✓ Caches →
→ Multicore protocols
Q+k

CPI/AMAT

- Trace 2
→ Trace 1
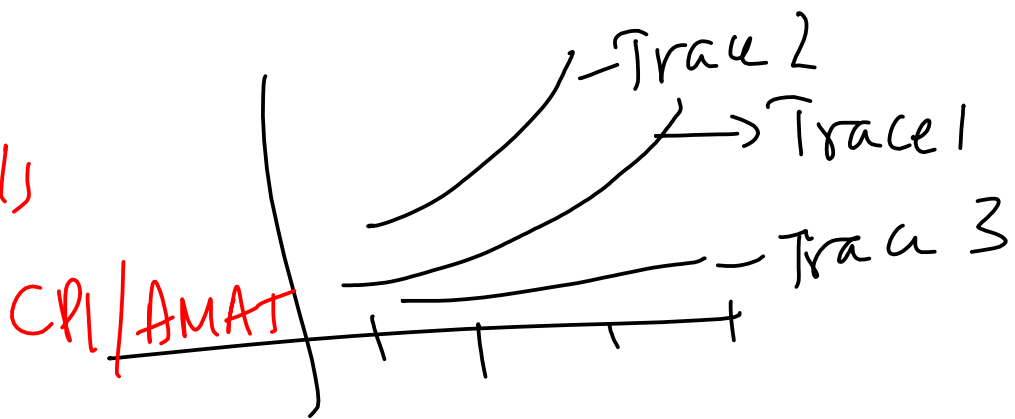- Trace 3

→ BPrediction

→ Exceptions

→ Unrolling

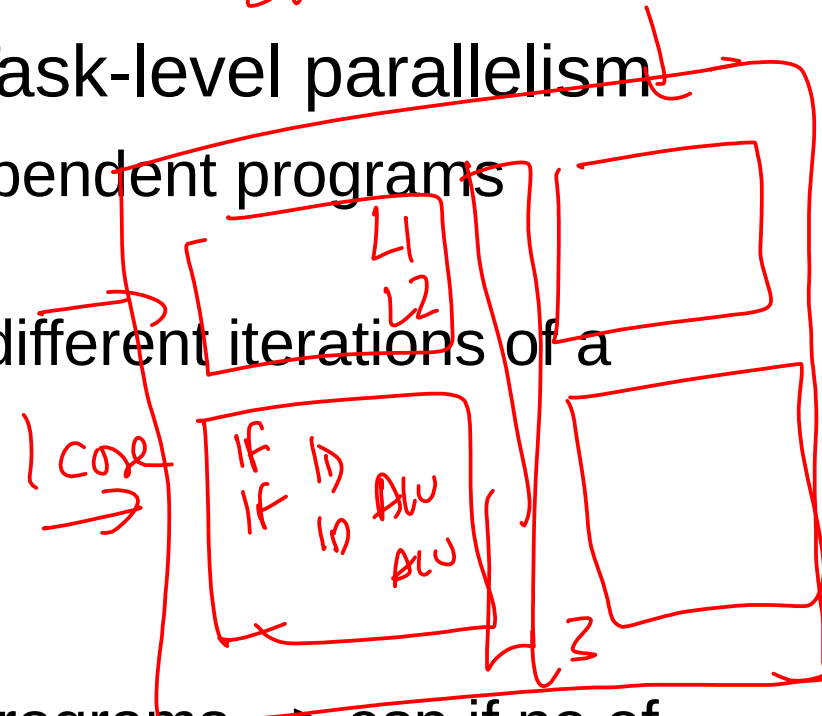→ Virtual memory

# Multicore architecture
# Cache coherency

CA- QA – H&P – Ch 5

# Motivation

- Better performance
- ILP --> Thread level parallelism/ Task-level parallelism
  - Use multiple processors to run independent programs simultaneously
  - Or tasks within the same program- different iterations of a loop
- Multicore processor
  - Multiple cores on the same die
  - Difficult to create parallel software programs --> esp if no of cores increase
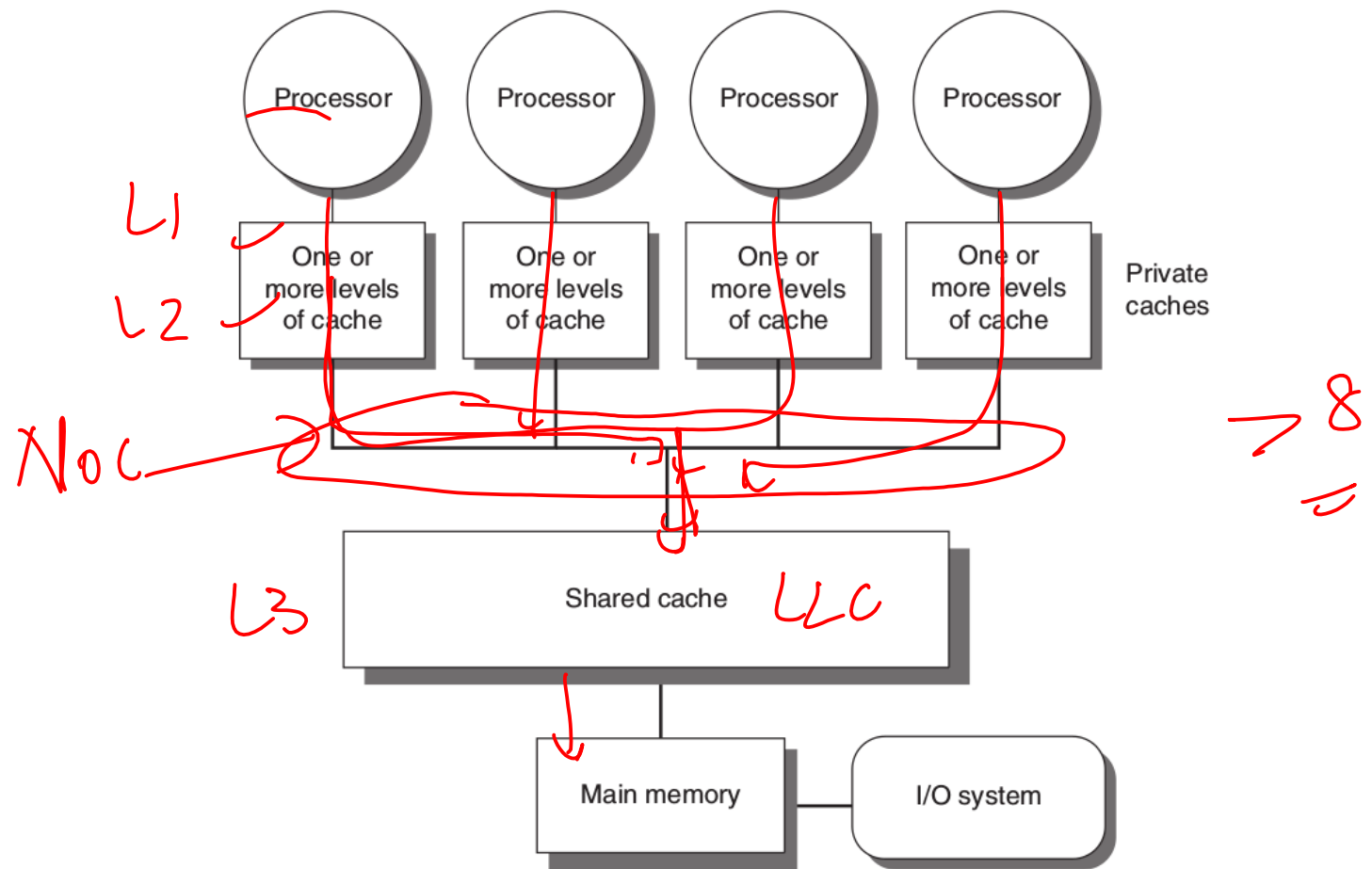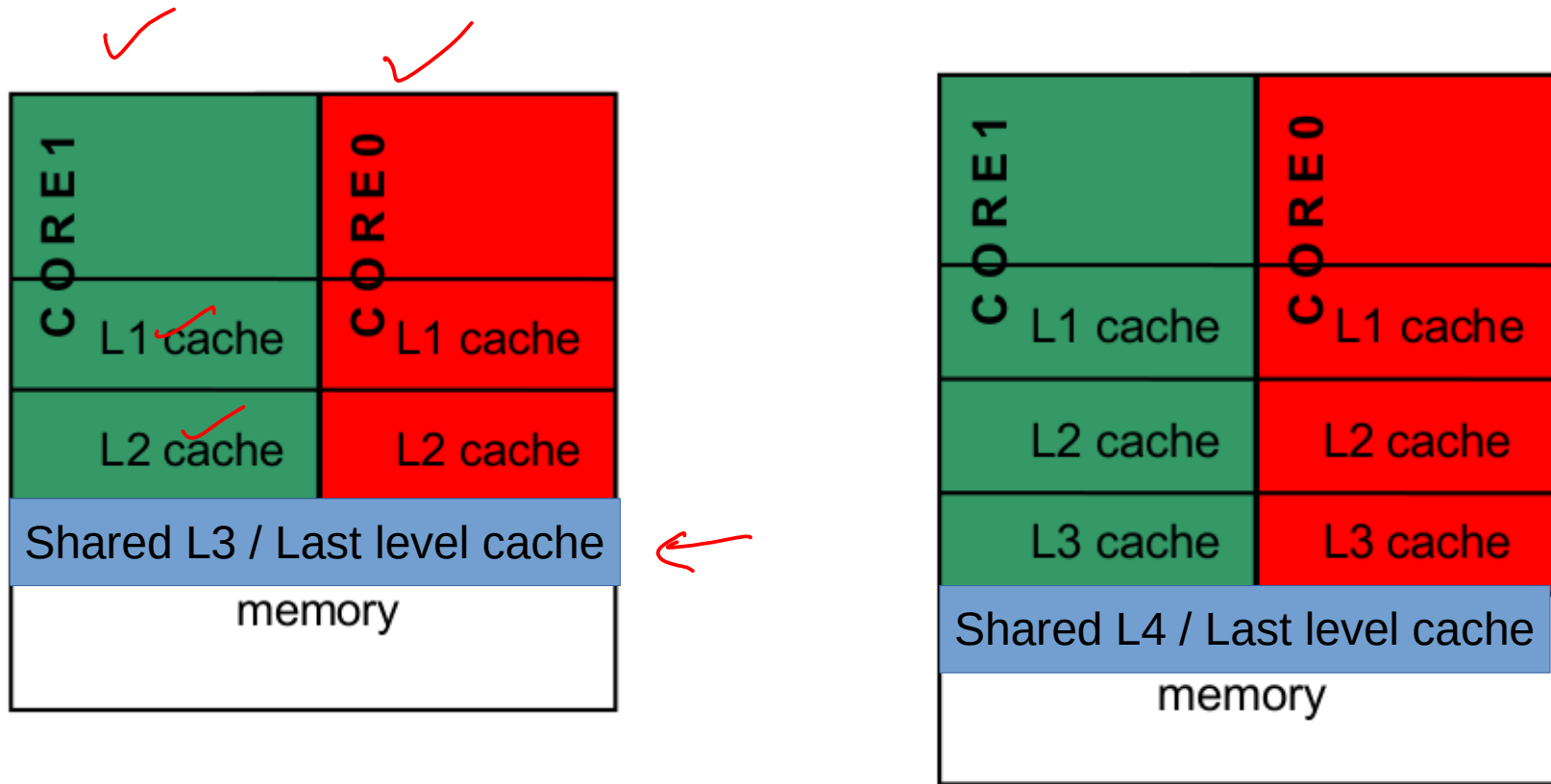
# Challenges

- Superscalar and OOO on a single processor use ILP --> No user intervention

- Multicore-->
  - Programs should be written to exploit parallelism
  - Tasks should be broken down into parallel small tasks
  - Balance the load
  - Communicate between tasks, synchronization
  - Compiler exploits data level parallelism and generates parallel threads
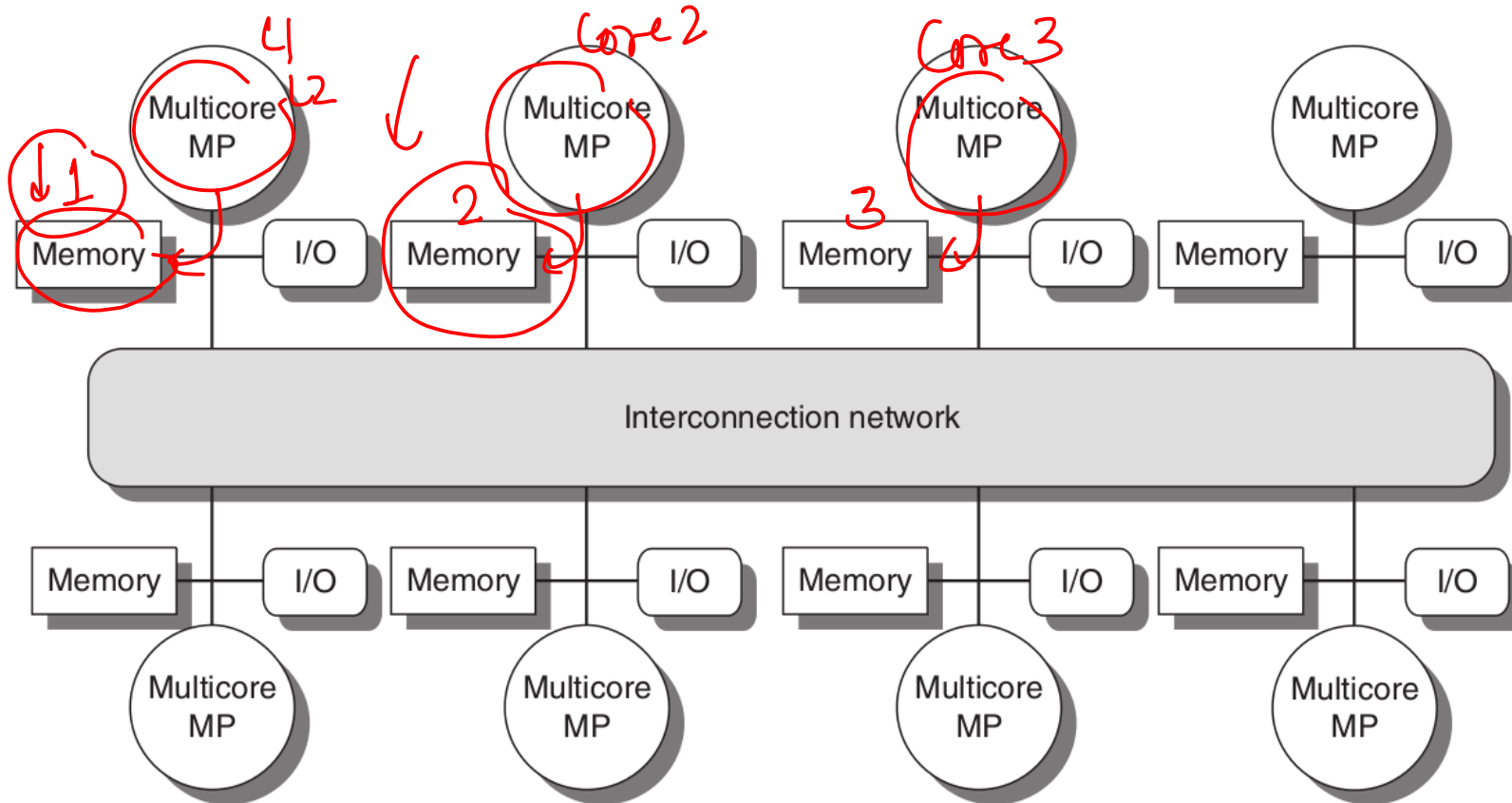
# Shared memory processors (SMP)



- Upto 8 cores
- Share same memory address space (centralized memory architectures)
  - Uniform memory architectures (UMA)
  - All processors have the same latency from memory

4

# Architectures



Local data: Local data used only by that core- can be in the local cache
Shared data: May be replicated in multiple caches

# Distributed memory processors



- More cores > 8. Share same memory address space.

- Individual memory with processors, to reduce long memory access latency with several cores

- Non-uniform memory access (NUMA)

- Focus on interconnect network design

- Communication between processors is a challenge

6

# Challenges in multi-cores

- Insufficient parallelism in programs
  - Write better algorithms or programs

- Long latency communications between processors
  - Architectural changes: Caches can store shared data: How do we maintain coherence between multiple copies
    - Local data (Not shared)
    - Shared data
  - Consistent view of memory
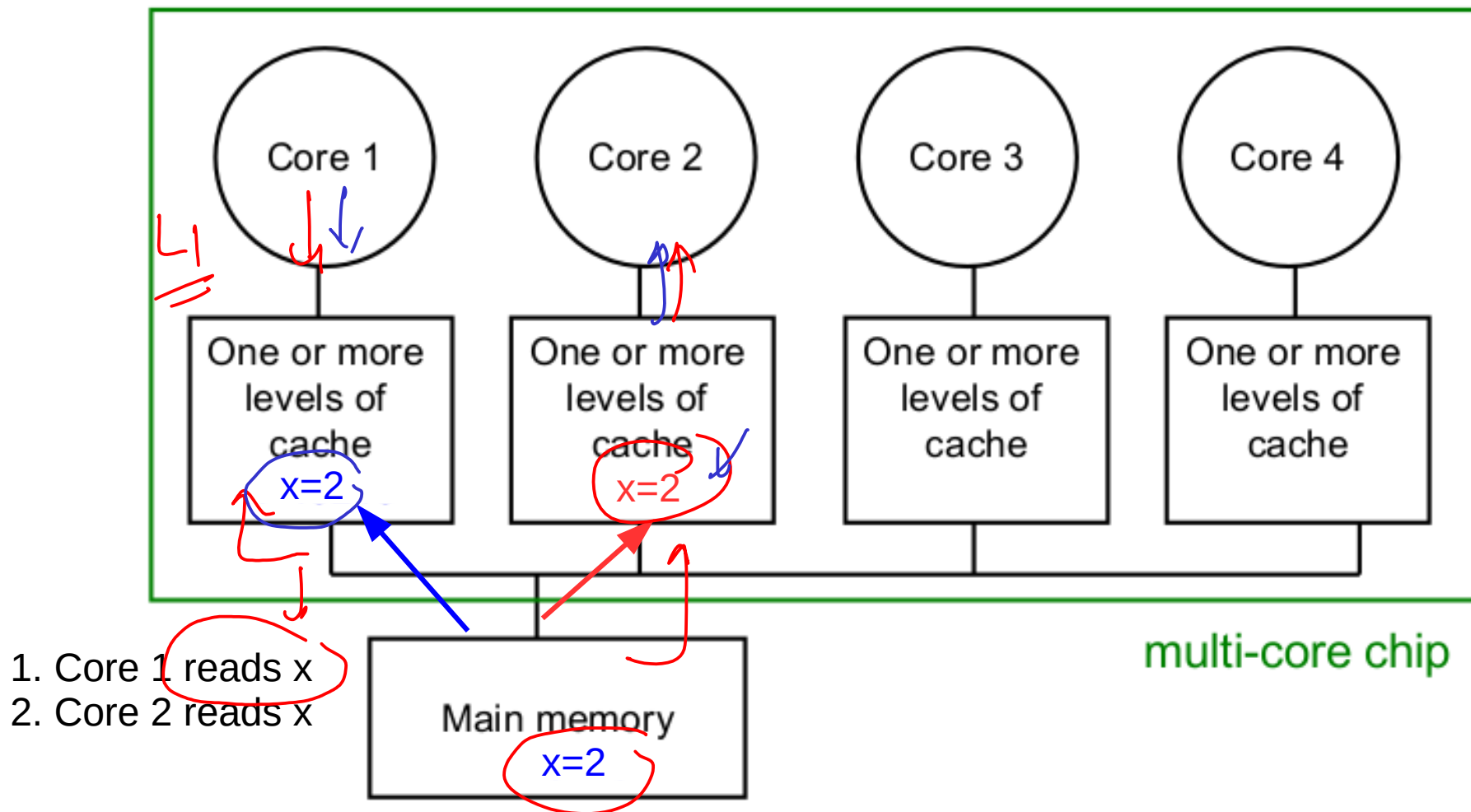  - Multi-threading
  - Prefetching

# Shared data

- Memory system is coherent if any read of a data item returns the most recently written value of that data item.

- Coherence: What is the value returned by a read

- Consistent: When is the value returned by the read

# Cache coherence in SMP
## Caching of shared data/address space

Reads



multi-core chip

1. Core 1 reads x
2. Core 2 reads x

Core 1    Core 2    Core 3    Core 4

One or more levels of cache    x=2

One or more levels of cache    x=2

One or more levels of cache
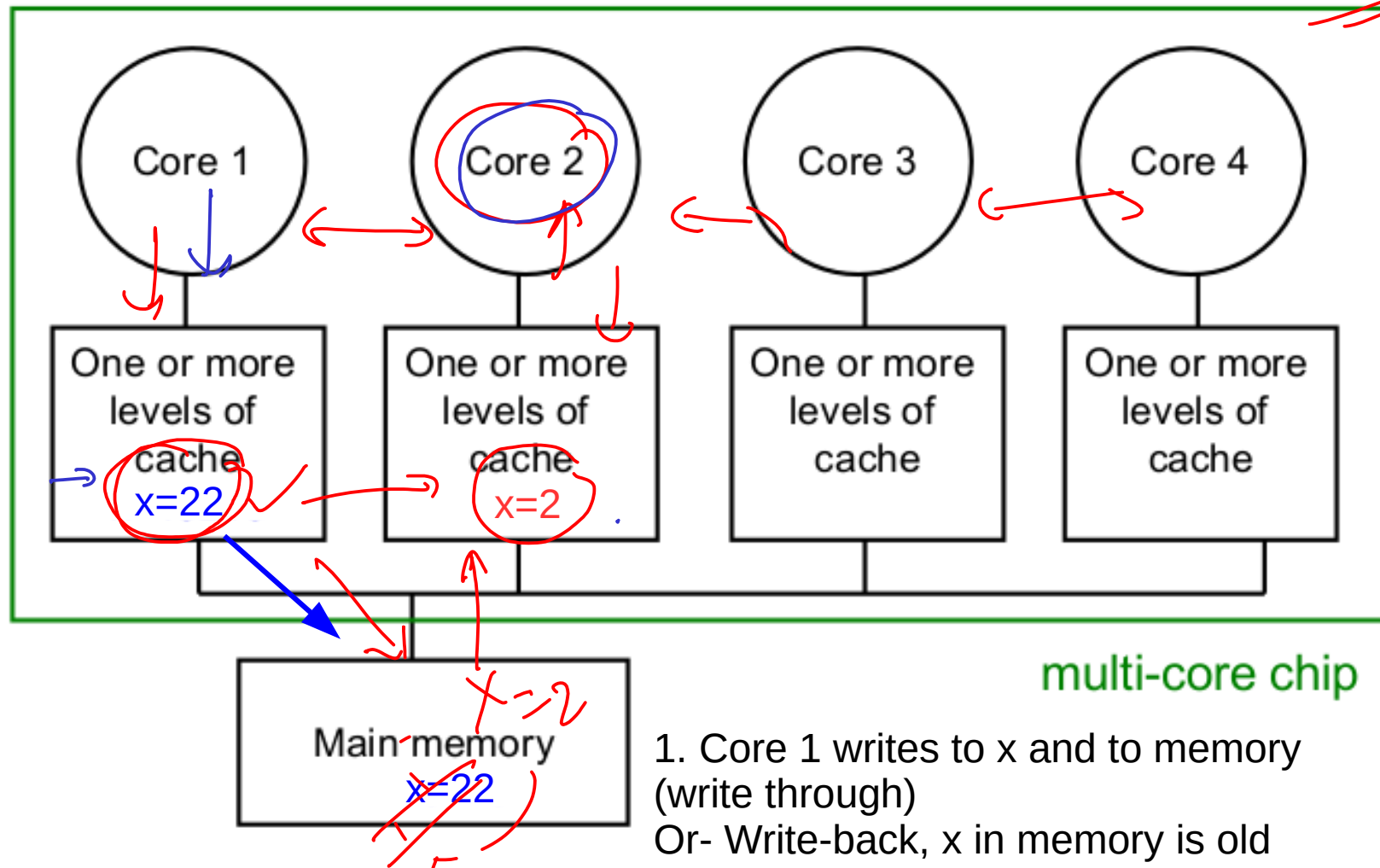
One or more levels of cache

Main memory   x=2
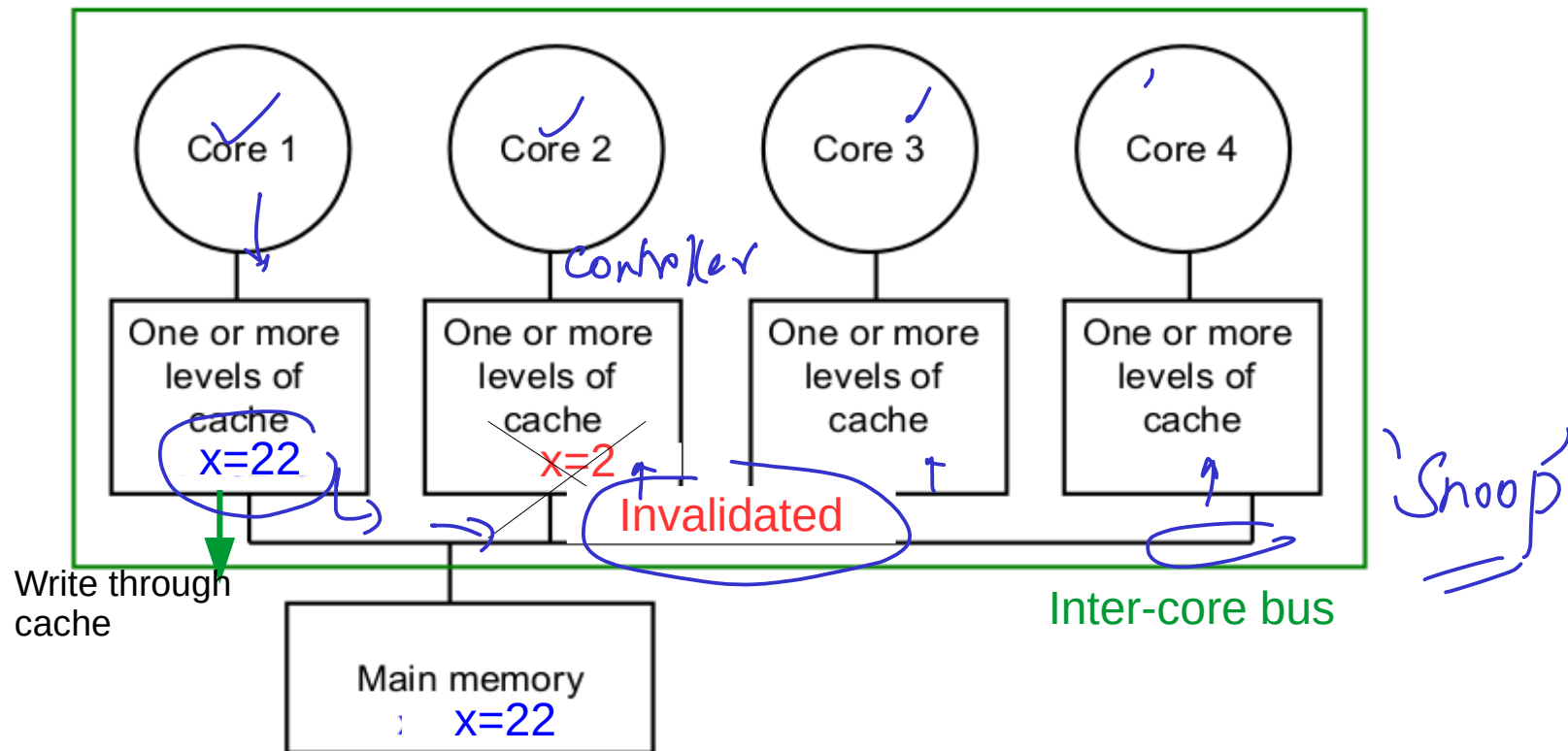
13

# Cache coherence
## Cache writes

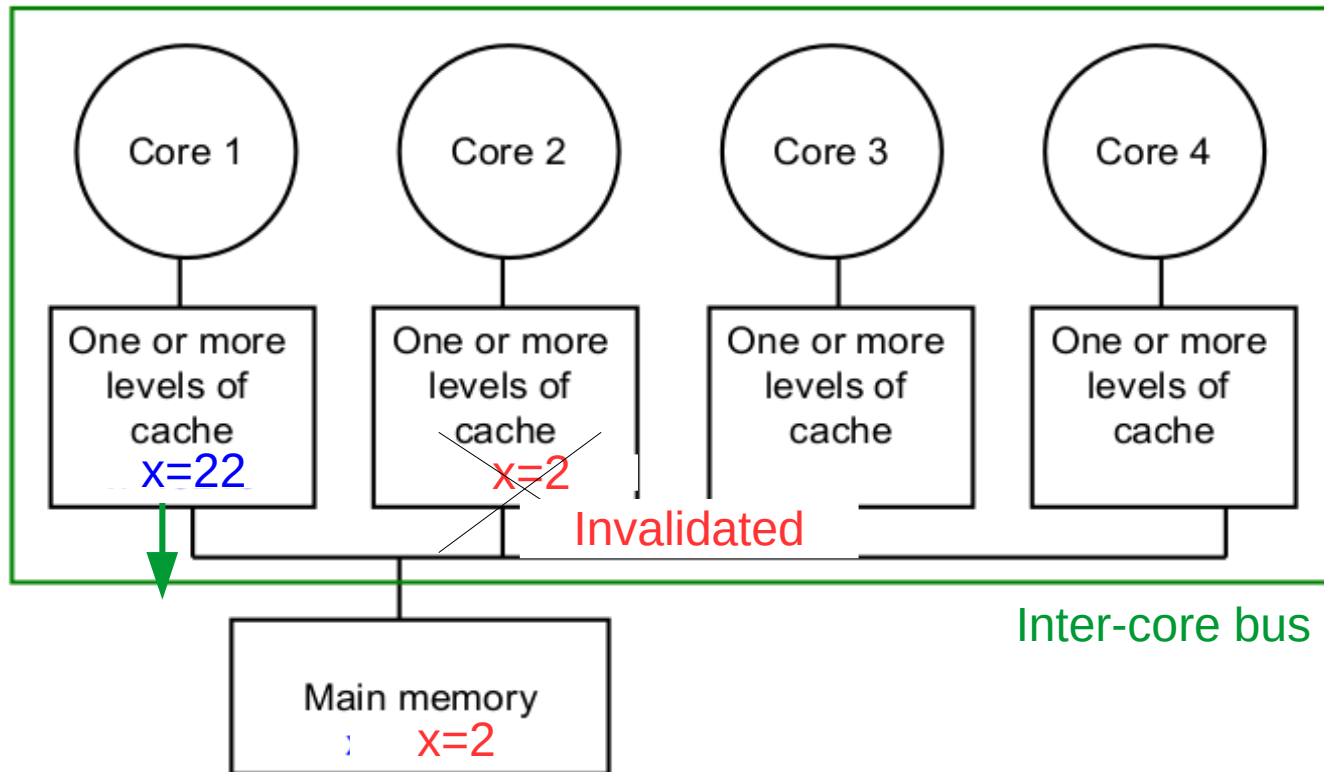Coherence: Processors read the latest value



1. Core 1 writes to x and to memory (write through)
Or- Write-back, x in memory is old
2. Core 2 reads x --> Wrong value x=2

# Option 1: Invalidation protocol
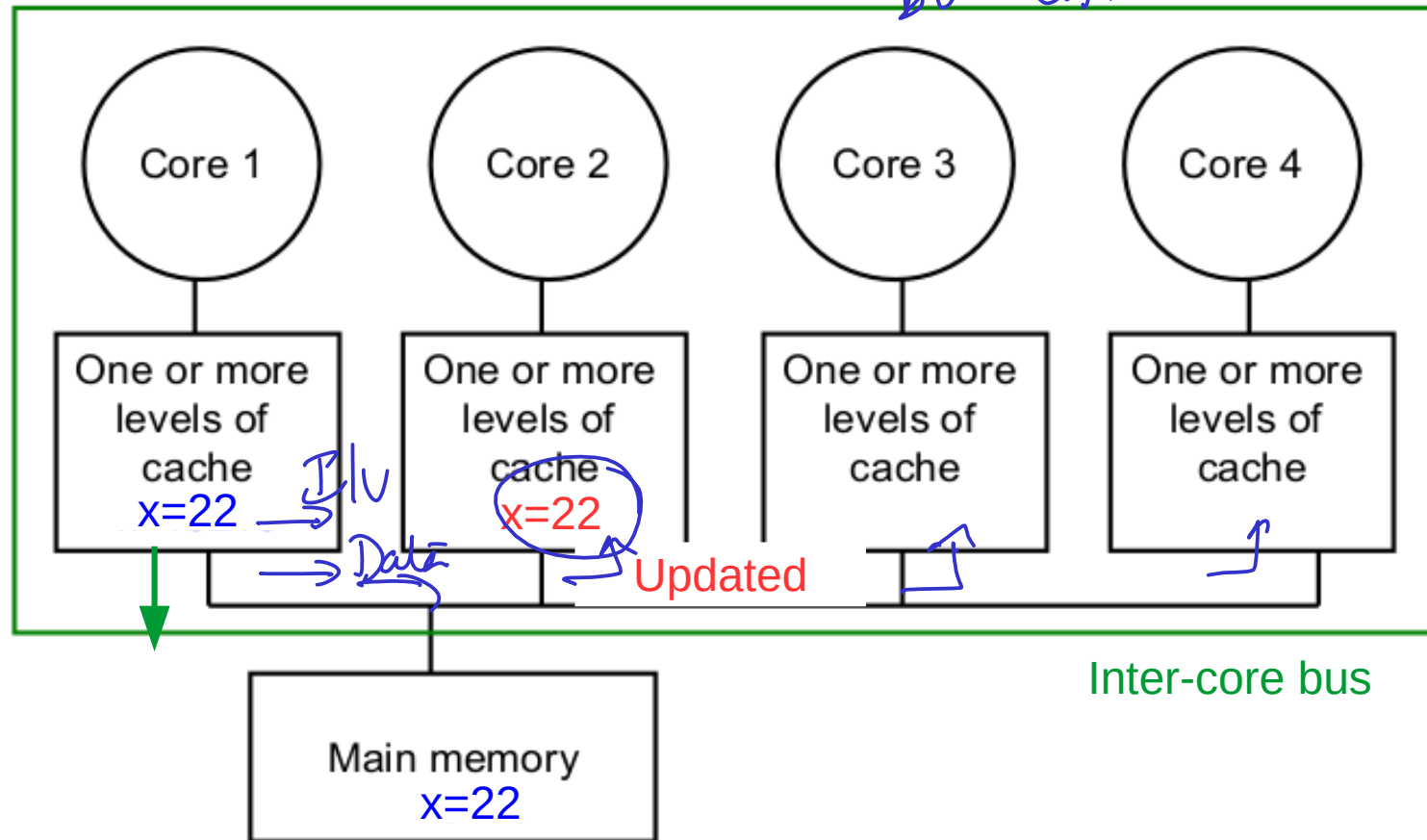


Write through cache

Inter-core bus

1. Core 1 gains control of the bus
2. Sends invalidation request on the bus --> copies of this block in other caches are invalidated (this should be done before writing. Else other cores might write to x)
3. All cores continuously "snoop" (monitor) the bus. Each core's cache controller invalidates its local copy
4. Core 1 has exclusive write permission to x and to memory (write through). Write x.
5. Core 2 reads x --> Cache miss--> loads new value from memory

15

# Write back: Invalidation protocol



1) Write back cache
2) Core 2- read request --> All other cores check if they have a local copy whose dirty bit=1
3) Read from memory is terminated.
4) Core 1 places the data on the bus, Core2 reads it.
5) Core 1 writes to memory when the line is evicted

# Option 2: Update protocol

*= Broadcast*



**Inter-core bus**

Core 1 · Core 2 · Core 3 · Core 4

One or more levels of cache — x=22 (Core 1)

One or more levels of cache — x=22 *Updated* (Core 2)

One or more levels of cache (Core 3)

One or more levels of cache (Core 4)
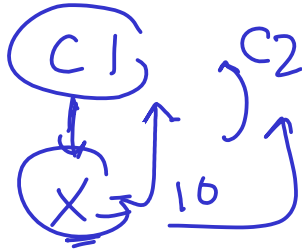
*Iu → Data*

Main memory
x=22

1. Core 1 writes to x and to memory (write through)
2. Core 1 sends the updated value on the bus (Broadcast)
3. All cores containing the data block update the value
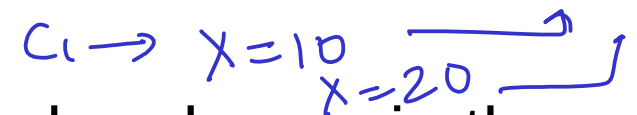4. Core 2 reads x -->  new value

17

# Invalidation vs Update

|  | Update | Invalidate |
|---|---|---|
| Writes from a core to the same address | Broadcast each write and update value in all cores<br><br>More bus traffic. | Invalidate first time, remaining writes do not send invalidation on the bus |
| Write to different words in the same block | Updates each word | First write invalidates all words in that block |
| C1 writes--> followed by C2 reads | Updates the value --> This works better | Invalidate causes a miss in C2 everytime. Not a good protocol here |

# Cache coherence Definition

- Core C1: Write to a location X, followed by a read from location X, returns the most recently written value (just like a uni-processor)

  - If there is no other write to X by another processor occurring between the write and the read

- C1 writes to X, C2 reads from X after "sufficient" time, C2 reads the most recent value from X: Coherent view of memory

- Writes to the same location are serialized and seen in the same order

  – C1 writes to Y and then C2 writes to Y, all cores must see the writes in the same order

# Cache coherence protocols

*Key idea: Track the state of any shared data block*

- ## Directory based protocol
  - Information of shared data is maintained in a directory. Every cache can track the sharing status.

- ## Snooping based protocol
  - Every cache maintains the status of each data block
  - Data is accessed through a shared bus between processors
  - Cache controllers monitor (snoop) the caches and the bus for the presence of the shared data

# Cache coherence protocols
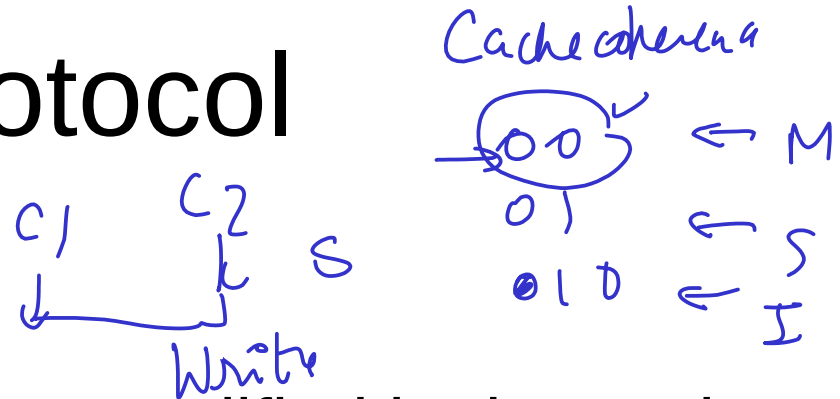# Invalidation based protocol
# MSI protocol

*Key idea: Track the state of any shared data block*

- Acknowledgements:
  - CA, Quantitative approach, Henessey and Patterson
  - Multicore Memory Caching Issues - Cache Coherency: PRACE Summer School 21-23 June 2012
  - Biswa Panda, IITK, Cache Coherence Protocols, CS422-Spring 2018
  - MSI Coherence - Georgia Tech - HPCA: Part 5

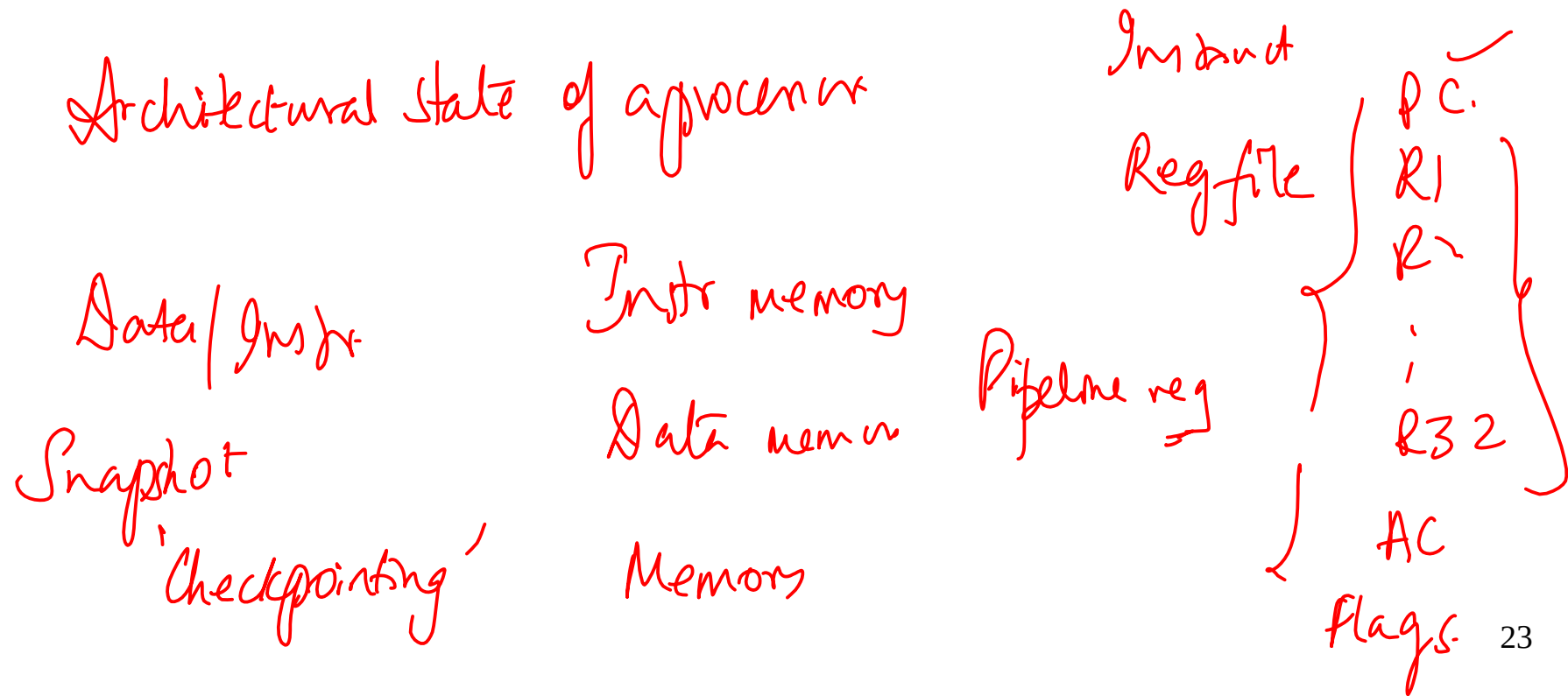# MSI protocol

3 states for a cache block:

- Modified :  The block has been modified in the cache. (Usually for writes). The data in the cache is then inconsistent with memory. Similar to Valid=1, dirty bit=1. Need to writeback before evicting

- Shared: This block is in one/more cache.  Memory-copy is up-to-date (Valid=1, dirty bit=0). The cache can evict the data without writing it to memory

- Invalid: This block is either not present in cache or has been invalidated by a bus request. Valid=0

# MSI protocol

Processor requests

- PrRd: Processor request to read a cache block.
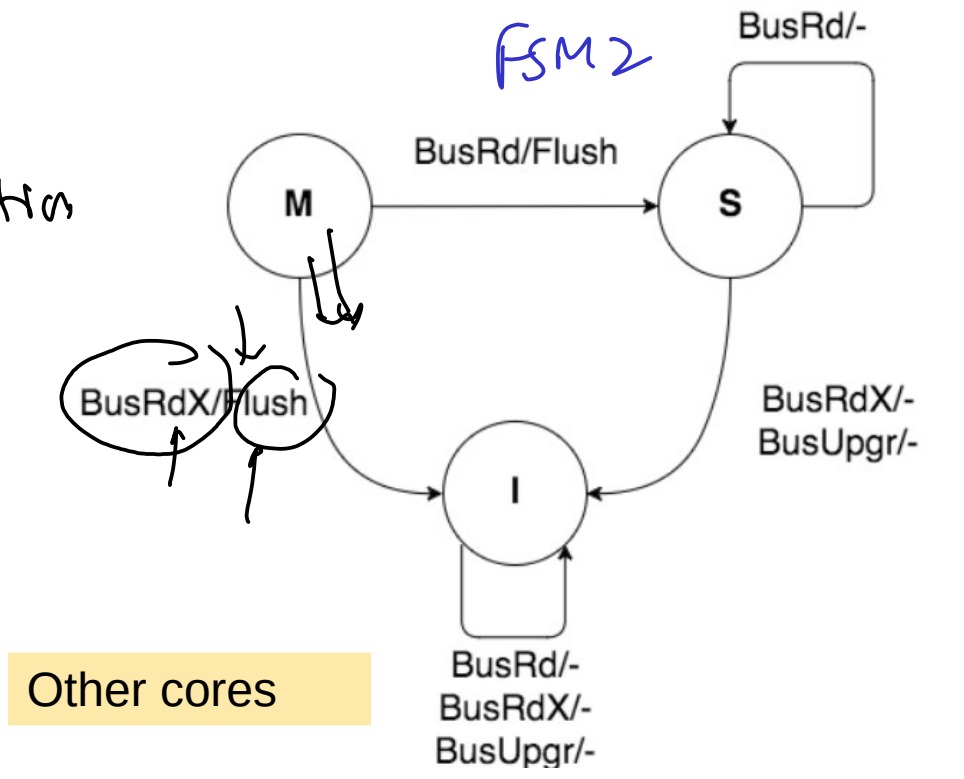- PrWr: Processor request to write a cache block.

Architectural State of a processor

Instruct

Reg file { P.C.
R1
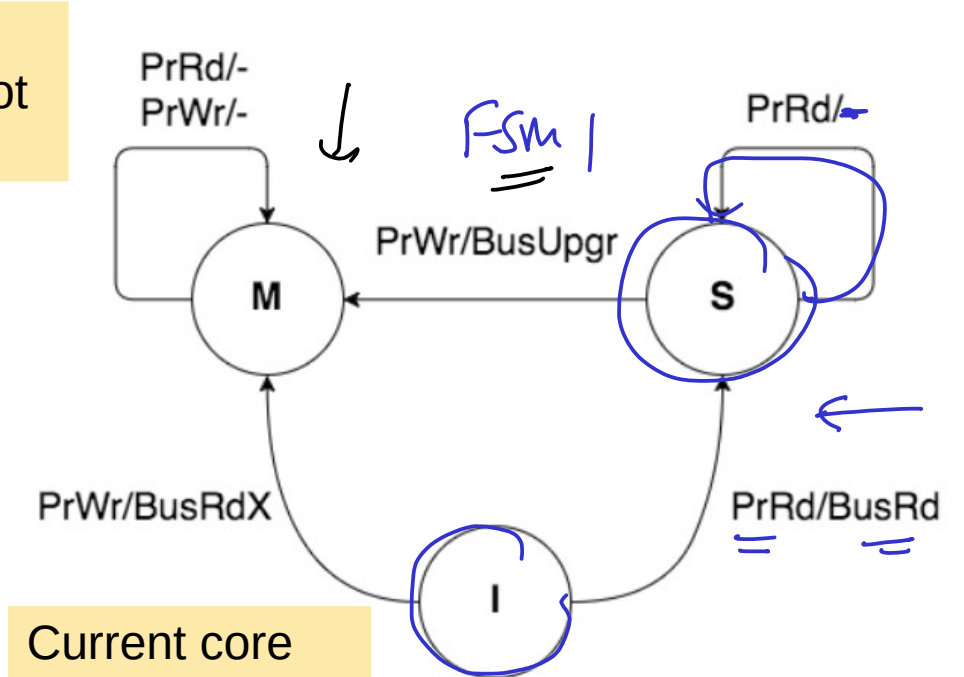R2
:
R32
AC
flags

Data/Instr.

Snapshot

'Checkpointing'

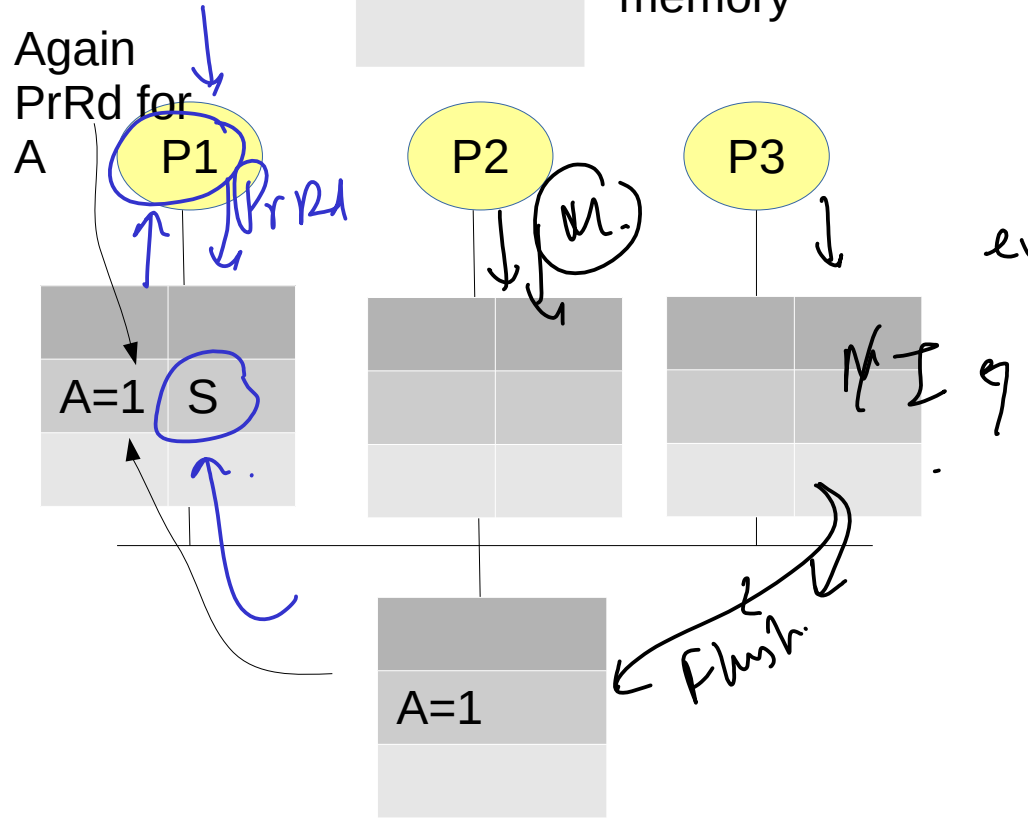Instr memory
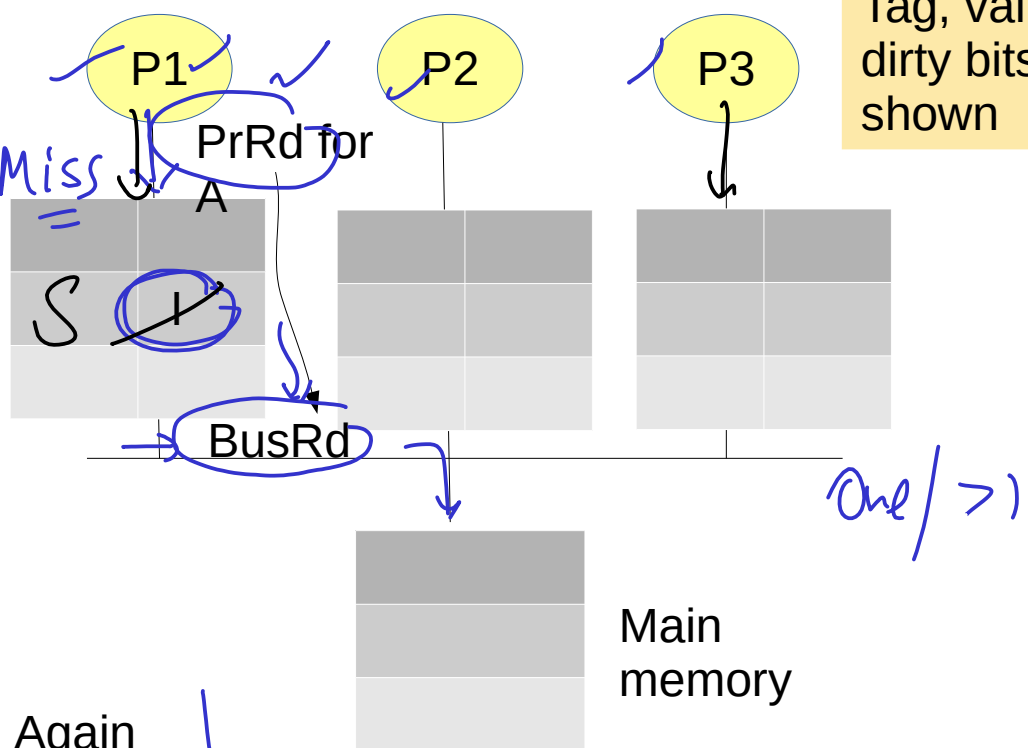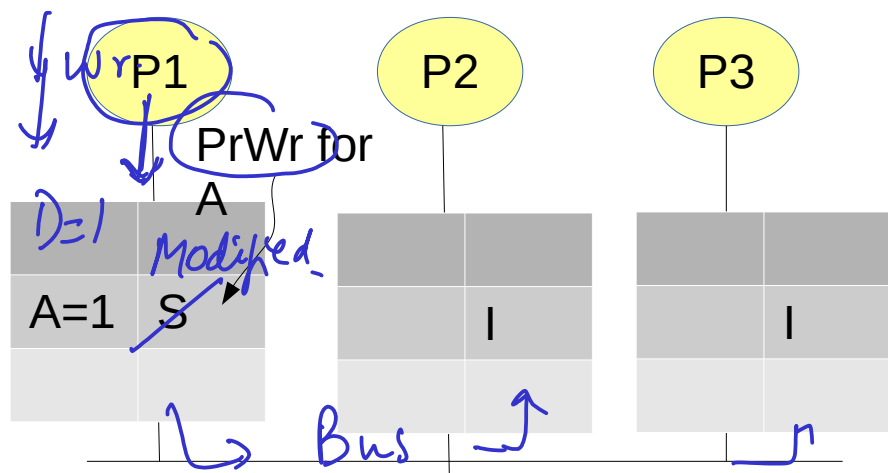
Data memory

Memory

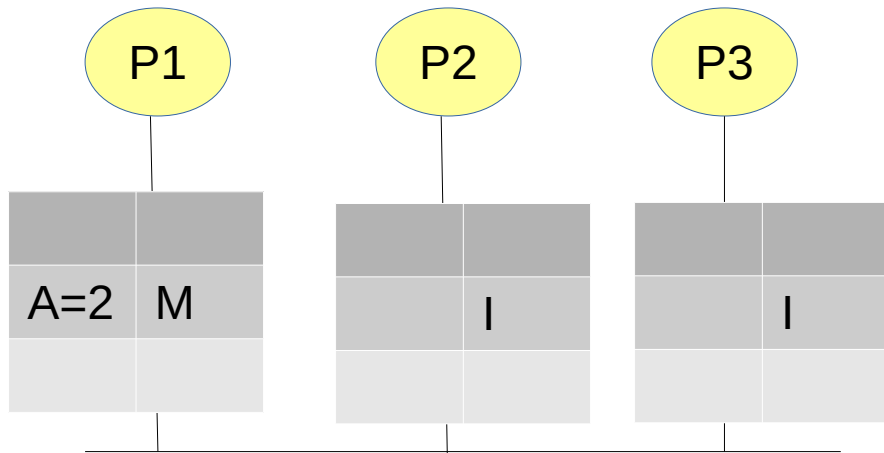Pipeline reg

# MSI protocol

Cache requests

- BusRd: When a read miss occurs in a processor's cache, it sends a BusRd request on the bus and expects data from memory

- BusRdX: When a write miss occurs in a processor's cache, it sends a BusRdX request on the bus which returns the cache block and gives exclusive write access to this cache. Needs to invalidate the block in the caches of other processors.

- BusUpgr: When there's a write hit in a processor's cache, it sends a BusUpgr request on the bus to invalidate the block in the caches of other processors. <u>BusUpgr and BusRdX can be combined</u>

- Flush/Writeback: Request indicates that a whole cache block is being written back to the memory

P1  P2  P3

Tag, valid, dirty bits not shown

Miss

PrRd for A

BusRd

One/>1

Main memory

Again PrRd for A

P1  P2  P3

PrRd

eviction

M I q

A=1  S

A=1

Flush

**FSM 1**

PrRd/- PrWr/-    PrRd/-

**M**    PrWr/BusUpgr    **S**

PrWr/BusRdX    PrRd/BusRd

**I**

Current core

**FSM 2**

BusRd/-

**M**    BusRd/Flush    **S**

BusRdX/Flush

BusRdX/- BusUpgr/-

**I**

BusRd/- BusRdX/- BusUpgr/-

Other cores
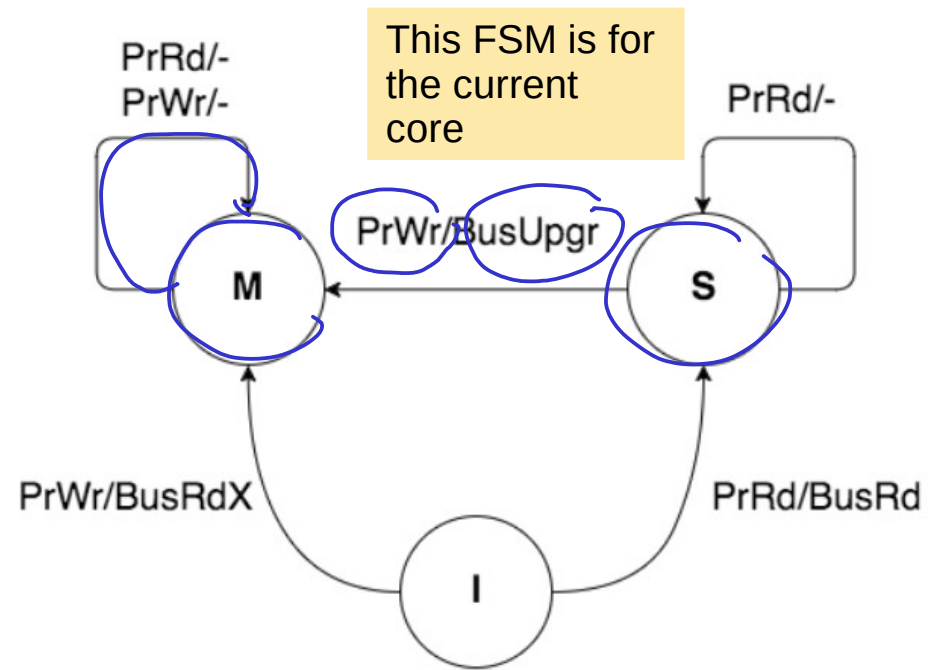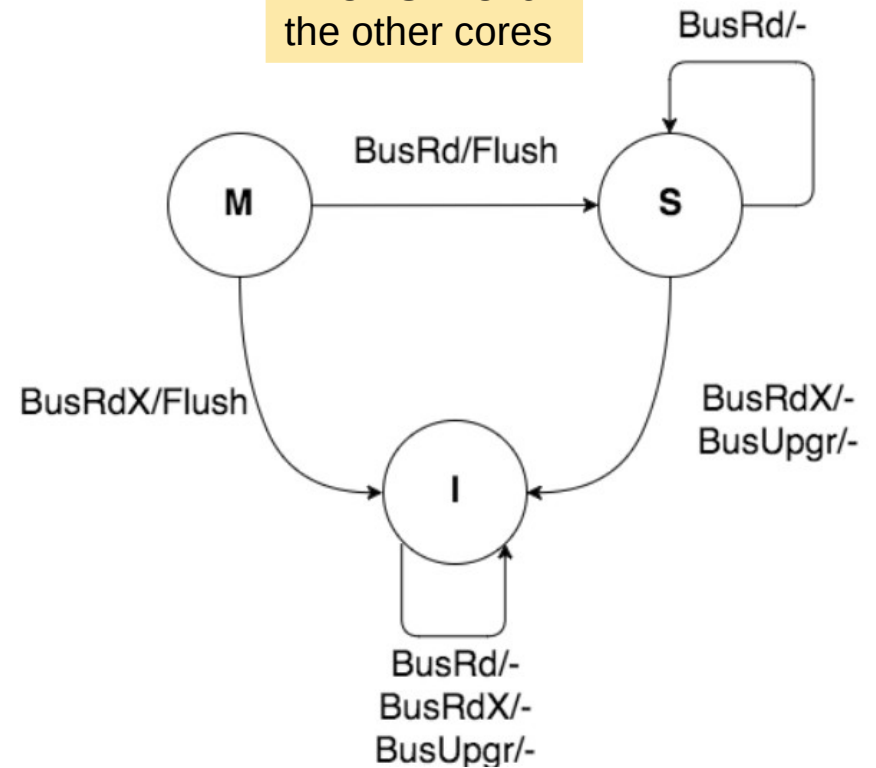
Wiki

This FSM is for the current core

This FSM is for the other cores

Write back cache assumed Cache hit

M means V=1, D=1

Wiki

P1  P2  P3

PrRd for A

A=2  M

M-D=1

S  D=0.

BusRd

Rd for the same cache line

I

A=1

S means V=1, D=0

P1  P2  P3

PrRd for A

A=2  S    A=2  S

Flush

**Writeback while changing from M state to S state (sinice S means dirty=0)**

A=2

Cache controller checks other caches: no direct memory read

→ FSM 1

PrRd/-
PrWr/-

PrRd/-

PrWr/BusUpgr

**M**          **S**

PrWr/BusRdX          PrRd/BusRd

**I**

2nd FSM

P1 sees BusRd on bus

BusRd/-

BusRd/Flush

**M**          **S**

BusRdX/Flush

BusRdX/-
BusUpgr/-

**I**

BusRd/-
BusRdX/-
BusUpgr/-

**Cache to cache transfer:**

Option 1:
- Ask Core 2 to abort/retry.
- Core 1 writes back to memory and then Core 2 reads from memory
- Large memory latency

Option 2: Core 1 intervention
- Core 1 informs memory that Core1 will respond to read request
- Core 1 will send the data to Core 2
- Since both cores are in Shared state (dirty=0), memory needs to pick this data and get updated

28

P1    P2    P3

PrWr for A

A=2  S    A=2  S

BusUpgr – exclusive access to cache block

A=2

D = 1
V = 1
Only core

PrRd/-
PrWr/-

M

PrWr/BusUpgr

PrRd/-

S

PrWr/BusRdX

PrRd/BusRd

I

P1    P2 Rd/wr    P3

PrRd for A

A=2  I    A=3  M

A=2

M means V=1, D=1

BusRd/-

M    BusRd/Flush    S

BusRdX/Flush

BusRdX/-
BusUpgr/-

I

BusRd/-
BusRdX/-
BusUpgr/-

Wiki

# P1 | P2 | P3

**Wr!**

PrRd for A

| A=2 | I |
|-----|---|

| A=3 | M |
|-----|---|

BusRd**X**

| | |
|---|---|
| A=2 | |

**Arch State of processor**

---

**P1**

→ PrRd/-
PrWr/-

**Exceptions**

PrRd/-

PrWr/BusUpgr

**M** ← **S**

PrWr/BusRdX

PrRd/BusRd

**I**

---

# P1 | P2 | P3

PrRd for A

| A=3 | S |
|-----|---|

| A=3 | S |
|-----|---|

| | |
|---|---|
| A=3 | |

P2 sees BusRd on bus

**P2**

**P3**

Flush --> Writeback

---

BusRd/-

BusRd/Flush

**M** → **S**

BusRdX/Flush

BusRdX/-
BusUpgr/-

**I**

BusRd/-
BusRdX/-
BusUpgr/-

Wiki

Alternately if P1 wanted to write

MSI

PrWr for A

A=2   I

BusRdX

A=2

P1   P2   P3

A=8   M

A=3   I

BusRdX

A=3

Flush--> since we invalidated the block Writeback

P2 sees BusRdX on bus

PrRd/-
PrWr/-

PrWr/BusUpgr

PrRd/-

M        S

PrWr/BusRdX        PrRd/BusRd

I

BusRd/-

M        BusRd/Flush        S

BusRdX/Flush        BusRdX/-
                   BusUpgr/-

I

BusRd/-
BusRdX/-
BusUpgr/-

Wiki

Alternately if P3 wanted to write

P1
P2
P3

PrWr for A

A=2  I

A=3  M

I  → M.

D=1
V←1

BusRdX

A=2

PrRd/-
PrWr/-

PrRd/-

PrWr/BusUpgr

M          S

PrWr/BusRdX          PrRd/BusRd

I

P2 sees BusRdX on bus

P1
P2
P3

Write

A=2  I

A=3  I

A=10  M    D=1

BusRdX

A=3

Flush--> since we invalidated the block Writeback

BusRd/-

BusRd/Flush

M          S

BusRdX/Flush

BusRdX/-
BusUpgr/-

I

BusRd/-
BusRdX/-
BusUpgr/-

# Eg 2

- <u>R1</u>, R2, W3, R2, W1, W2, R3, R2

David Henty, EPCC, Edinburgh

# Run through this sequence

- R1, R2, W3, R2, W1, W2, R3, R2

David Henty, EPCC, Edinburgh

# Run through this sequence

- R1, R2, <u>W3</u>, R2, W1, W2, R3, R2

David Henty, EPCC, Edinburgh

# Run through this sequence

- R1, R2, W3, <u>R2</u>, W1, W2, R3, R2

David Henty, EPCC, Edinburgh

# Run through this sequence

- R1, R2, W3, R2, W1, W2, R3, R2

David Henty, EPCC, Edinburgh

# Run through this sequence

- R1, R2, W3, R2, W1, W2, R3, R2

David Henty, EPCC, Edinburgh

# Run through this sequence

- R1, R2, W3, R2, W1, W2, <u>R3</u>, R2

David Henty, EPCC, Edinburgh

# Run through this sequence

- R1, R2, W3, R2, W1, W2, R3, <u>R2</u>

# Summary: Invalidation based protocol

- Ensure that a processor has exclusive access to a data item before it writes that item.

- It invalidates other copies on a write.

- Exclusive access ensures that no other readable or writable copies of an item exist when the write occurs: All other cached copies of the item are invalidated.

- To perform an invalidate, the processor acquires bus access and broadcasts the address to be invalidated on the bus.

- All processors continuously snoop on the bus, watching the addresses. The processors check whether the address on the bus is in their cache. If so, the corresponding data in the cache are invalidated

- Two processors trying to gain bus access: First one gains priority

# Challenges and solutions

- Beyond 8 cores snooping bandwidth at the caches is a bottleneck

- Every cache must examine every miss and request placed on the bus

- Directory based protocol is a solution:
  - Intel i7 in Xeon series places a directory with cache status at the last level cache
  - AMD Opteron: Mix of snoop and directory. Broadcasts to find shared copies, uses acknowledgement for ordering operations

- To increase the communication bandwidth between processors and memory, we can use multiple buses and interconnection networks, such as crossbars or small point-to-point networks.

# MSI disadvantage

- A processor reads a cache line and then writes to it, when no other core has that cache line

- I--> S --> M, Bus Rd --> BusUpgr

- BusUpgr- to tell all other processors that it is about to write to the block. Not useful if no other processors have that block

- To remove the need to issue this BusUpgr/ bus transactions, MESI adds E.

# MESI protocol

- Modified: Valid=1,Dirty=1 (Writes). Only copy, only updated copy.

- Exclusive: Valid=1, Dirty=0, Only copy, clean copy, consistent with memory (Reads)

- Shared: Valid=1, Dirty=0, This block is in one more cache. The cache can evict the data without writing it to memory

- Invalid: Valid=0, This block is either not present in cache or has been invalidated by a bus request.

P1  P2  P3

PrRd for A

E / I

CC    BusRd    CC    CC

$I \to S \to M$

Bus Upgrd

P1  P2  P3

A=1  E

A=1

$I \to S \to M$

PrRd/-
PrWr/-

PrRd/-

WB = -
D=1

PrWr/-

M    E

1st FSM
P1

PrWr/BusRdX

PrRd/BusRd(!C)

PrWr/BusUpgr

S    I

PrRd/BusRd(C)

addr 20    I S S

PrRd/-

2nd

M    BusRdX/Flush    E

I M I I
I - I
M

R2    BusRd/Flush    BusRd/FlushOpt    BusRdX/FlushOpt

R3

S    I

BusRdX/FlushOpt
BusUpgr/-

BusRd/-
BusRdX/-
BusUpgr/-

BusRd/FlushOpt

$I \to S \to M$

Wiki

P1    P2    P3

PrWr for A

A=1   E   M
D=1

No bus transaction during E --> M

A=1

P1    P2    P3

A=2   M

A=1

J → E → M

PrRd/-
PrWr/-

PrWr/-

M        E

PrWr/BusRdX

PrRd/BusRd(!C)

PrWr/BusUpgr

S        I

PrRd/BusRd(C)

PrRd/-

WB

I M I I

Main mem

BusRdX/Flush

M        E

BusRd/Flush    BusRd/FlushOpt    BusRdX/FlushOpt

S        I

BusRdX/FlushOpt
BusUpgr/-

BusRd/FlushOpt

BusRd/-
BusRdX/-
BusUpgr/-

47

Wiki

P1

P2

P3

PrRd for A

A=2  M

BusRd

A=1

SSI
$MI
→ IMI

PrRd/-
PrWr/-

PrRd/-

PrWr/-

M

E

PrWr/BusRdX

PrRd/BusRd(!C)

PrWr/BusUpgr

S

I

PrRd/BusRd(C)

PrRd/-

P1

P2

P3

PrRd for A

A=2  S

A=2  S

Writeback

A=2

Flush – no direct memory read

P1 sees BusRd on bus

BusRdX/Flush

M

E

BusRd/Flush

BusRd/FlushOpt

BusRdX/FlushOpt

S

I

BusRdX/FlushOpt
BusUpgr/-

BusRd/-
BusRdX/-
BusUpgr/-

BusRd/FlushOpt

48

Wiki

P1    P2    P3

PrWr for A

A=2   S

A=2   S

BusUpgr

A=1

PrRd/-
PrWr/-

PrWr/-

PrRd/-

M    E

PrWr/BusRdX

PrRd/BusRd(!C)

PrWr/BusUpgr

S    I

PrRd/BusRd(C)

PrRd/-

P1    P2    P3

PrRd for A

A=2   I

A=3   M

A=2

BusRdX/Flush

M    E

BusRd/Flush    BusRd/FlushOpt    BusRdX/FlushOpt

S    I

BusRdX/FlushOpt
BusUpgr/-

BusRd/-
BusRdX/-
BusUpgr/-

BusRd/FlushOpt

What happens if P1 reads now?

49

Wiki

PrRd
for A

P1  P2  P3

A=2  I    A=3  M

BusRd

A=2

P1  P2  P3

PrRd
for A

A=3  S    A=3  S

Flush +
Writeback

A=3

P2 sees
BusRd on
bus

PrRd/-
PrWr/-

PrWr/-

M    E

PrWr/BusRdX

PrRd/BusRd(!C)

PrWr/BusUpgr

S    I

PrRd/BusRd(C)

PrRd/-

BusRdX/Flush

M    E

BusRd/Flush    BusRd/FlushOpt    BusRdX/FlushOpt

S    I

BusRdX/FlushOpt
BusUpgr/-

BusRd/-
BusRdX/-
BusUpgr/-

BusRd/FlushOpt

50

Wiki

# MSI vs MESI

- MSI: A processor reads a cache line and then writes to it --> bus transactions?
  - I--> S --> M, Bus Rd --> BusUpgr
  - BusUpgr- to tell all other processors that it is about to write to X. Useless as no other processors have that line
  - To remove the need to issue this BusRdX, MESI adds E.
- MESI: If a cache line is in E state, then it is the only cache to contain that line. If the processor wants to write to that line, no bus operations need to occur.
  - Saves lot of bus transactions if a single processor is doing serial read after writes

# Valid/invalid transitions

✓ addr 20 ⊂

- M -- > E    X

- E --> M    Rd/lw ⟶ Sw/wr ✓

- E-->I    ✓

- S-->E ⟹ X

- I-->S ✓ Read ⊤ M/S

- E-->E ⟹ Rd-Rd

I S S   addr 20

# Valid/invalid transitions

- M -- > I--> E
- E --> M: On a local write
- E-->I: If another core writes to same line
- S-->I--> E
- I-->S: Read, while another core has the same block
- E-->E: Local reads

# ?

- Assume processor P1 and P2 trying to access cache line X. Say, P1 has X in E state, P2 has it in M state. If P2 requests for a cache write and it is a hit. What will be the new states of the 2 cache lines in a MESI protocol?

- Assume processor P1 and P2 trying to access cache line X. Say, P1 has X in M state, P2 has it in I state. If P2 requests for a cache write. What will be the new states of the 2 cache lines in a MESI protocol?

# Problem with MSI or MESI

- Core 1 in M state, Core 2 is in I state
  - Core 2 Read
  - Core 1 does **Memory write**
  - Core 1--> S, Core 2--> S
  - **Memory writes on cache-cache transfers**
- Core 1 and Core 2 in S state. If Core 3 wants to read, who provides the data?
  - If one of the cores have to respond, which one? And how do we decide?
  - Memory responds --> **Memory read**
- **Memory reads/writes are expensive**
- Solution: A single "owner" of the data to respond to cache reads (avoid memory reads) and writes to memory later

*[handwritten annotations:]*

Read
M    I
S    S
Cache-Cache transfer

Read
S  S  P

# MOESI protocol

- Owned: Combines M and S states. Valid=1, Dirty=1. Not the only copy. Shared read access. Owner responds for reads and also writes back to memory

- Modified: Valid=1, Dirty=1. Only copy, only updated copy. (Writes)

- Shared: Valid=1, Dirty=0, This block is in <span style="color:red">more than one cache</span>. Clean copy. Shared read access

- Exclusive: Valid=1, Dirty=0, Only copy, clean copy (Reads)

- Invalid: Valid=0, This block is either not present in cache or has been invalidated by a bus request.

# MOESI

P1   P2   P3

PrRd for A

A=2 M

I

BusRd

No WB

A=1

MOESS

O behaves like S state, except that it is responsible for providing data to Core 3.

Writes back to memory when it gets replaced or gets invalidated.

P1   P2   P3

PrRd for A

A=2 O

A=2 S

No memory write

A=1

P1   P2   P3

PrRd for A

A=2 O

A=2 S

A=2 S

P1 provides data to P3

A=1

# MOESI

O behaves like S state, except that it is responsible for providing data to Core 3.

Writes back to memory when it gets replaced or gets invalidated.

P1  P2  P3

A=2  O    A=2  S    A=2  S

Wr

??

A=1

P1  P2  P3

PrRd for A

A=1

59

# MOESI

E I I
@ S I

P1 reads or writes — M

P1 requests to write

P1 reads — E

Exclusive read miss

Write miss

I I I
I M I

Other processors request

Other processors request to read

Other processors request to write

Other processors request to write

Other processor request to write

P1 requests to write

P1 requests to write

Other processors request to read

Read by any processor

S

Shared read miss

Other processors request to write

I

Other processors request to write

O

Read by any processor

S S I
O S I

M S I
J ≥ 1

M E S I
= Only one
≥ 1

M
↓
S
↓
I

M
↓
Owner
↳ 3rd cœr

M O E S I

M
↓
I O S S
I I M

E I I
S S
E I
S S I
S S S

- P1. M-->O when other processors read
- P1. O--> M when P1 writes
- P1. O --> I when other processor writes
- P1. O --> O, when other processors read

60

Researchgate

# Summary

| M | Exclusive R/W access | Dirty=1 | Respond with data and update memory |
|---|---|---|---|
| E | Exclusive R/W access | Dirty=0 | Reach E state upon a read from I |
| O | Shared access (only read) | Dirty=1 | Respond with data and update memory when replaced |
| S | Shared access (only read) | Dirty = 0 | Not responsible for providing data or updating memory |
| I | Invalid block | | |

|  | MSI | MESI | MOESI |
|---|---|---|---|
| RD Mem A | I --> | I --> | I --> |
| Wr Mem A |  |  |  |

Change in states and bus transactions??

Directory
Snoop
MOESI

|            | MSI                                  | MESI                                          | MOESI                                         |
|------------|--------------------------------------|-----------------------------------------------|-----------------------------------------------|
| RD Mem A   | I-->S                                | I-->E                                         | I-->E                                         |
| Wr Mem A   | S-->M (send invalidation on bus)     | E-->M (No need to send invalidation on bus)   | E--> M (No need to send invalidation on bus)  |

# Distributed Shared-Memory and Directory-Based Coherence

CA- QA - H&P Ch – 5

HPCA Part 5:
https://www.youtube.com/watch?v=xjRDejGF26M&list=PLAwxTw4SY aPkr-vo9gKBTid_BWpWEfuXe&index=110

"A Primer on Memory Consistency and Cache Coherence.pdf"
Synthesis lectures on CA

# Disadvantages of snoop based protocols

- Snooping protocol: Communication with all caches on every cache miss and on writes

- Huge bus bandwidth for synchronization > 8 cores

- Solution: Non-broadcast network.  Distributed memory, and eliminate the need for broadcast on every miss

# Single directory

- A directory keeps the state of every block that may be cached, which caches have copies of the block

- Directory is maintained in the last level shared cache. Eg., Intel i7 maintains this in the inclusive L3 and optionally does snoop when the bus is free (Combines directory + snoop)

- Cache controllers send all requests to the directory. For eg., if the data is present in Core2's cache, the data is forwarded to Core2 by the directory

# Distributed directory

- Directory is distributed such that, the coherence protocol knows where to find the directory information

- Each directory tracks the caches that share the memory addresses of the memory in that node.

# Same cache states

- Shared—One or more nodes have the block cached, and the value in memory is up to date (and in all the caches).

- Uncached—No node has a copy of the cache block.

- Modified—Exactly one node has a copy of the cache block, and it has written the block, so the memory copy is out of date. The processor is the owner of the block.

- State transitions are similar to snooping protocol

- Requests do not go on the bus. Goes to the directory

# State transitions

Core 1
Cache1

Core 2
Cache2

Directory entry for
Block A

Presence bit for each of the 8 cores
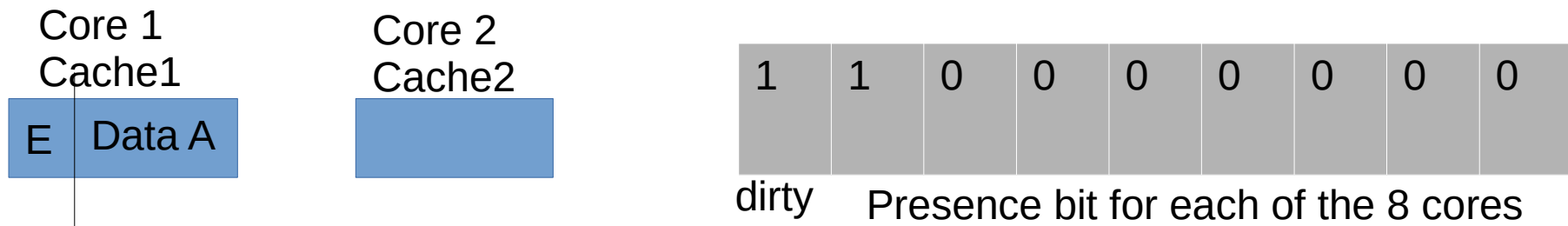
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|

dirty

Core 1: Read block A --> Read miss --> Read request goes to the directory slice which has the particular entry --> Indicates not present in any core

Data fetched from memory sent to Cache 1 --> Change Cache 1 from Invalid to Exclusive state (if doing MESI)--> Change directory state. (Dirty=1, since E can become M without broadcasting)

Advantage: Parallely Cache 2 could be interacting with another directory without bus bottleneck

Core 1
Cache1

| E | Data A |
|---|--------|

Core 2
Cache2

| 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|

dirty    Presence bit for each of the 8 cores

69

# State transitions

Core 1
Cache1

| E | Data A |

Core 2
Cache2

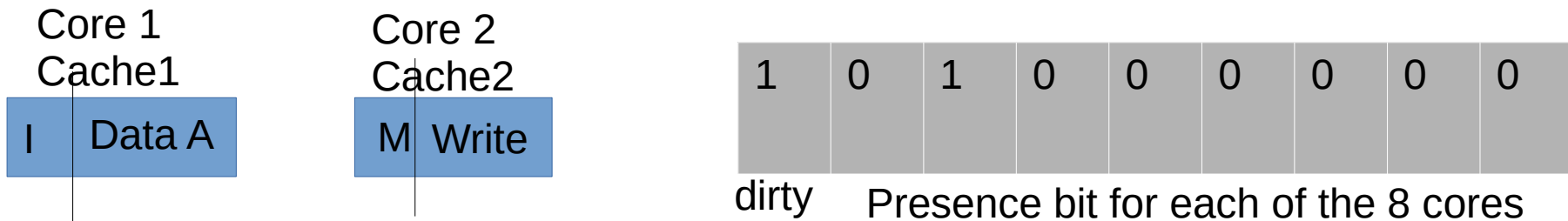| 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

dirty    Presence bit for each of the 8 cores

Core 2: Write to block A --> Sent to directory --> Block is in Cache 1. Directory forwards write request to Cache 1. Cache 1 **can invalidate and send acknowledgement back to the directory**

Core 1
Cache1

| I | Data A |

Core 2
Cache2

| M | Write |

| 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |

dirty    Presence bit for each of the 8 cores
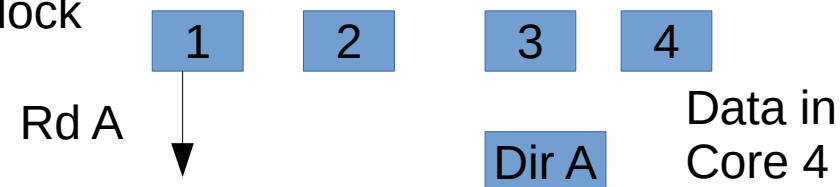
Local node is the node where a request originates
Home node: memory location/directory entry of address reside
[Physical address is distributed so the node that contains the memory and directory
for a given physical address is known]
A remote node is the node that has a copy of a cache block

| 1 | 2 | 3 | 4 |

1: Local, 3: Home, 4: Remote
P = requesting node, A = requested address

Rd A

Dir A

Data in
Core 4

| Message type | Source | Destination | Message contents | Function of this message |
|---|---|---|---|---|
| Read miss | Local cache | Home directory | P, A | Node P has a read miss at address A; request data and make P a read sharer. |
| Write miss | Local cache | Home directory | P, A | Node P has a write miss at address A; request data and make P the exclusive owner. |
| Invalidate | Local cache | Home directory | A | Request to send invalidates to all remote caches that are caching the block at address A. |
| Invalidate | Home directory | Remote cache | A | Invalidate a shared copy of data at address A. |
| Fetch | Home directory | Remote cache | A | Fetch the block at address A and send it to its home directory; change the state of A in the remote cache to shared. |
| Fetch/invalidate | Home directory | Remote cache | A | Fetch the block at address A and send it to its home directory; invalidate the block in the cache. |
| Data value reply | Home directory | Local cache | D | Return a data value from the home memory. |
| Data write-back | Remote cache | Home directory | A, D | Write-back a data value for address A. |

# A 4-core example with directory

- HPCA Part 5 Directory example:
  https://www.youtube.com/watch?v=lZZYILcQ68Y&list=PLAwxTw4SYaPkr-vo9gKBTid_BWpWEfuXe&index=113