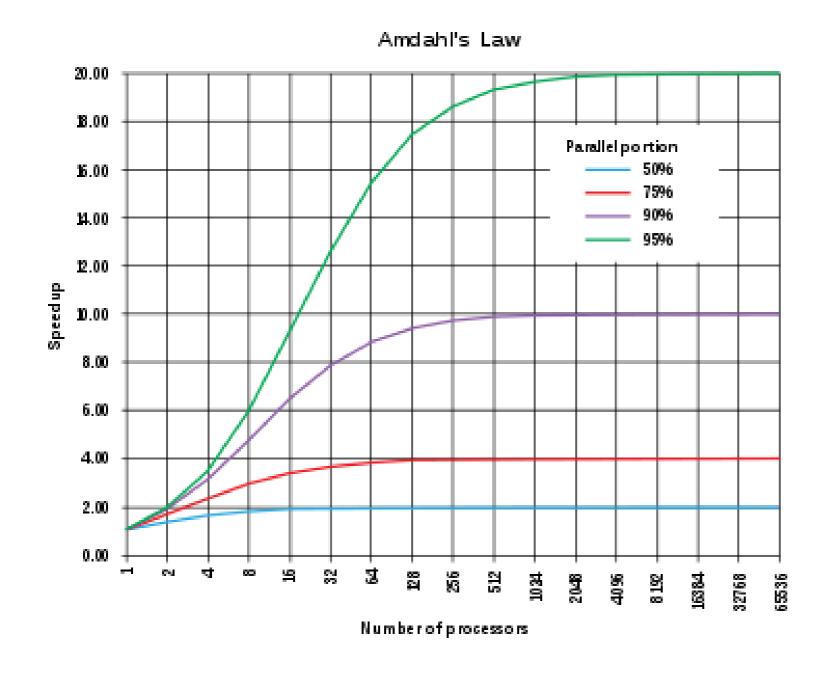


- •682.6 mm² die size
- •76 CPU cores
- •7,100,000,000 transistors



## Amdahl's Law

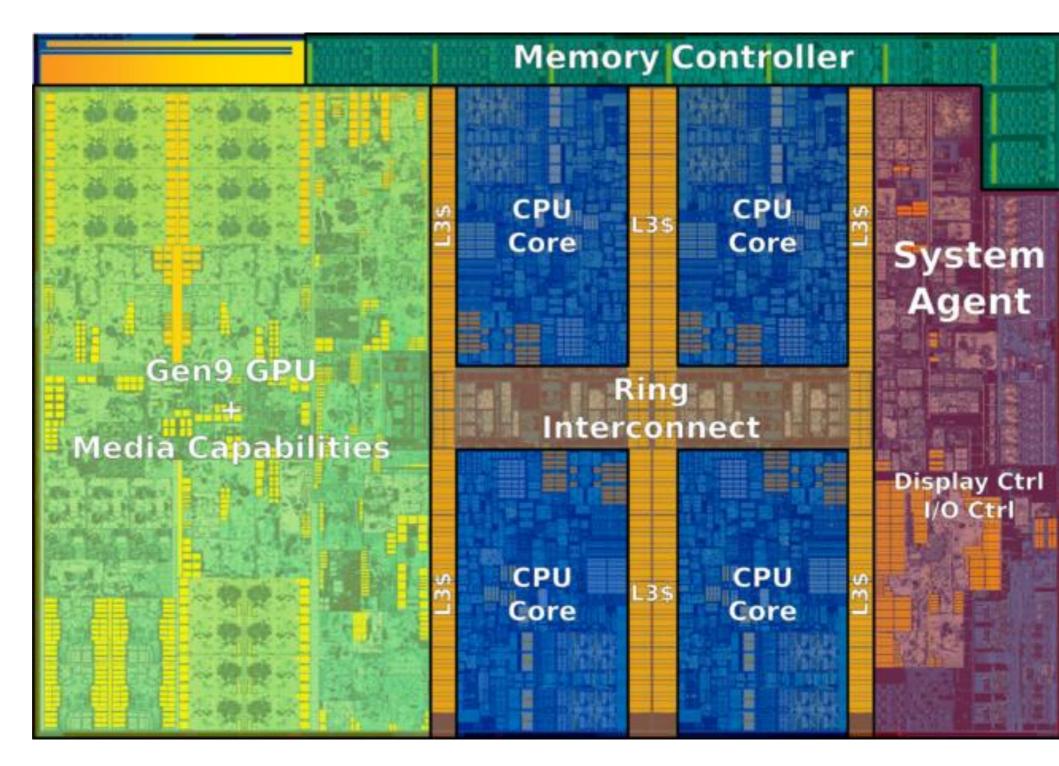
1/((1-p) + p/s)



Amdahl's Corollary:
Speedup is limited by fraction of the *program* that is parallelizable

#### •14 nm process

- •11 metal layers
- •~1,750,000,000 transistors
- •~9.19 mm x ~11.08 mm
- •~101.83 mm² die size
- •4 CPU cores + 24 GPU EUs



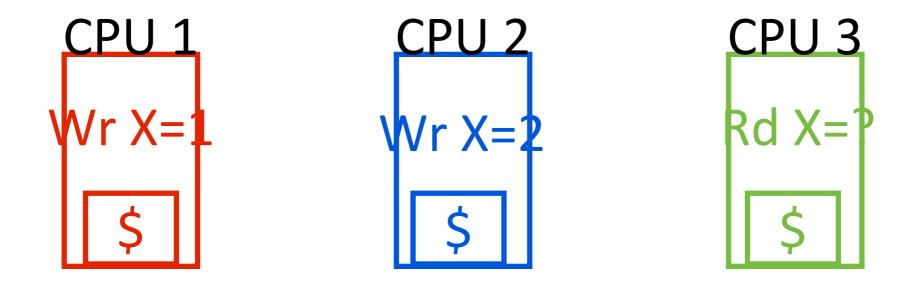


"Coherence seeks to make the caches of a shared-memory system as functionally invisible as the caches in a single-core system. Correct coherence ensures that a programmer cannot determine whether and where a system has caches by analyzing the results of loads and stores."

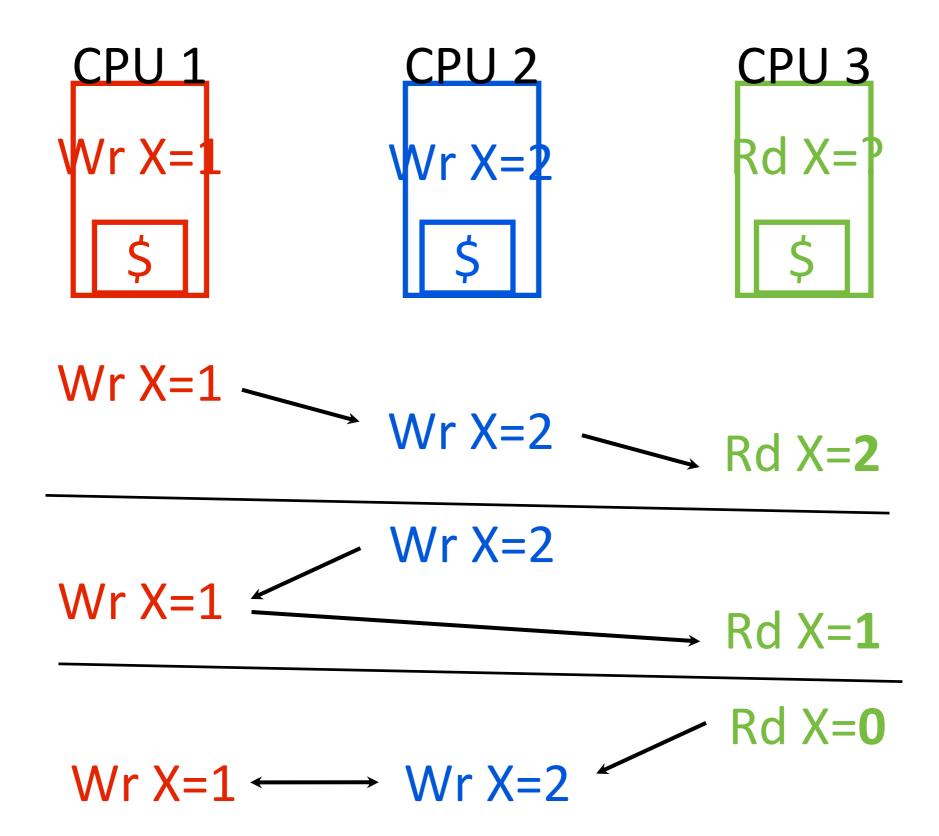
Excerpt from "Primer on Memory Consistency and Cache Coherence"

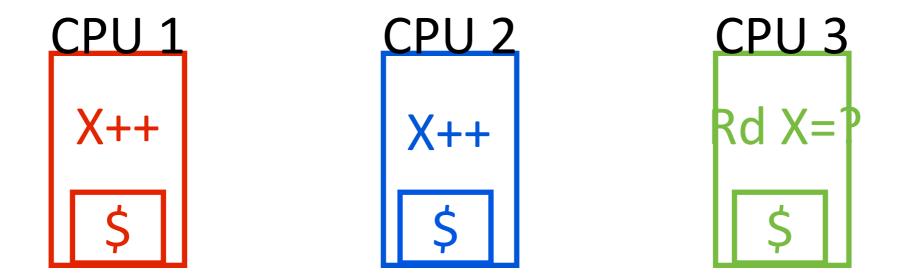
Mark Hill, 2011

# Cache Coherence

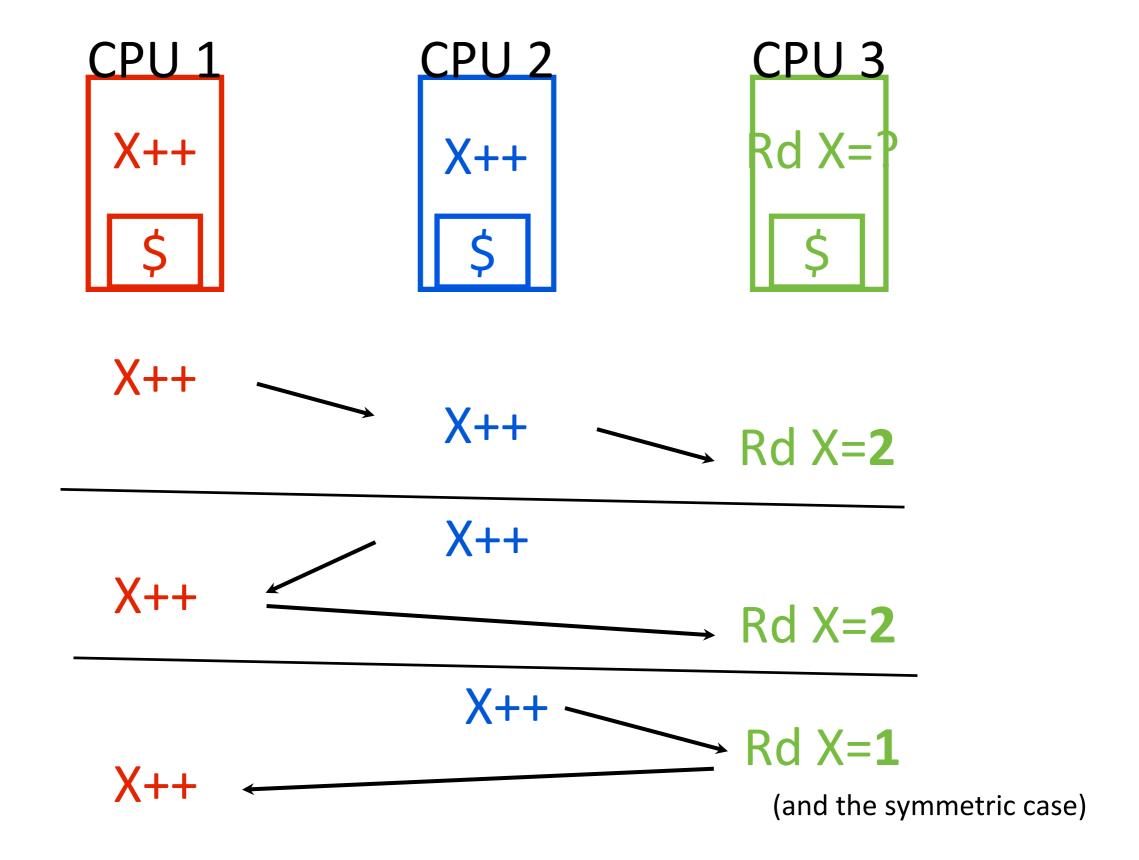


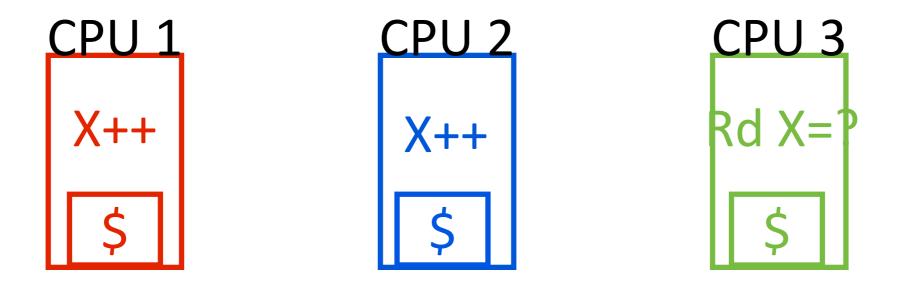
What is the behavior of this parallel program? (X initially 0)



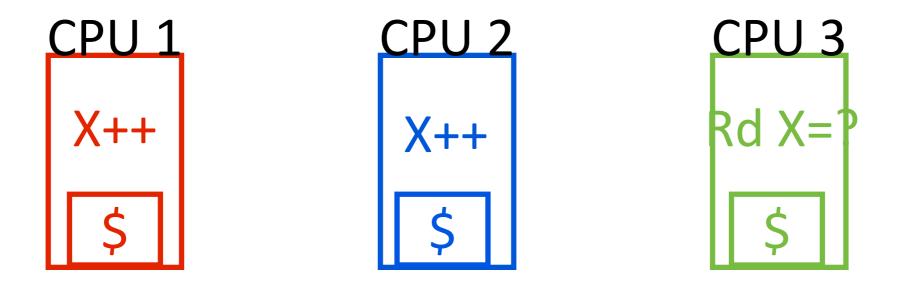


What about this example? (X initially 0)

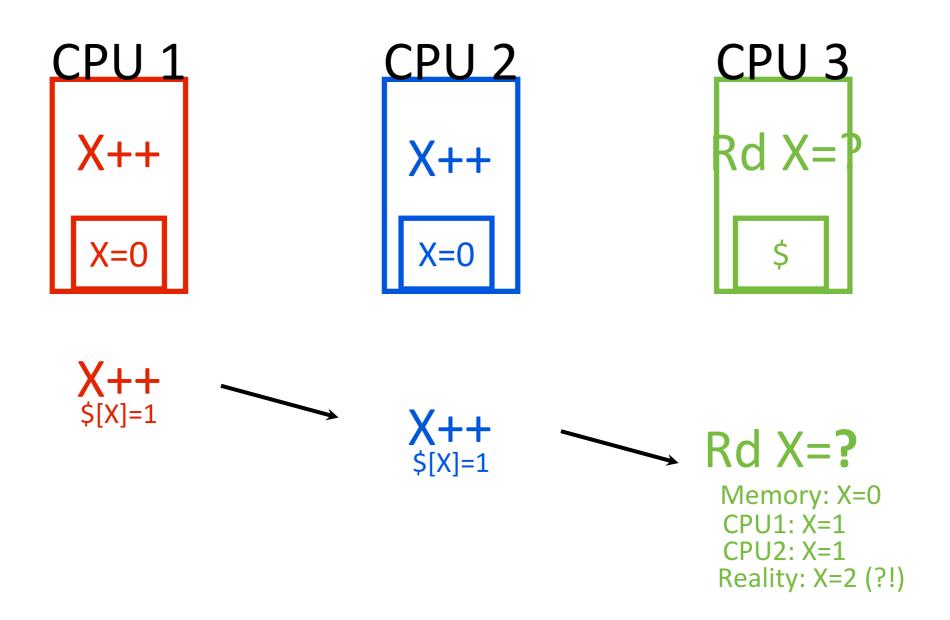




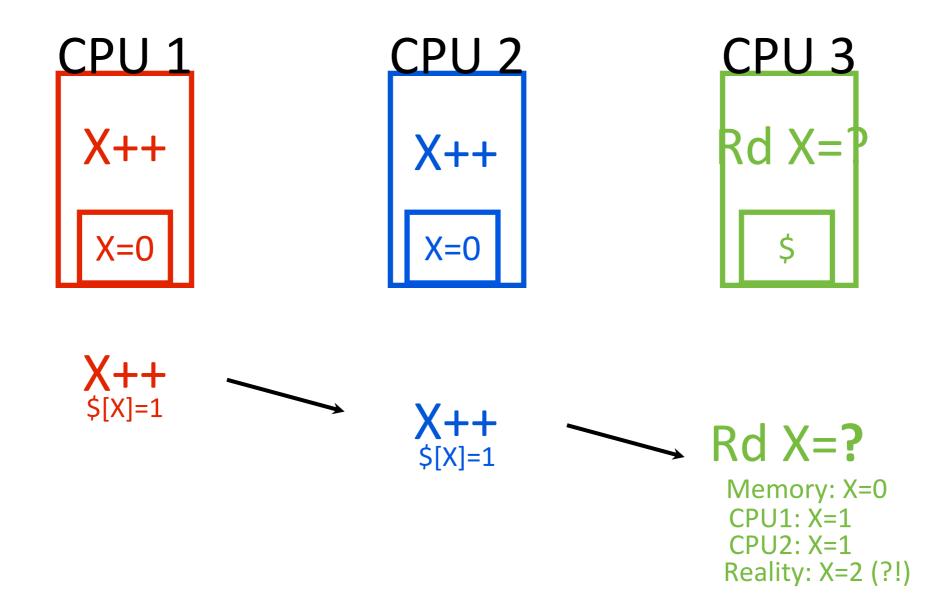
What assumptions are we making about the system to produce the results 0, 1, and 2?



We assume the updates see one anothers' results! (Why wouldn't they?)



So what the heck do we do now?



### Never let this happen. Caches should be coherent.

"coherence ensures that a programmer cannot determine whether and where a system has caches by analyzing the results of loads and stores"

# Informally Defining Coherence

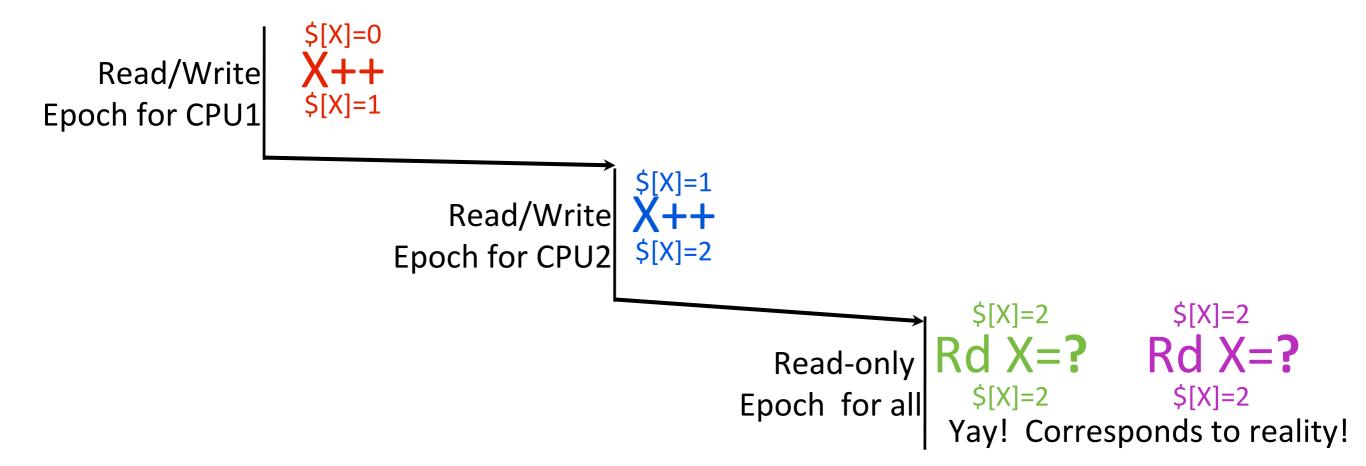
"Coherence serializes all reads with all updates to the same location by different CPUs/caches, so that each read sees the result of the most recent update by any other"

"Single Writer/Multiple Reader (SWMR) Invariant

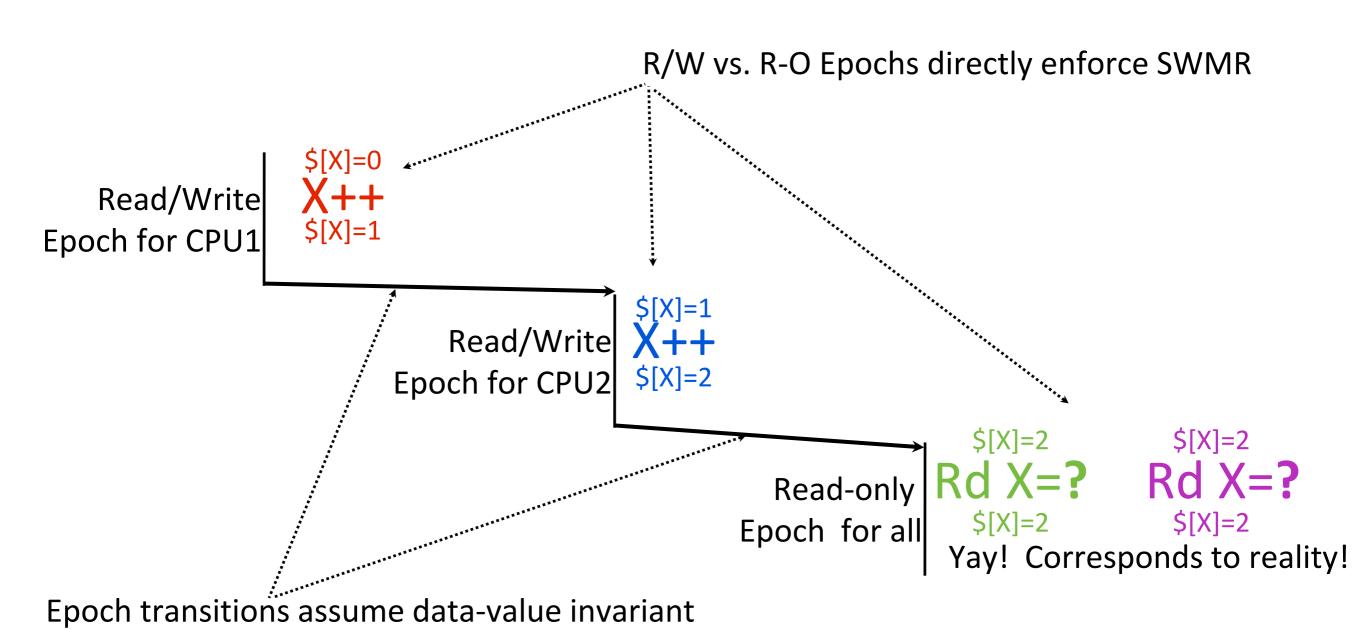
+

Data-Value Invariant"

# Epoch Model



# Epoch Model



# What do we need to implement the Epoch Model?

Need to add concept of R/W epoch vs. R-O epoch

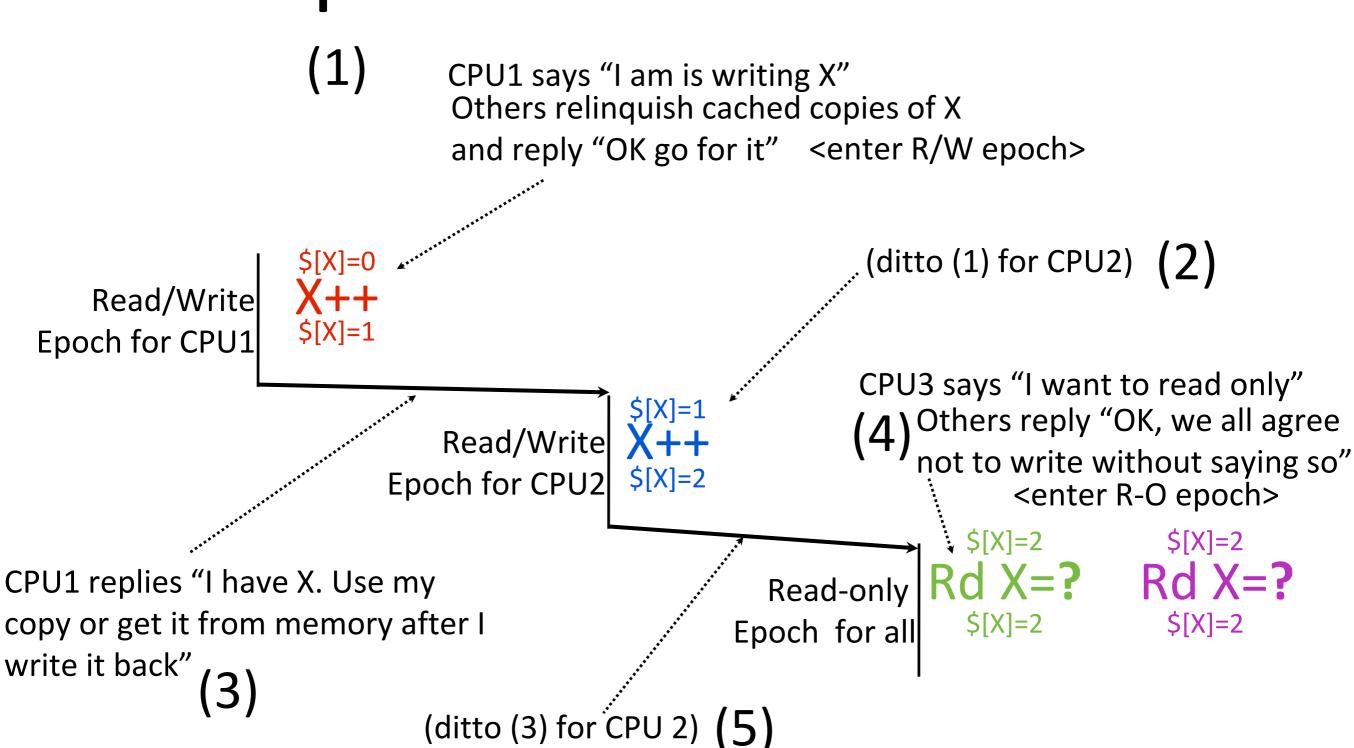
Need to add gadget that correctly moves data between epochs

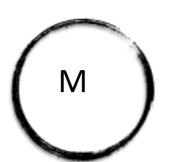
Add state to each cache line saying whether it is R-O or R/W

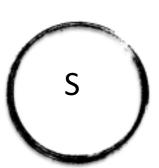
Add protocol actions to move lines from state to state based on (1)local memory operations; and (2)other CPUs' memory operations

Add support to get data from (1)local cache; (2)a remote cache; or (3)main memory, depending on line's protocol state

# High-level sketch of protocol in action

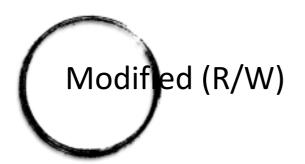


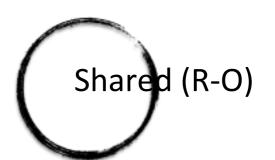




Per-line coherence states

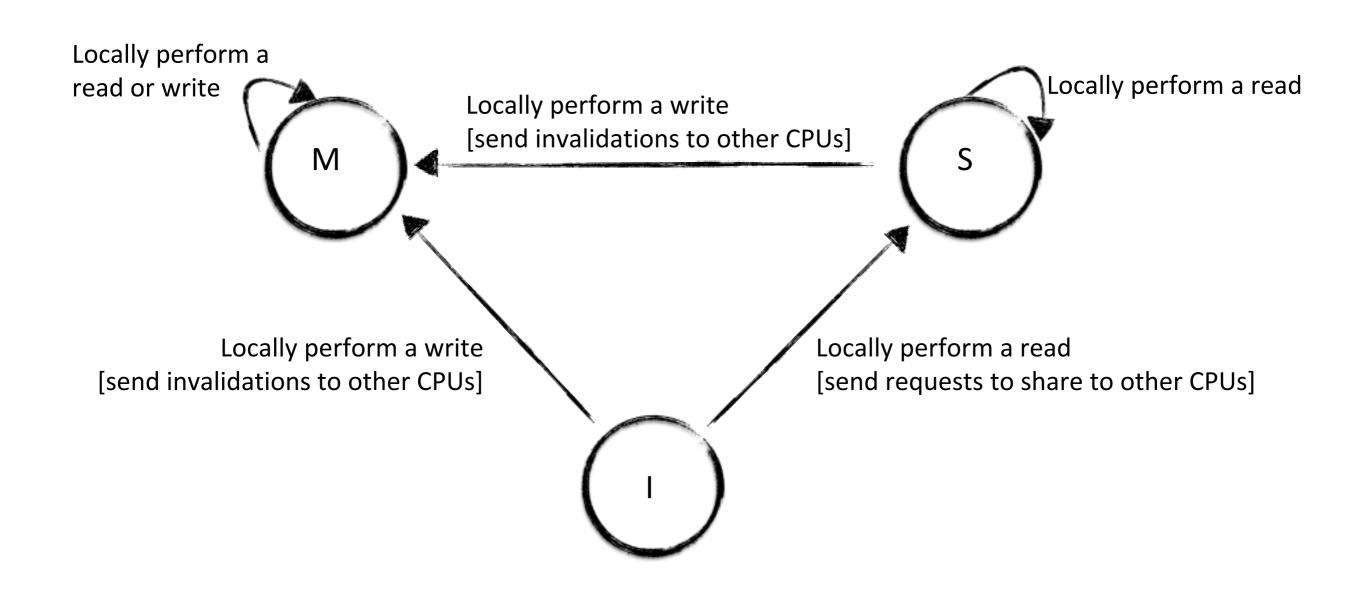




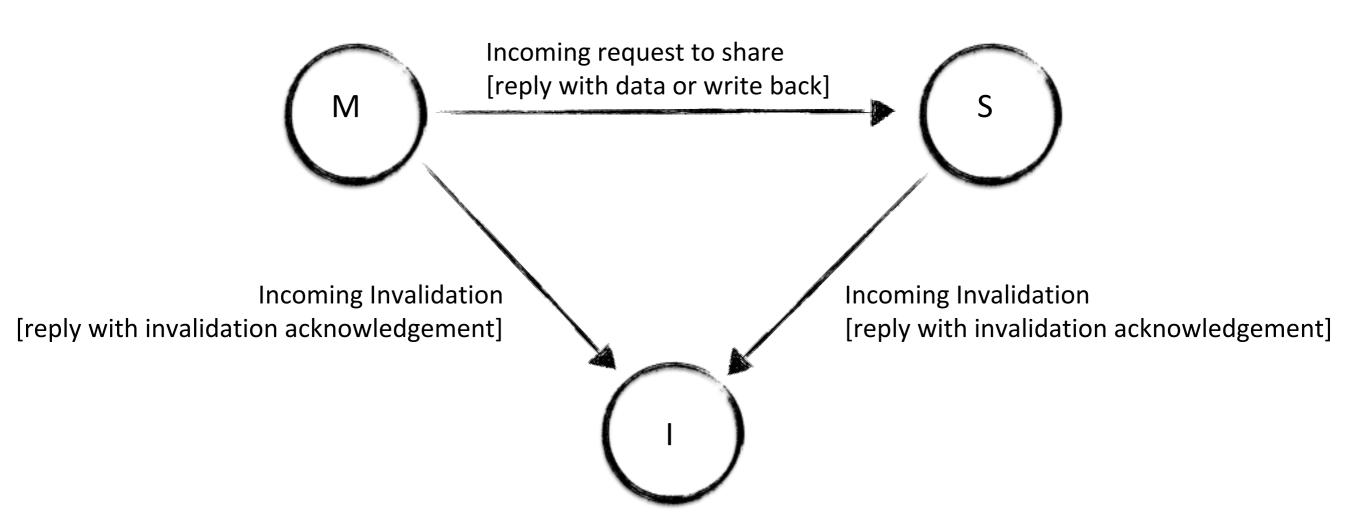




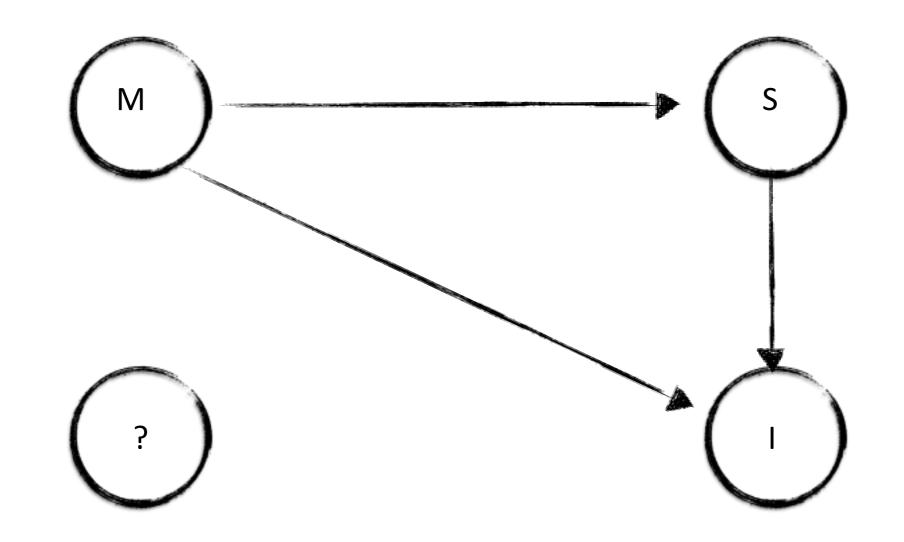
#### Local operations perspective



# Cache Coherence Protocol Remote operations perspective

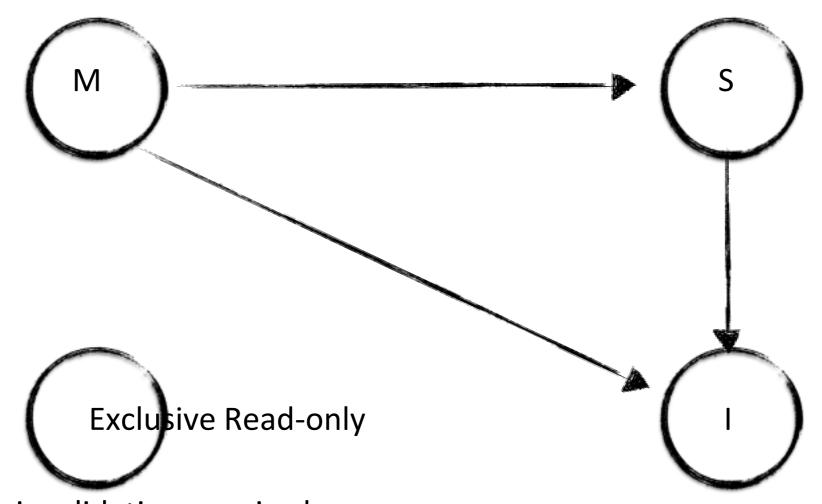


# Can we design another state?

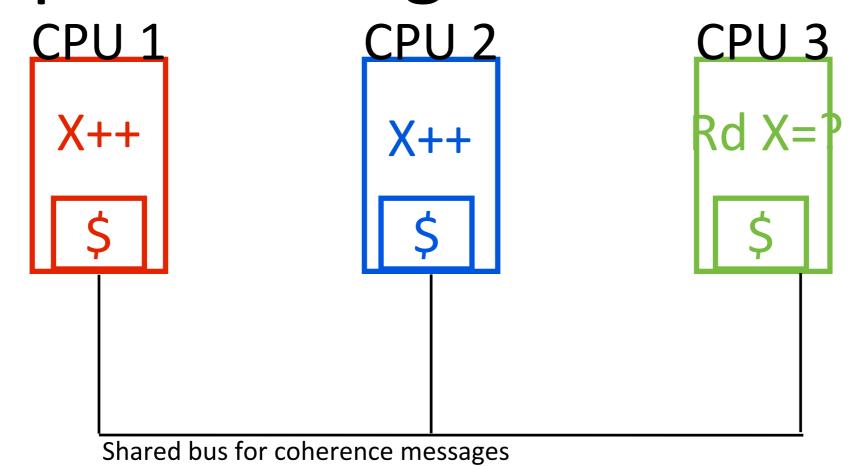


What should we optimize?

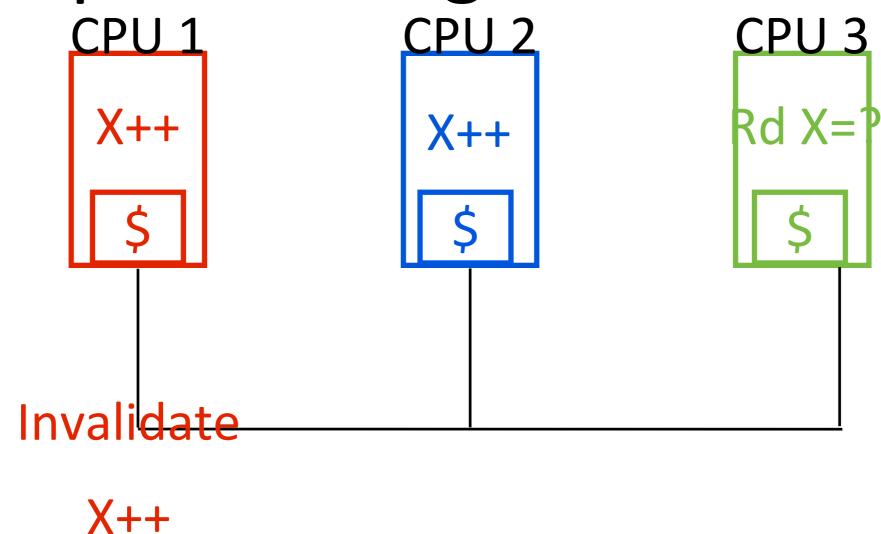
# Can we design another state?

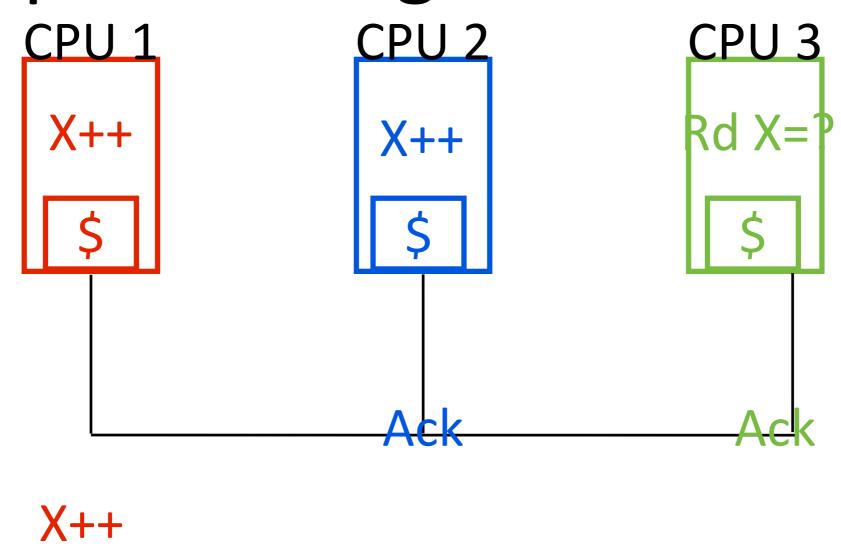


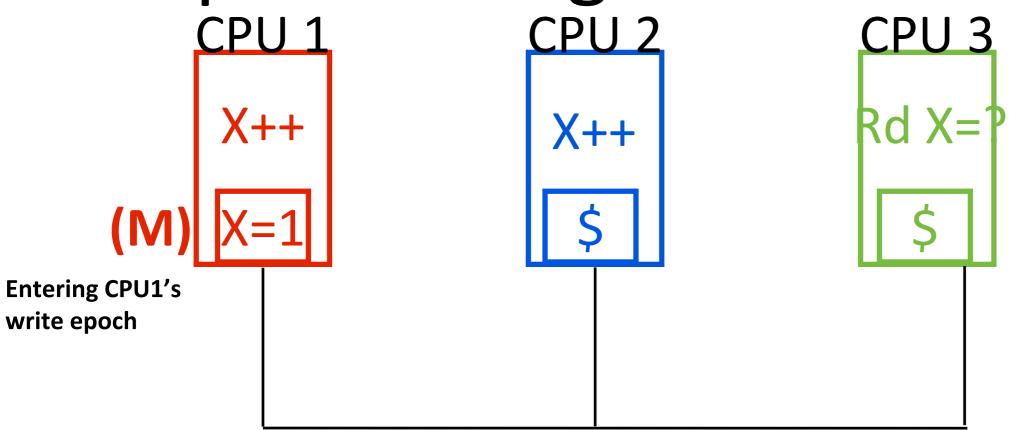
(Benefit: no invalidation required to transition from E->M, like from S->M)

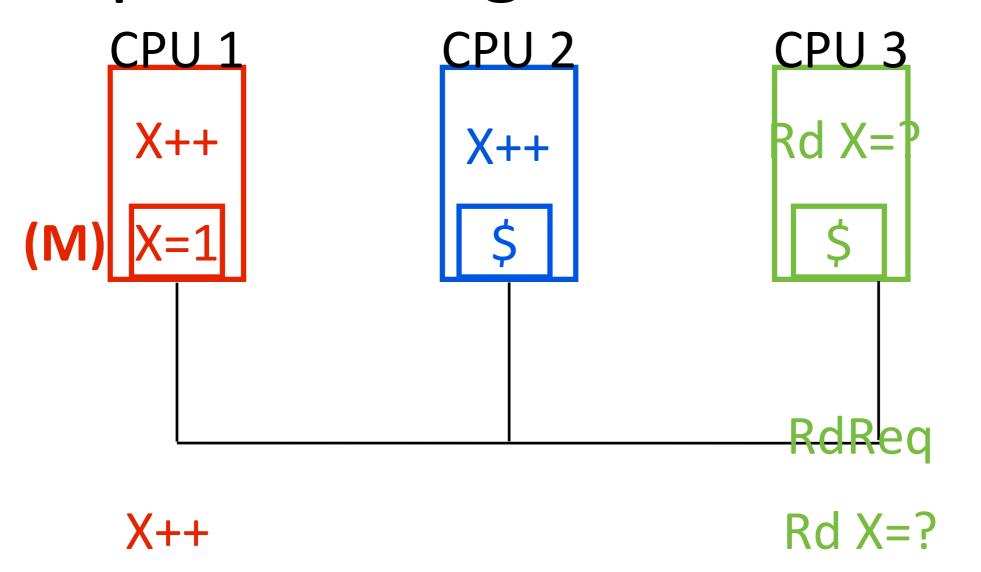


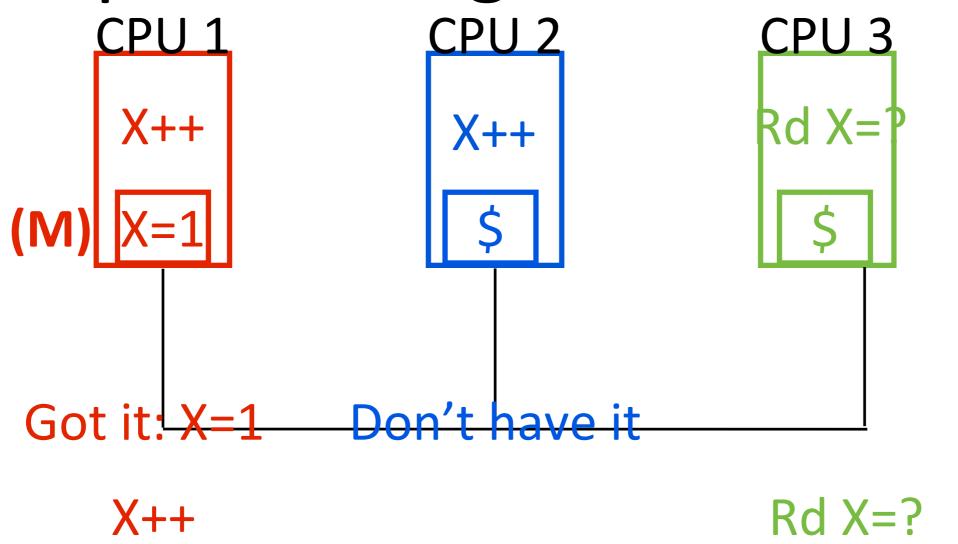


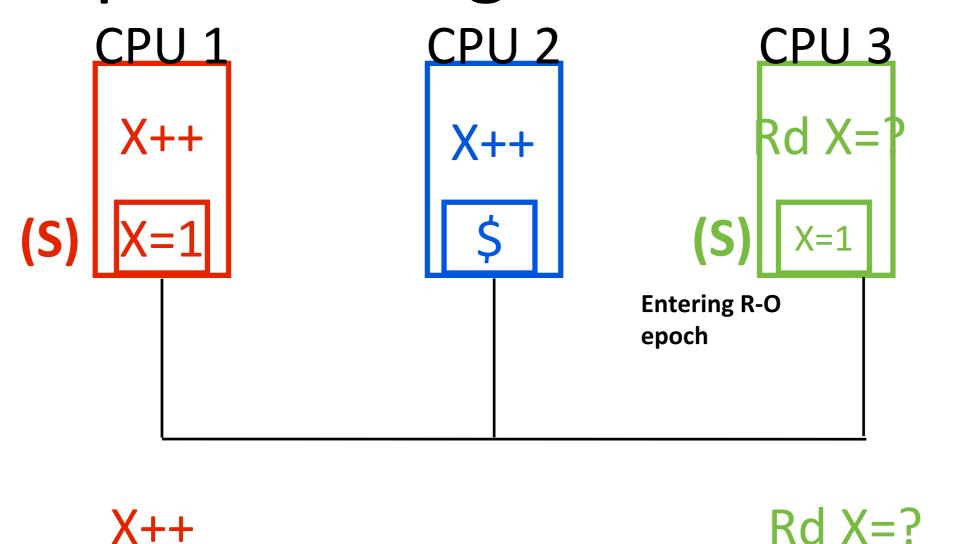


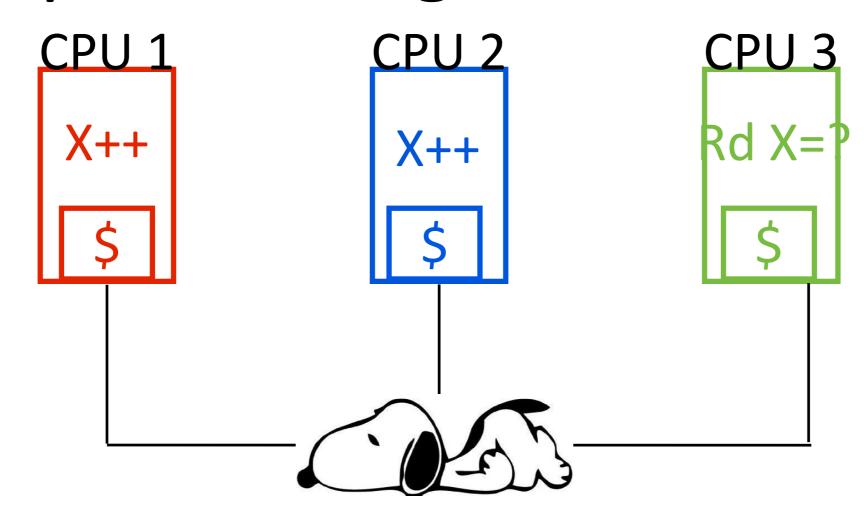




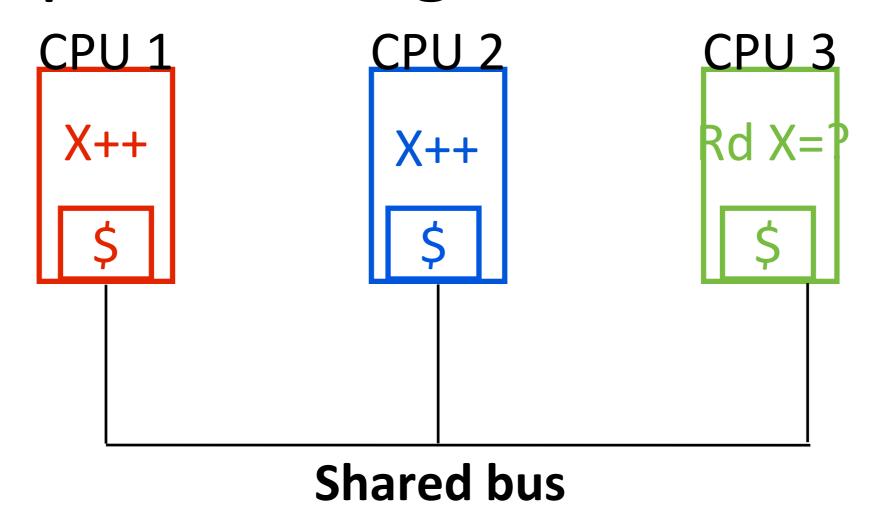








What sucks about Snoopy?



Bus limits scalability due to congestion and complex message arbitration

Figure 1-1. Uncore Sub-system Block Diagram of Intel Xeon Processor E5-2600 Family

(System Config

Controller)

Core 7

Core 6

Core 5

Core 4

Pbox

(Physical

Layer)

PCle x16

PCle x16

LLC

(Slice of 20M Last

Level Cache)

LLC

(Slice of 20M Last

LLC

(Slice of 20M Last

Level Cache)

LLC

(Slice of 20M Last

Level Cache)

Four Intel® SMI

Channels

IIO (Integrated IO)

(LLC Cohe

Engin

**CBox** 

(LLC Coheren

CBox

(LLC Coheren

CBox 4

Engine)

LLC Cohe

Engine)

Engine)

**PCU** 

(Power

Controller)

LLC

(Slice of 20M Last

Level Cache)

Pbox

Core 0

Core 1

Core 2

Core 3

QPI

(Packetizer)

CBox 0

CBox 1

CBox 2

CBox 3

LC Coherence

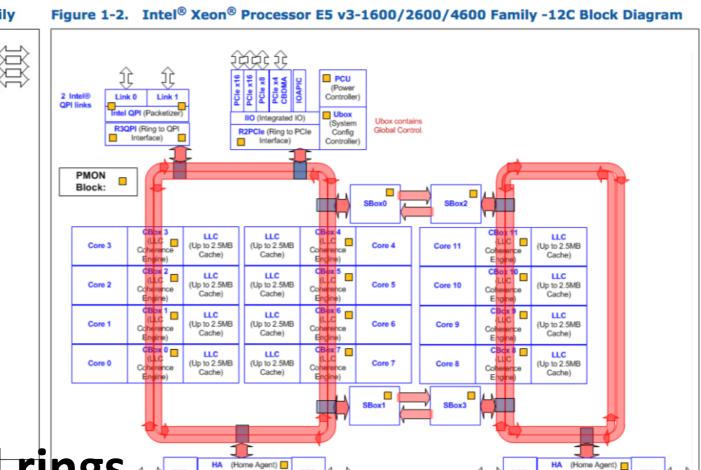
Coherence

.C Coherence

C Coherence

2 Intel® (=

QPI links (\( \)

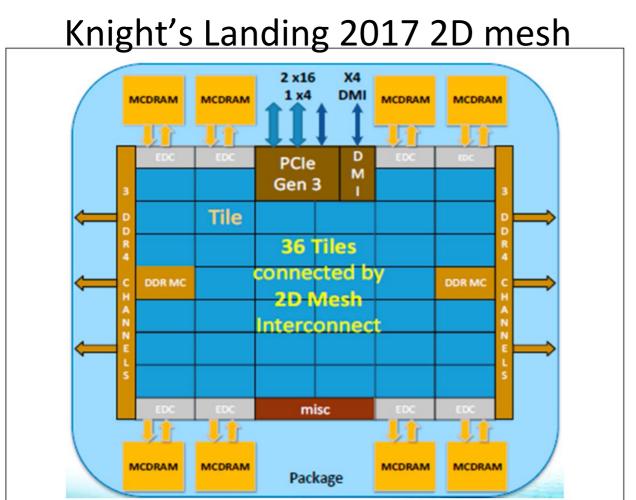


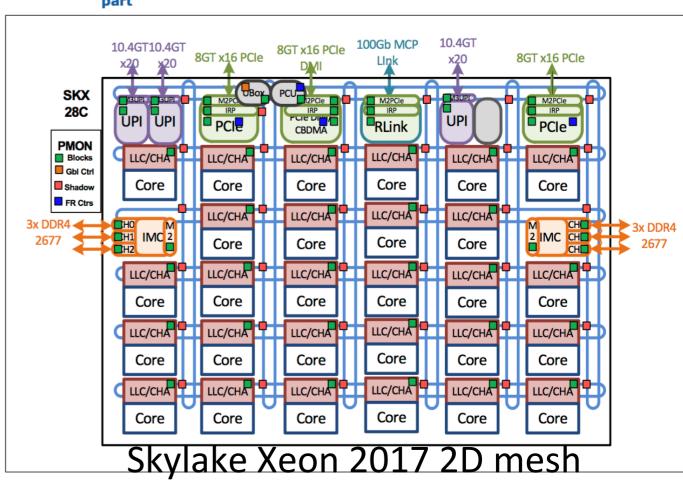
Sandybridge: bi-directional rings

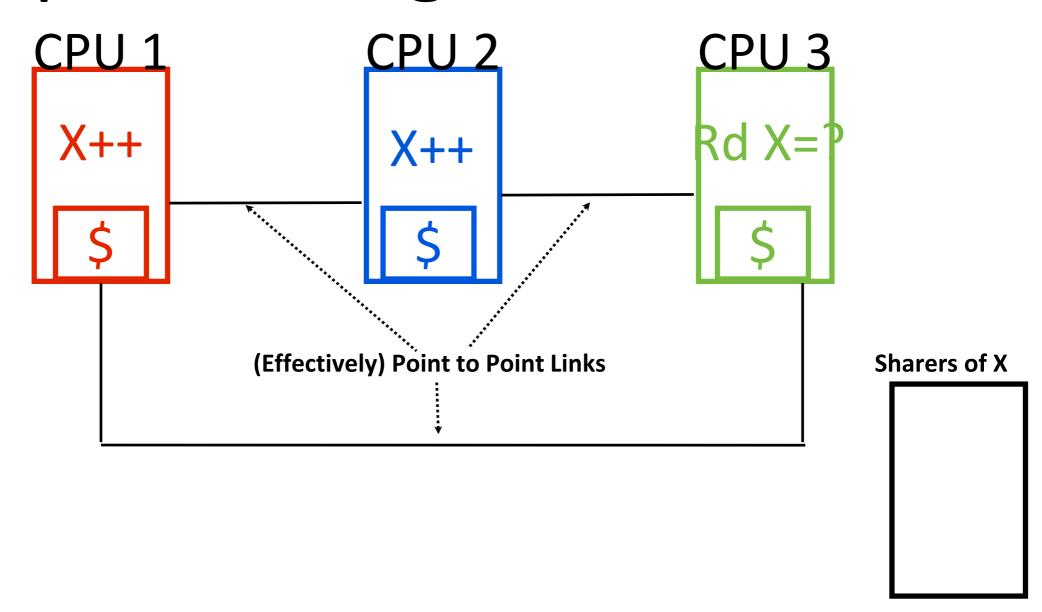
(Memory

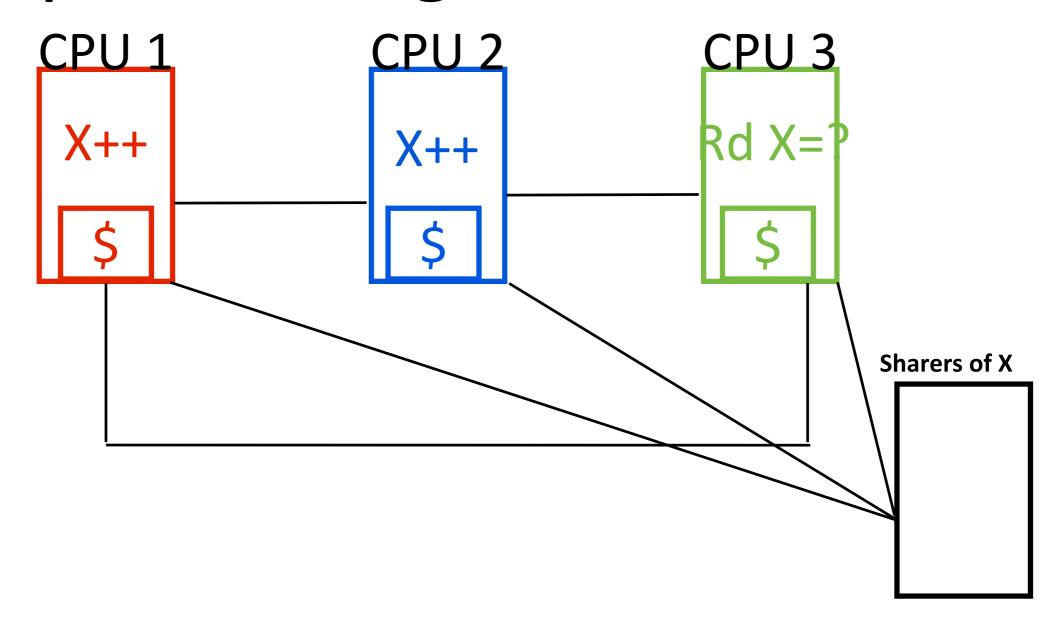
(Physical

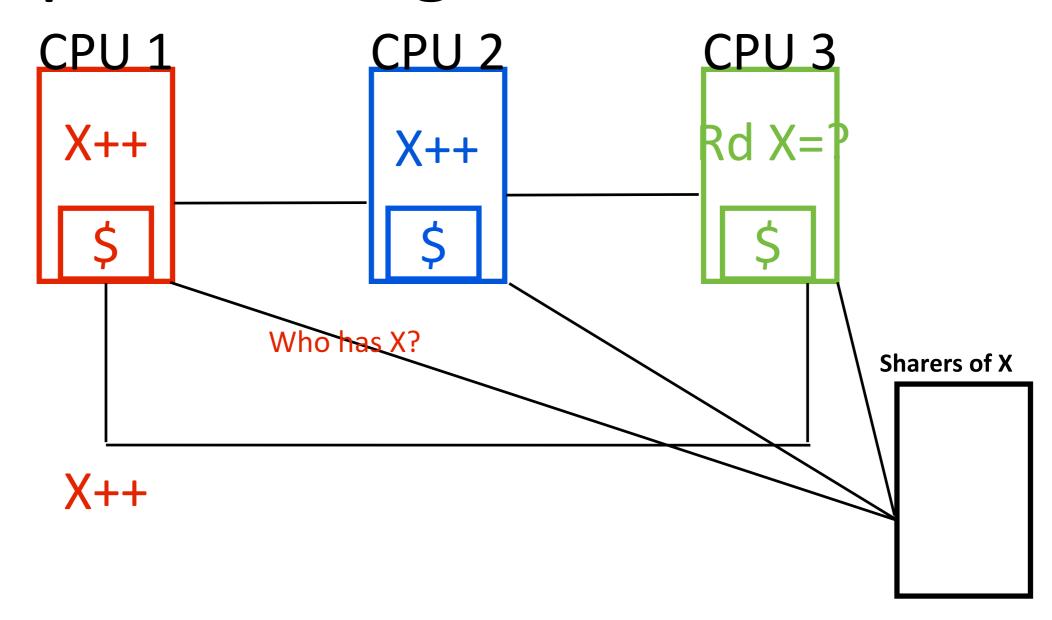
Figure 1-1. Intel® Xeon® Processor Scalable Memory Family - Block diagram for a 28C part

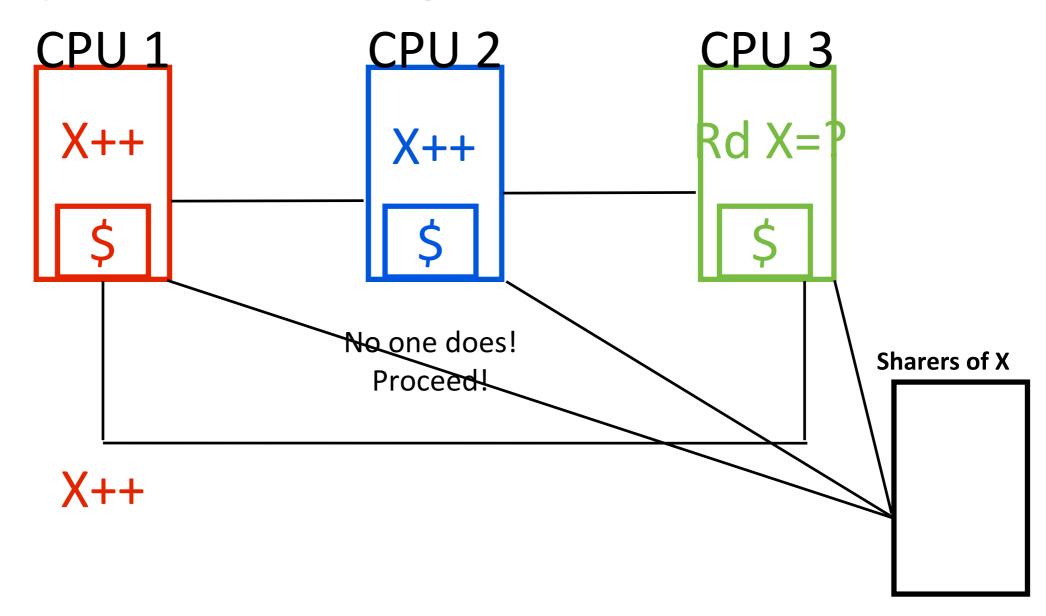


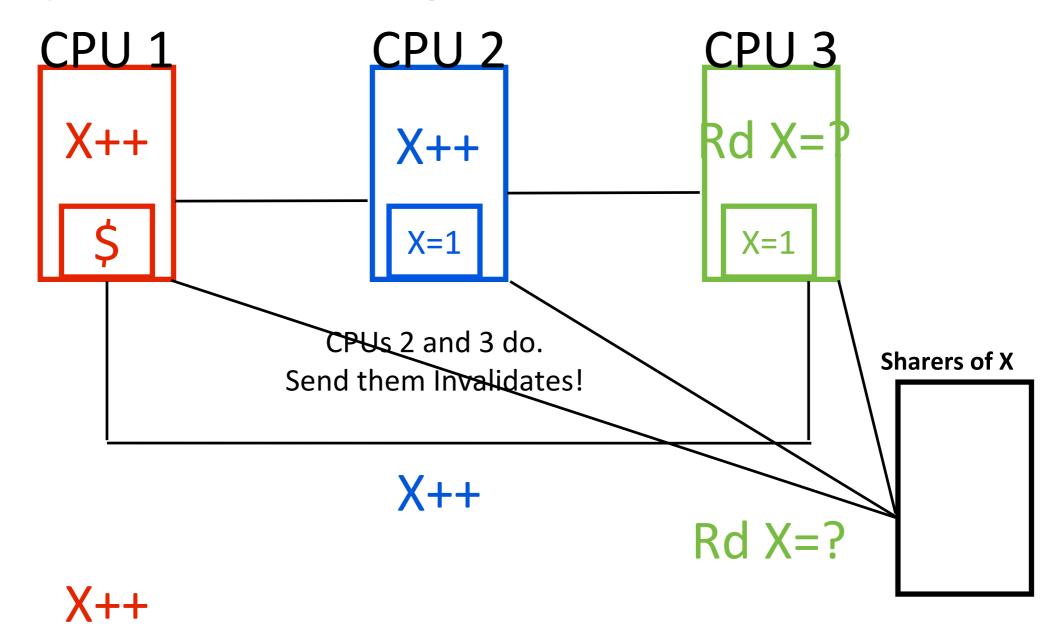


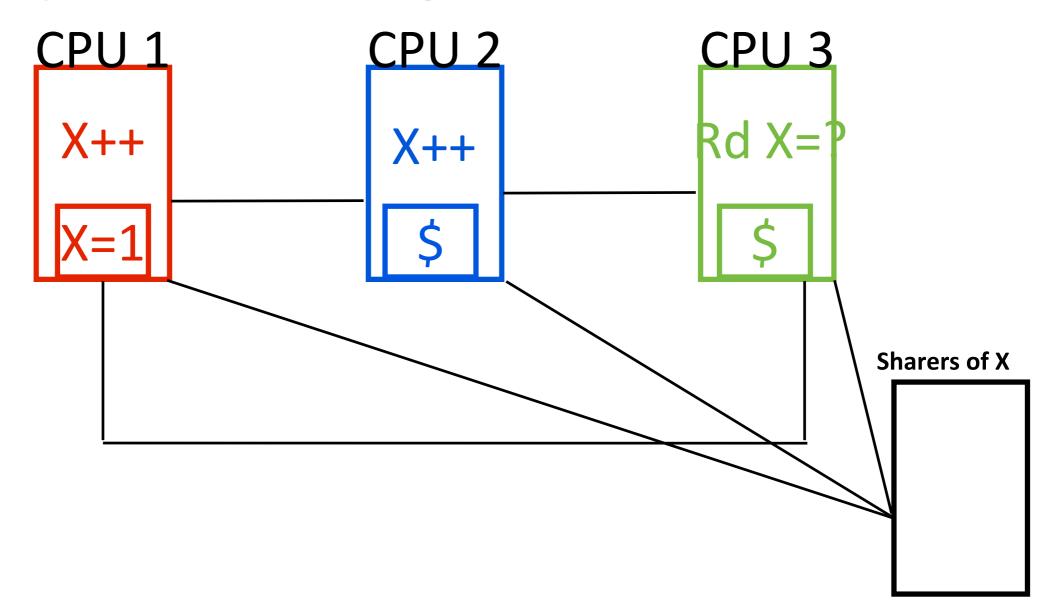




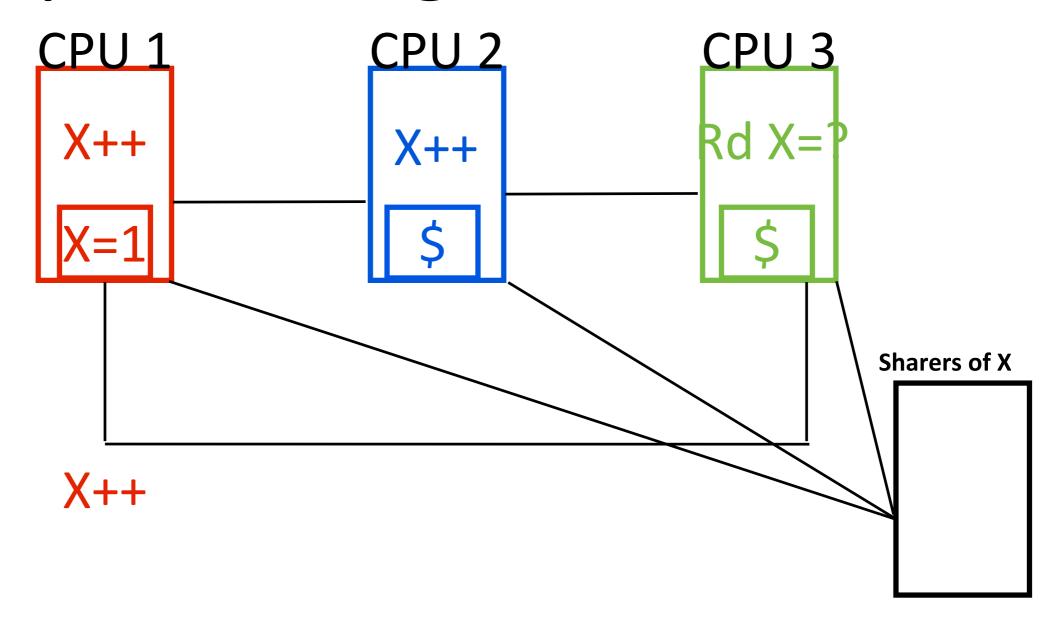




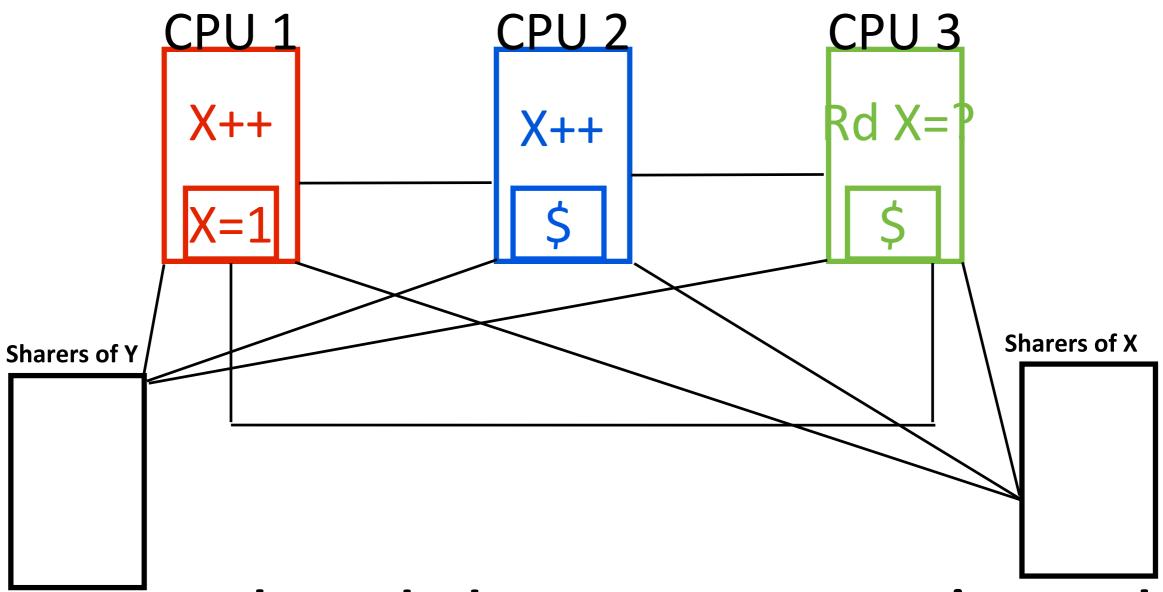




Benefit: No broadcast on shared bus

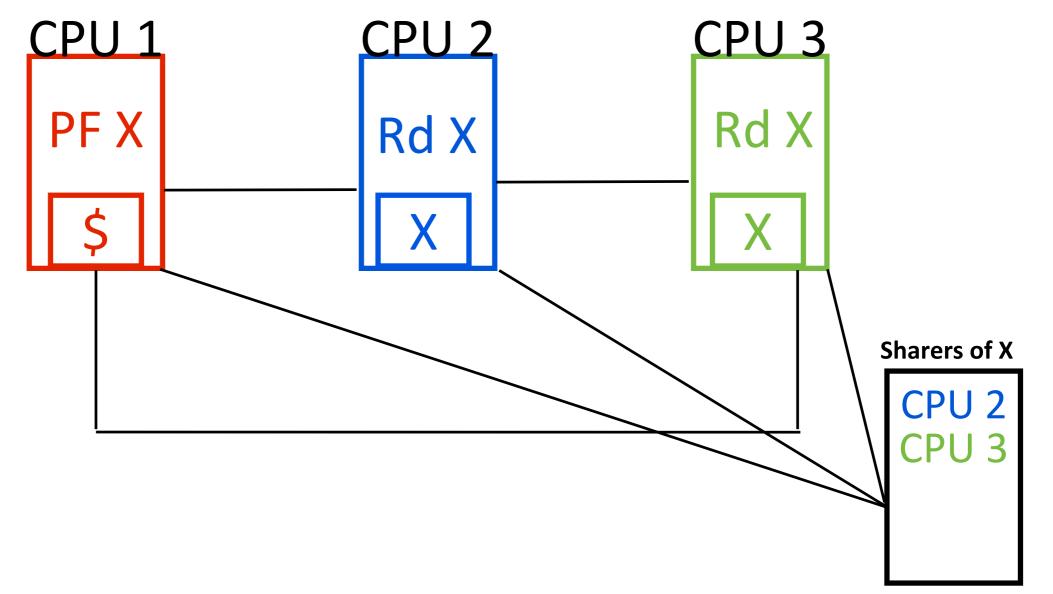


Drawbacks?



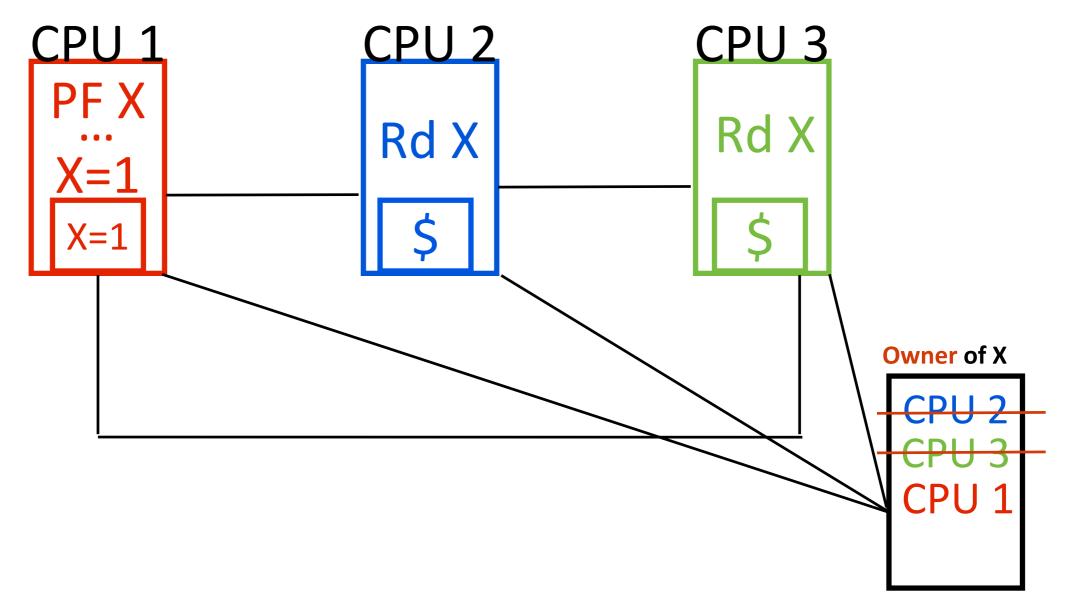
Centralized directory won't scale (In Practice: Distribute Directory)

### Optimization: Non-binding Prefetch



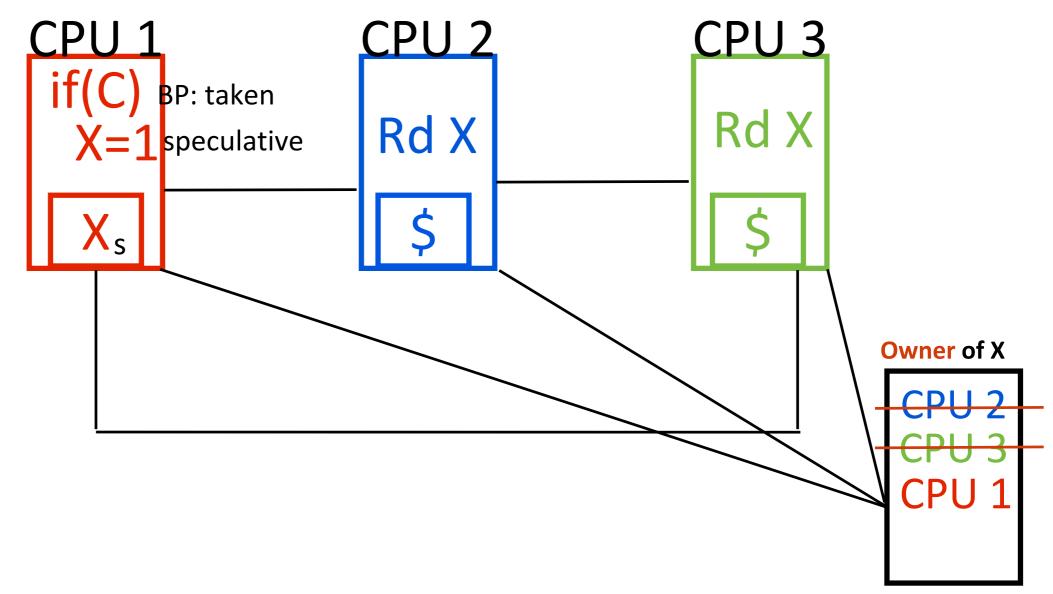
Prefetch instruction preemptively changes coherence state

### Optimization: Non-binding Prefetch

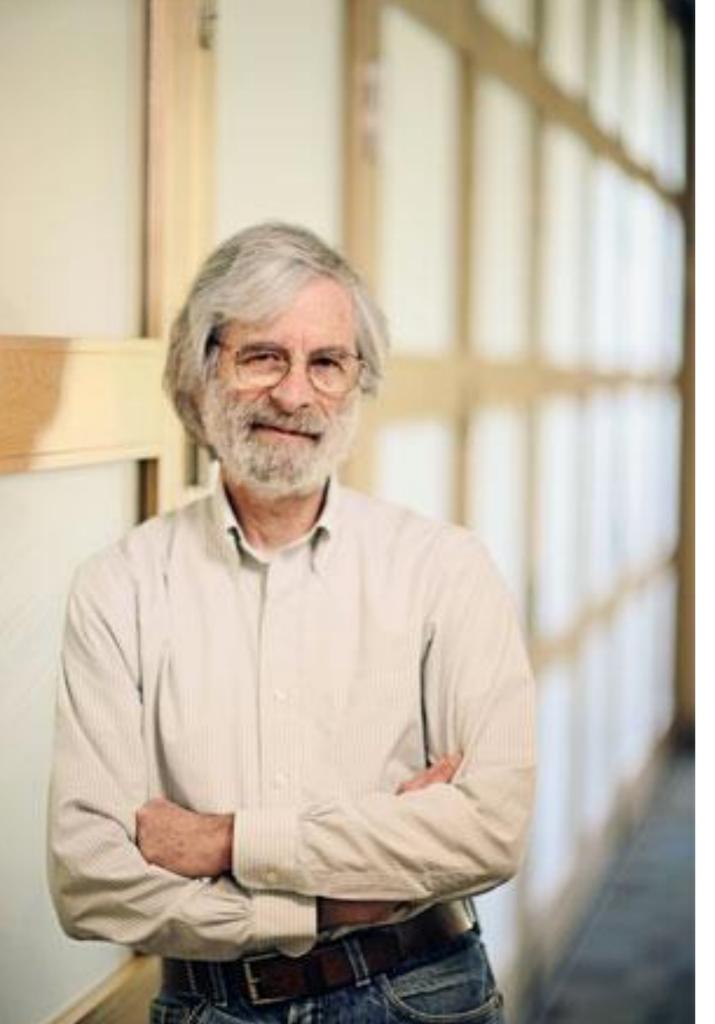


Benefit?

### Optimization: Speculation



Speculative operations that squash behave like non-binding pre-fetch



"computers execute operations in a different order than is specified by the program. A correct execution is achieved if the results produced are the same as would be produced by executing the program steps in order. For a multiprocessor computer, such a correct execution by each processor does not guarantee the correct execution of the entire program."

Excerpt from "How to Make a Multiprocessor Computer That Correctly Executes Multiprocess Program"
LESLIE LAMPORT, 1979

## Memory Consistency

# Memory Consistency Model

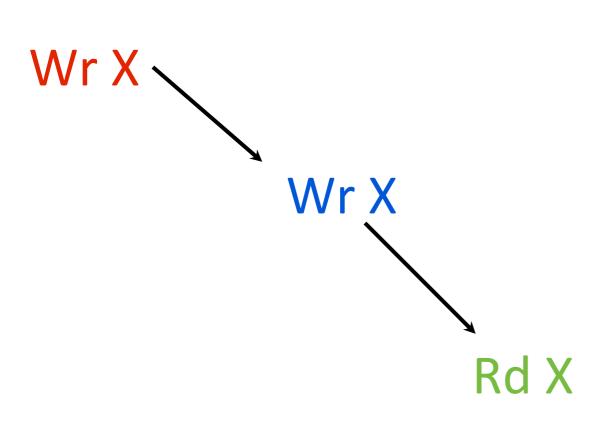
Informal Definition:

"Defines the value a read operation may read at each point during the execution"

"Defines the set of legal observable orders of memory operations during an execution"

"Defines which reorderings of memory operations are permitted"

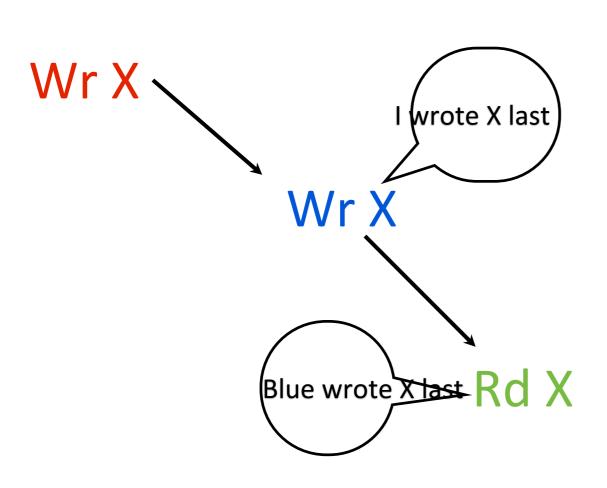
### Review: Coherence



#### 2 Invariants:

- "One Writer or
   One or More Readers"
- 2) "Reading X gets the value Rd X of the last write to X"

### Review: Coherence

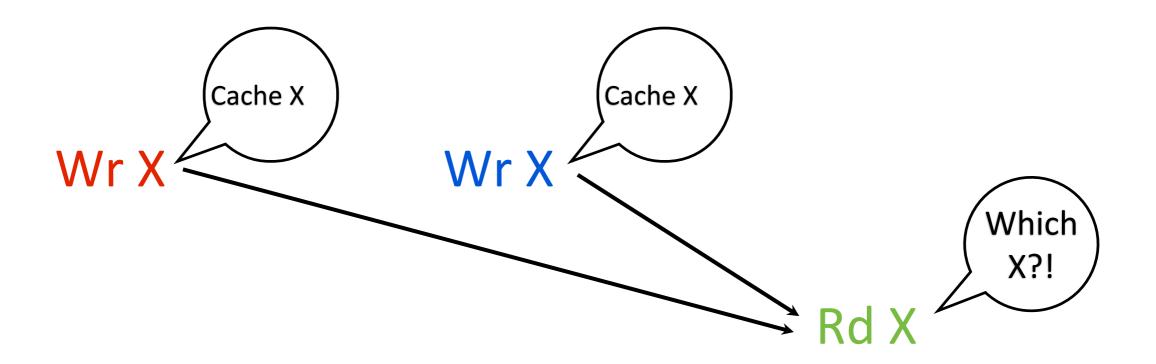


#### 2 Invariants:

- "One Writer or
   One or More Readers"
- 2) "Reading X gets the value of the last write to X"

### Without Coherence

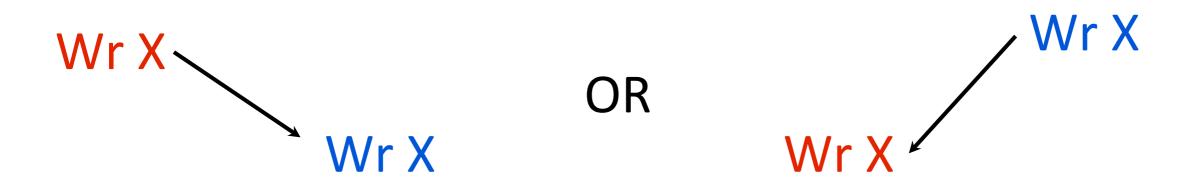
(The coherence invariants prevent this from happening)



Processors can't decide who wrote last.

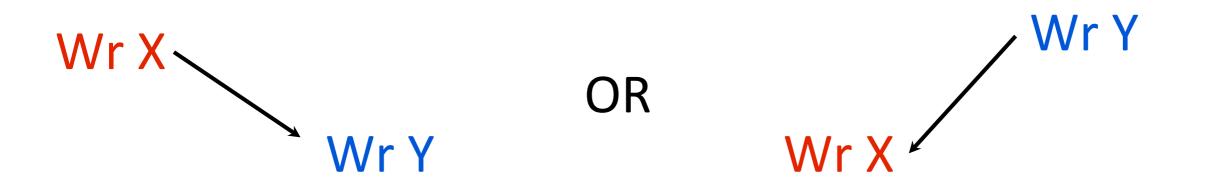
Green is hosed.

## Coherence is Ordering



Coherence defines the set of legal orders of accesses to a single memory location

## Consistency is Ordering



Consistency defines the set of legal orders of accesses to multiple memory locations

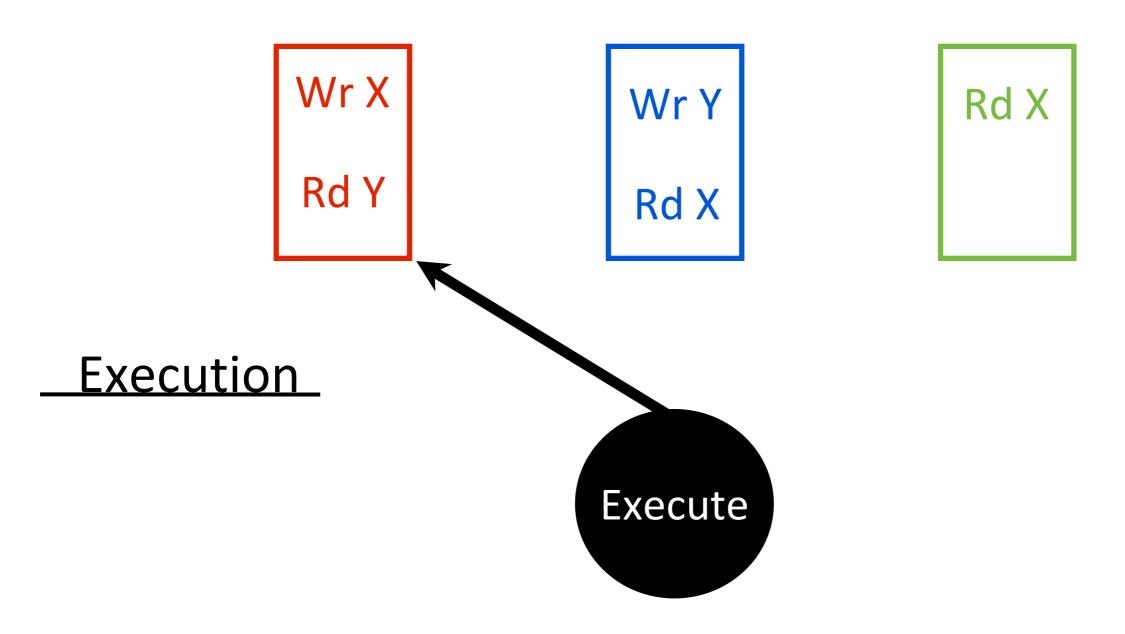
### Sequential Consistency (SC)

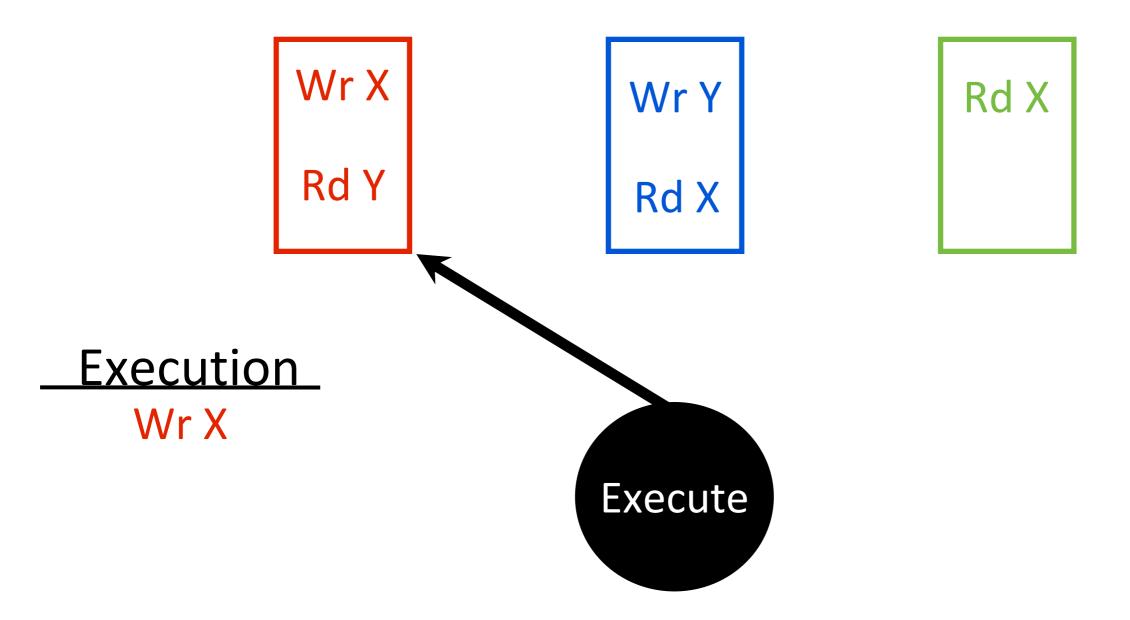
The simplest, most intuitive memory consistency model

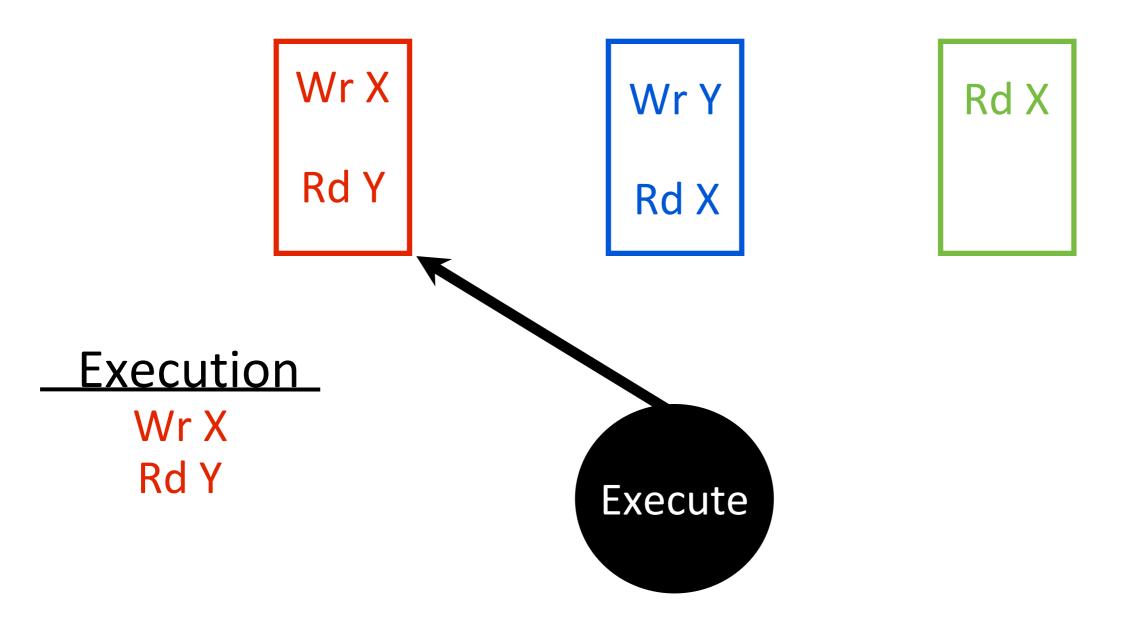
Two Invariants to SC:

Instructions are executed in program order

All processors agree on a total order of executed instructions





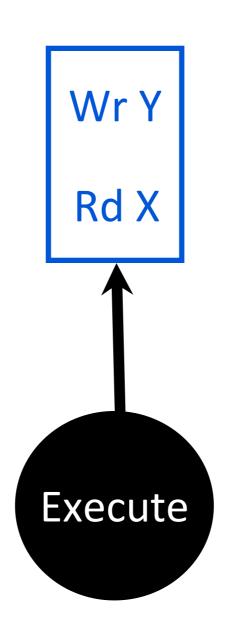


Wr X Rd Y

**Execution** 

Wr X Rd Y

WrY



Rd X

Wr X Rd Y

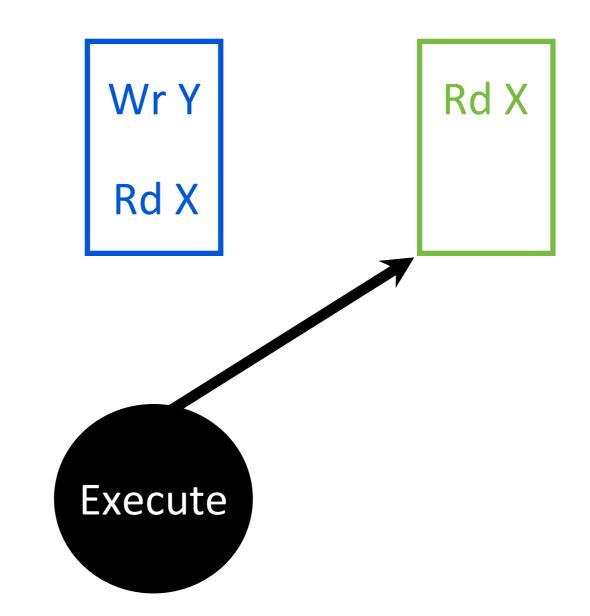
**Execution** 

Wr X

Rd Y

WrY

Rd X







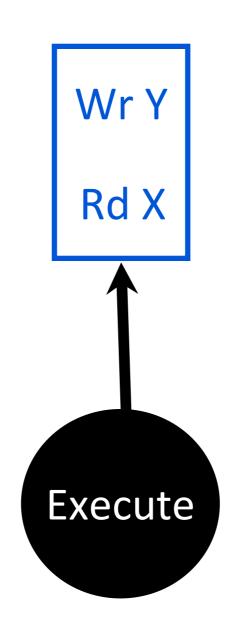
Wr X

Rd Y

WrY

Rd X

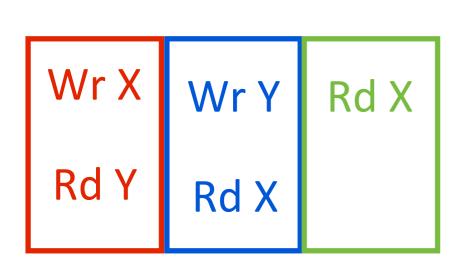
Rd X



Rd X

# Who cares?.... You care!

SC is the most complex model that we can ask **programmers** to think about.



SC prohibits all reordering of instructions (Invariant 1)

## Why are Instructions Reordered?

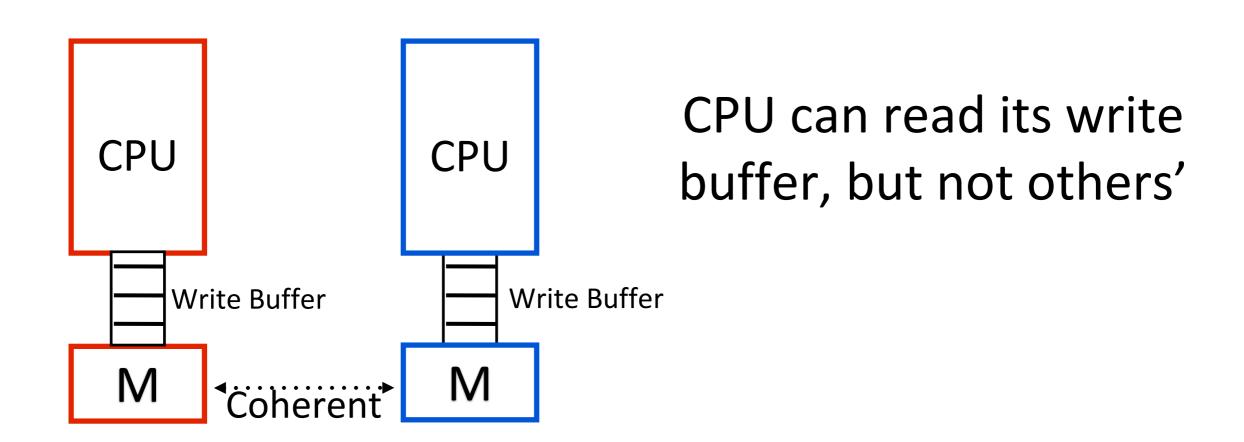
And when does it matter?

## Why are Instructions Reordered?

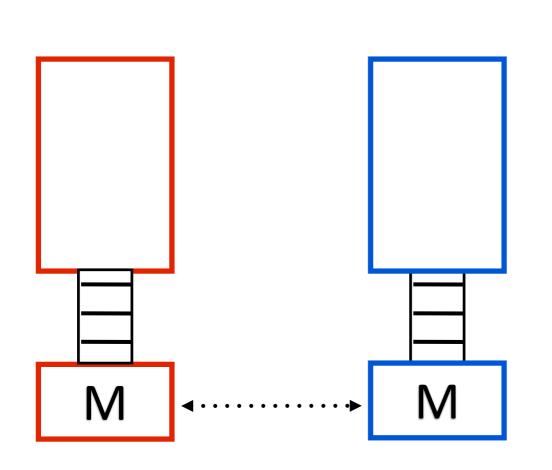
Optimization.

Elsewhere?

### Reordering #1: Write Buffers

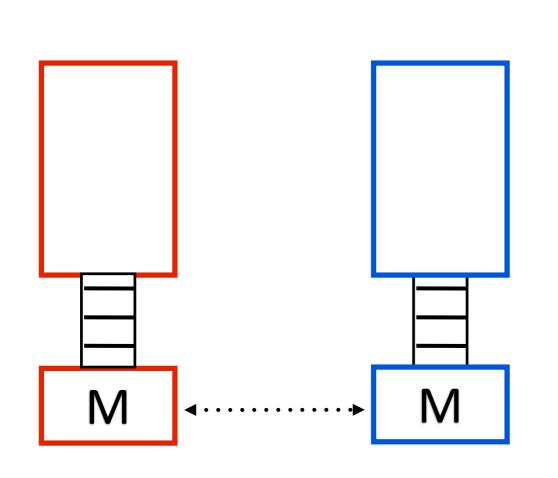


Buffered writes eventually end up in coherent shared memory



 $\frac{Program}{Initially X == Y == 0}$   $X=1 \qquad Y=1$ 

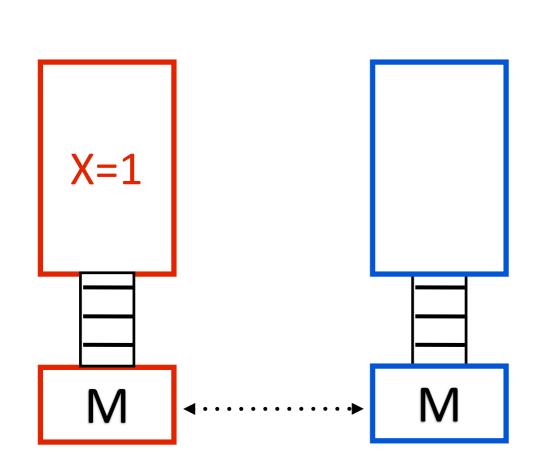
Is r1==r2==0 a valid result?



Program
Initially X == Y == 0  $X=1 \qquad Y=1$   $r1=Y \qquad r2=X$ 

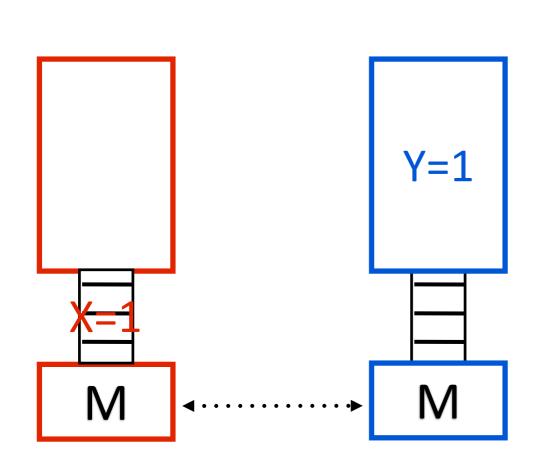
Is r1==r2==0 a valid result?

r1 == r2 == 0 is **not** SC, but it can happen with write buffers



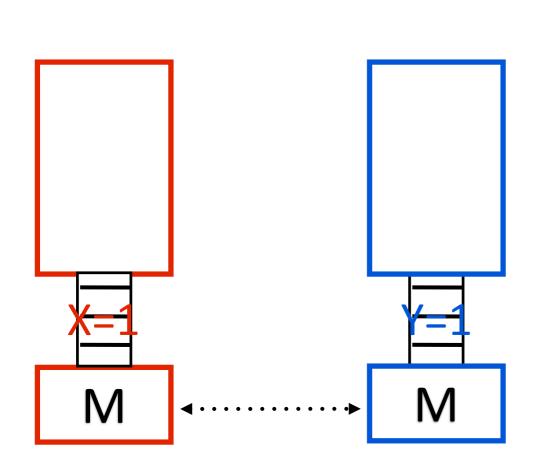
 $\frac{Program}{Initially X == Y == 0}$ 

Execution



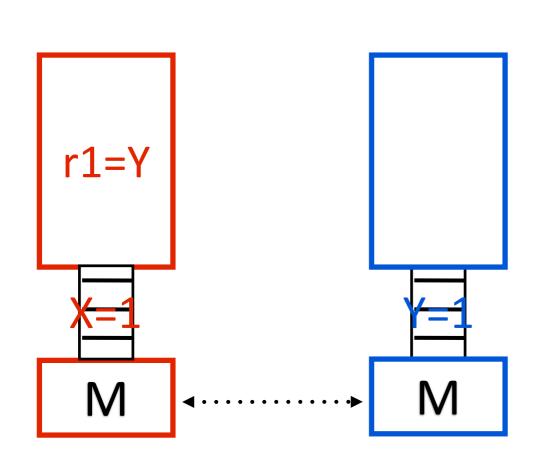
 $\frac{Program}{Initially X == Y == 0}$ 

r1=Y r2=X



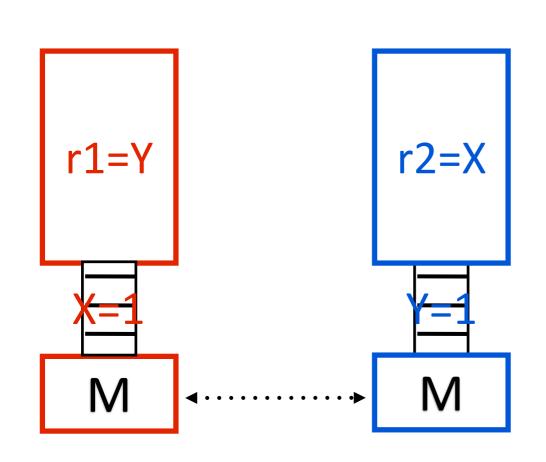
 $\frac{Program}{Initially X == Y == 0}$ 

r1=Y r2=X

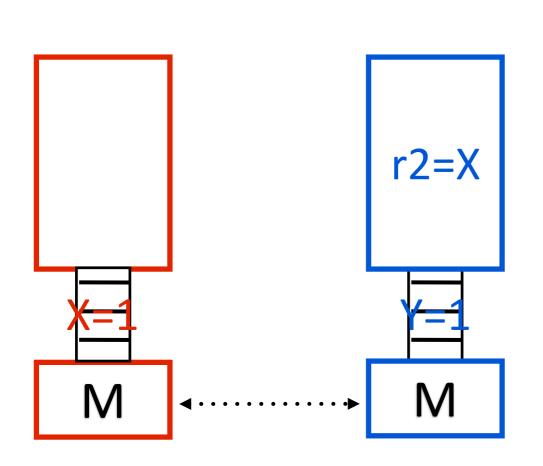


 $\frac{Program}{Initially X == Y == 0}$ 

r2=X

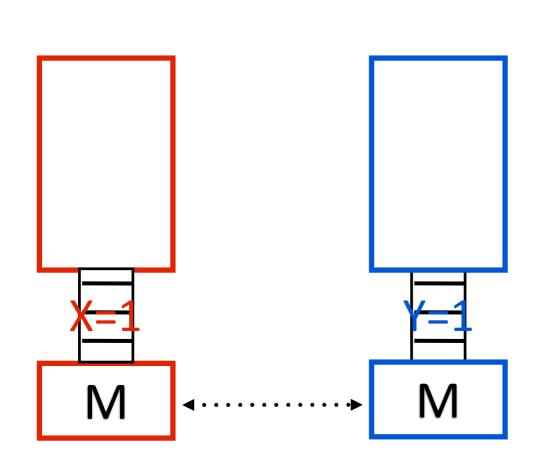


 $\frac{Program}{Initially X == Y == 0}$ 



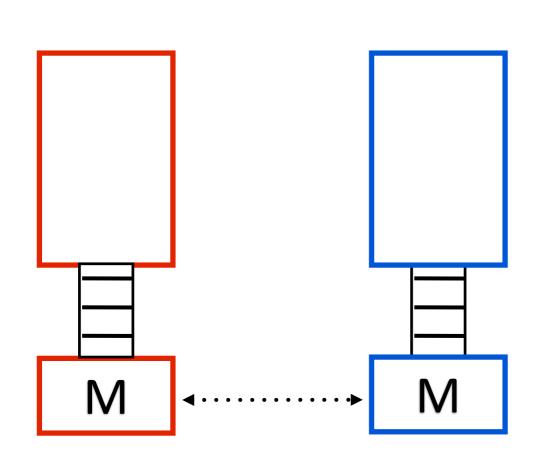
 $\frac{Program}{Initially X == Y == 0}$ 

<u>Execution</u> r1=Y [r1 <- 0]



 $\frac{Program}{Initially X == Y == 0}$ 

Execution r1=Y [r1 <- 0] r2=X [r2 <- 0]

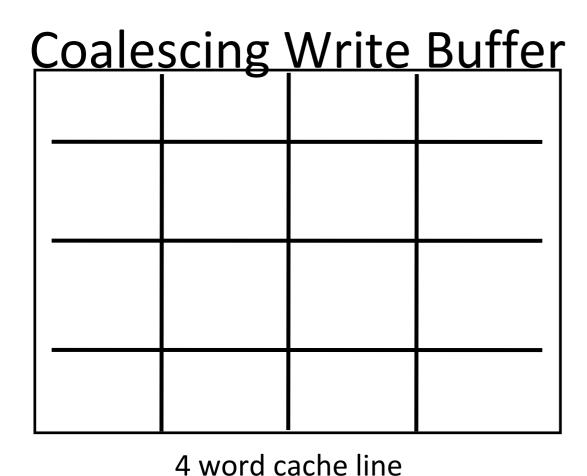


WBs let reads finish before older writes

 $\frac{Program}{Initially X == Y == 0}$ 

Execution

```
r1=Y [r1 <- 0]
r2=X [r2 <- 0]
X=1
Y=1 (Not SC!)
```



Program
X,Z in same \$ line
X=1

Y=1 Z=1

**Coalescing Write Buffer** 

X=1		

Program
X,Z in same \$ line

X=1

Y=1

Z=1

Coalescing Write Buffer

X=1		
		Y=1

Program

X,Z in same \$ line

X=1

Y=1

Z=1

Coalescing Write Buffer

X=1		
		Y=1
	Z=1	

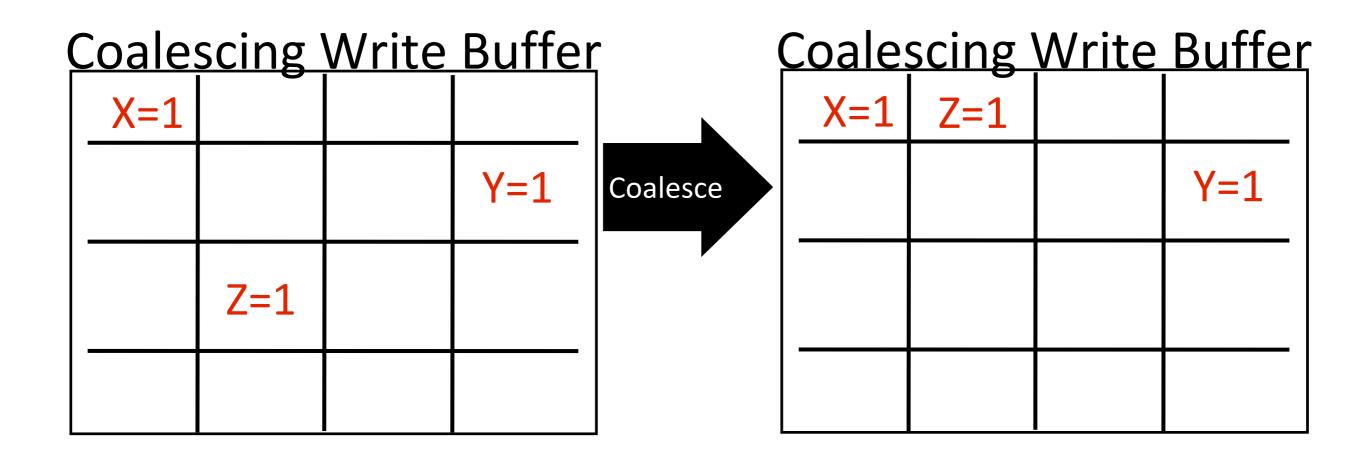
Program

X,Z in same \$ line

X=1

Y=1

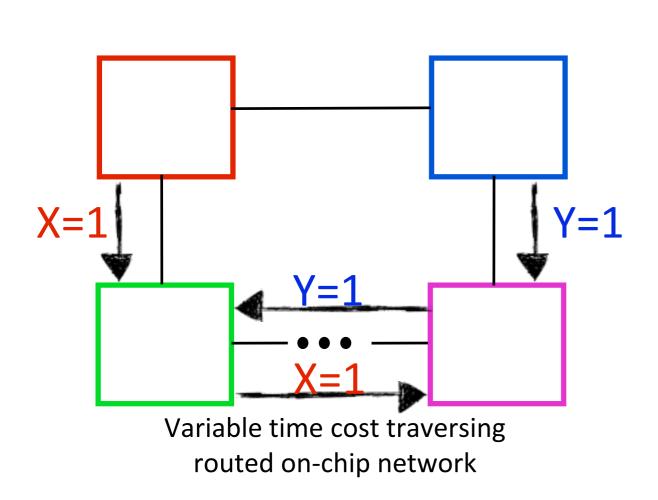
Z=1



Combining the write to X & Z saves bandwidth, but **reorders** Z=1 and Y=1

## Reordering #3: Interconnect

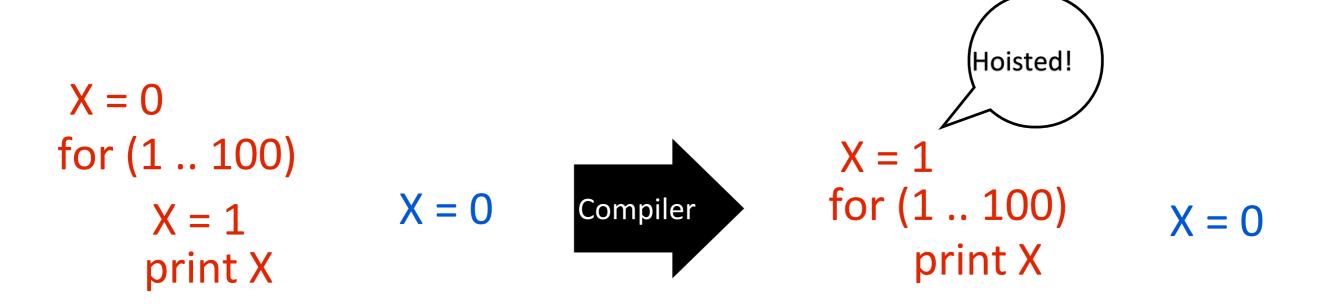
X=1



#### <u>Program</u>

#### **Execution**

## Reordering #4: Compilers



The compiler hoists the write out of the loop, permitting new (non-SC) results (e.g., "1 0 0 0 0 0 0...")

## When is Reordering a Problem?

When Executions Aren't SC

When a memory operation happens before itself

#### <u>Execution</u>

#### Happens-Before Graph

$$X=1$$
  $Y=1$ 

$$r1=Y$$
  $r2=X$ 

When a memory operation happens before itself

## 

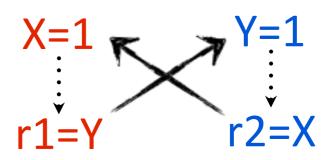
#### Happens-Before Graph

Program Order HB Edge

When a memory operation happens before itself

#### <u>Execution</u>

#### Happens-Before Graph



Program Order HB Edge

↓ Causal Order HB Edge

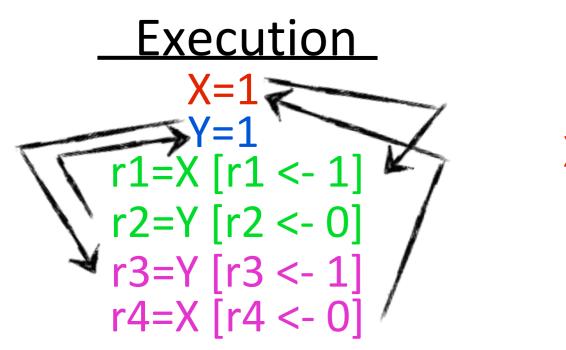
When a memory operation happens before itself

# 

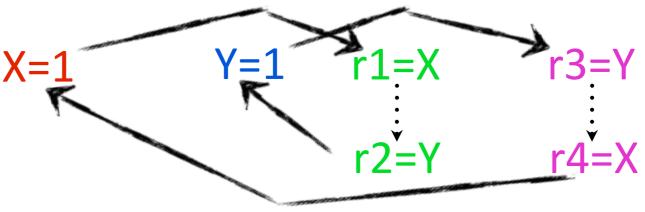
Happens-Before Graph

If there is a cycle in the happens-before graph, the execution is not SC

When a memory operation happens before itself



Happens-Before Graph



If there is a cycle in the happens-before graph, the execution is not SC

## So... are Computers Wrong?!

SC is how programmers think.

SC prohibits all reordering of instructions

WBs let reads finish before older writes

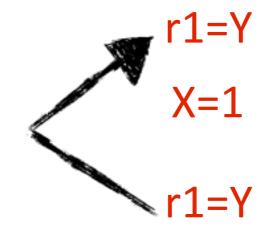
Combining writes saves bandwidth but reorders writes

## Relaxed Memory Consistency

Relaxed Memory Models permit reorderings, unlike SC

## **x86-TSO** (intel x86s)

"The Write Buffer Memory Model"

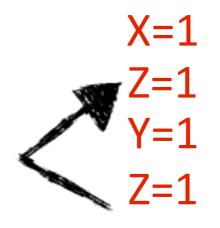


Relaxes W->R order

Total Store Order - loads may complete before older stores to different locations complete.

## PSO(SPARC)

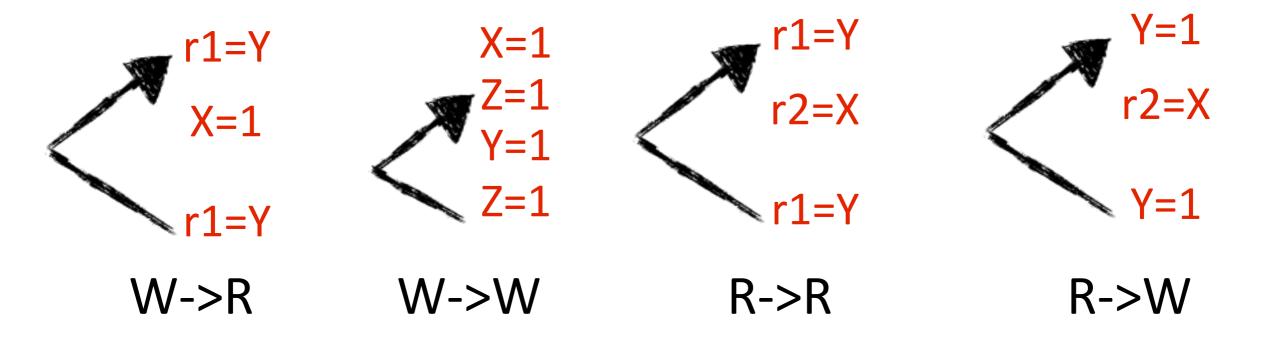
"The Write Combining Memory Model"



Relaxes W->W order

Partial Store Order - loads and stores may complete before older stores to different locations complete.

### In General



Starting with PSO and relaxing R->R and R->W yields Weak Ordering or Release Consistency (alpha)

Depending on the implementation

## SC and Relaxed Consistency

SC is required for correctness and programmer sanity

+

Reordering is required\* for performance

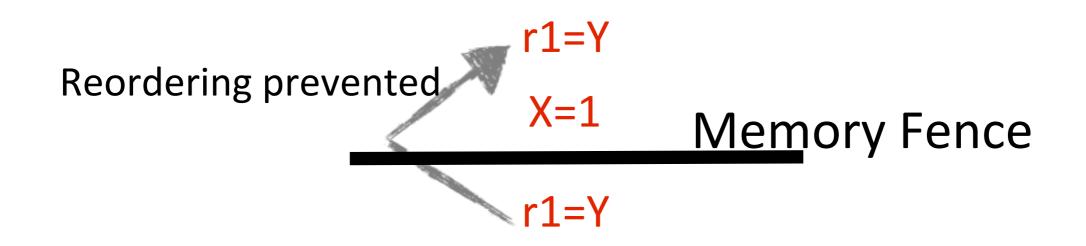
Goal: Ensure SC executions while permitting Relaxed Consistency reorderings

\*Usually; the MIPS memory model is **SC** (surprising!)

## How to ensure SC, but permit reordering?

## Synchronization Prevents Reordering

Memory fences are another type of synchronization

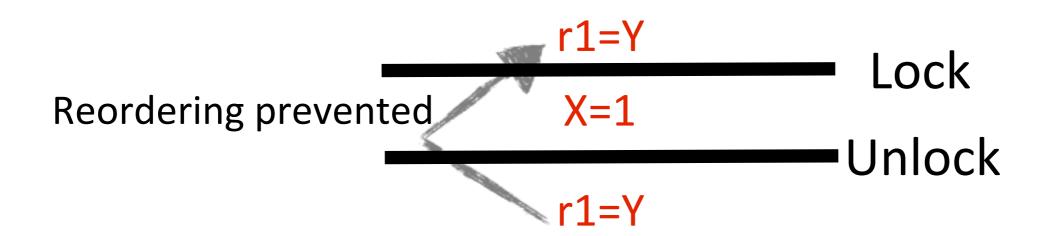


Fence implementation depends on reordering implementation

TSO: Stall reads until write buffer is empty

## Synchronization For Real Programmers

Memory fences are wrapped up in locks, etc.

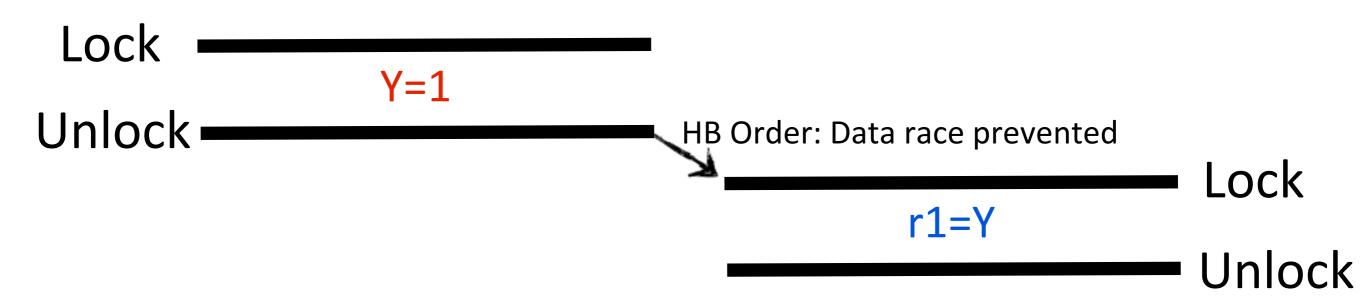


Direct use of fences possible, but inadvisable.

USE A LIBRARY.

#### Data Races

Synchronization imposes happens-before on otherwise unordered operations



Data Race: Unordered operations to the same memory location, at least one a write

## Memory Models across the System Stack

#### Language

Java/C++: SC for data-race-free programs

#### Compiler

Conservative with reordering when d-r-f can't be proved

#### Architecture

Usually very weak for max optimization (lots of reordering)

Note: fences from "above" ensure SC