# Multicores and Cache coherence

**Arkaprava Basu**

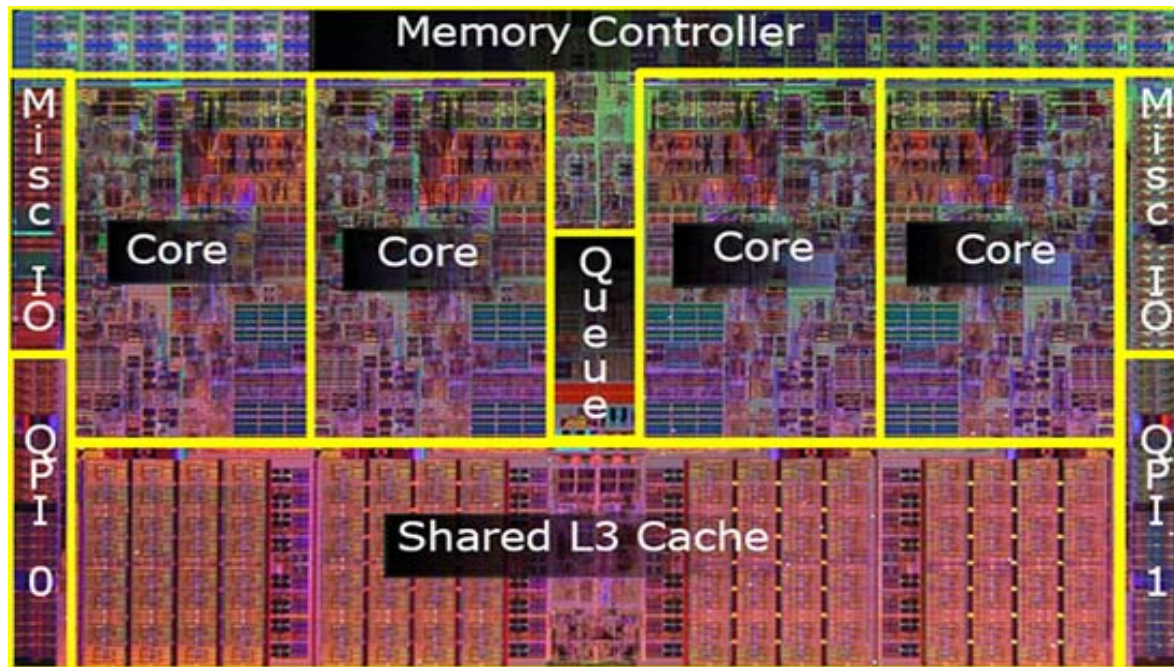**https://www.csa.iisc.ac.in/~arkapravab/**

# Acknowledgements

- Some of the slides in the deck were provided by Luis Ceze (Washington), Nima Horanmand (Stony Brook), Mark Hill, David Wood, Karu Sankaralingam (Wisconsin), Abhishek Bhattacharjee(Yale).
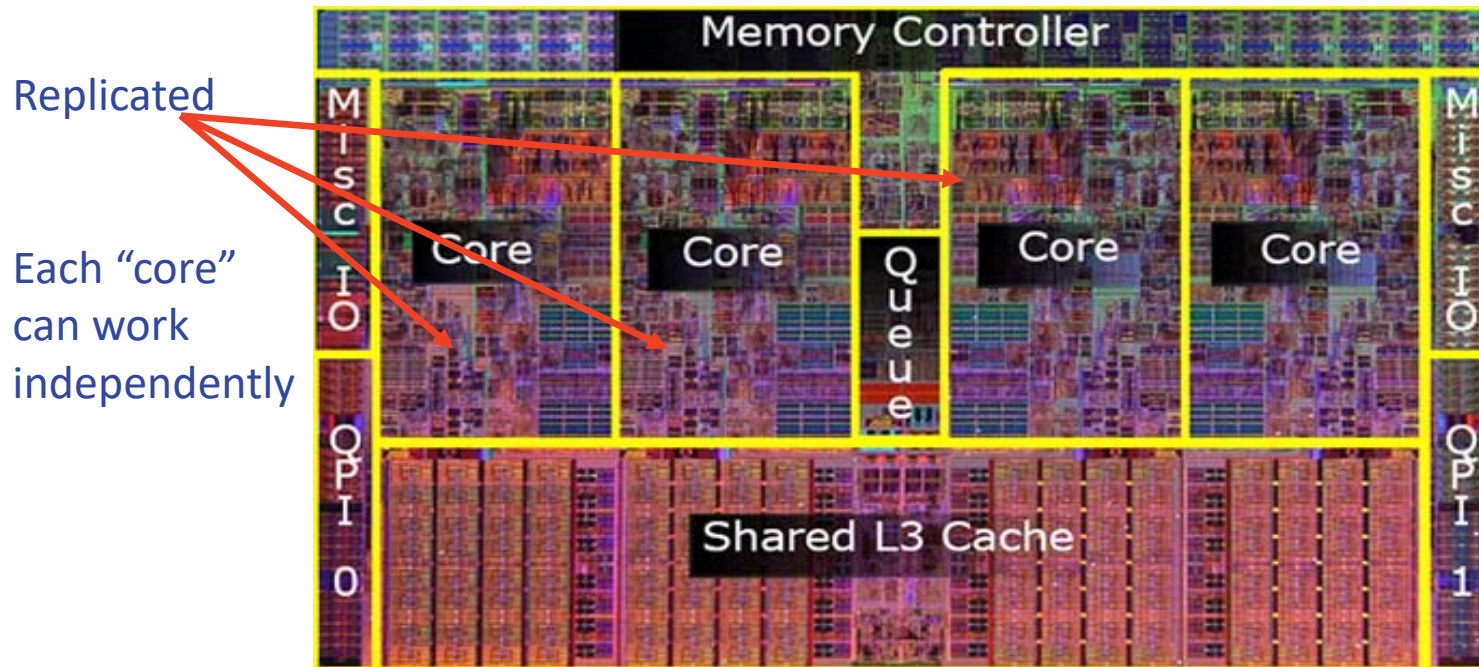
# What are Multicore chips?

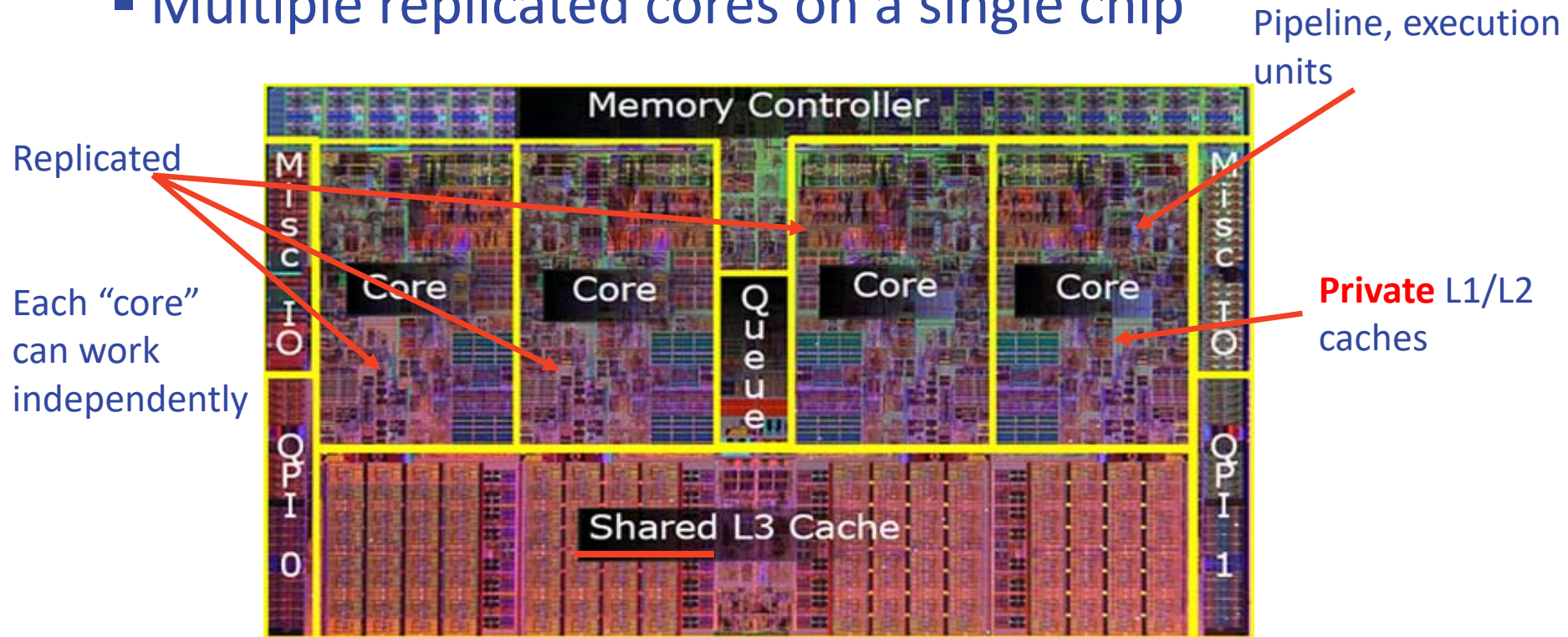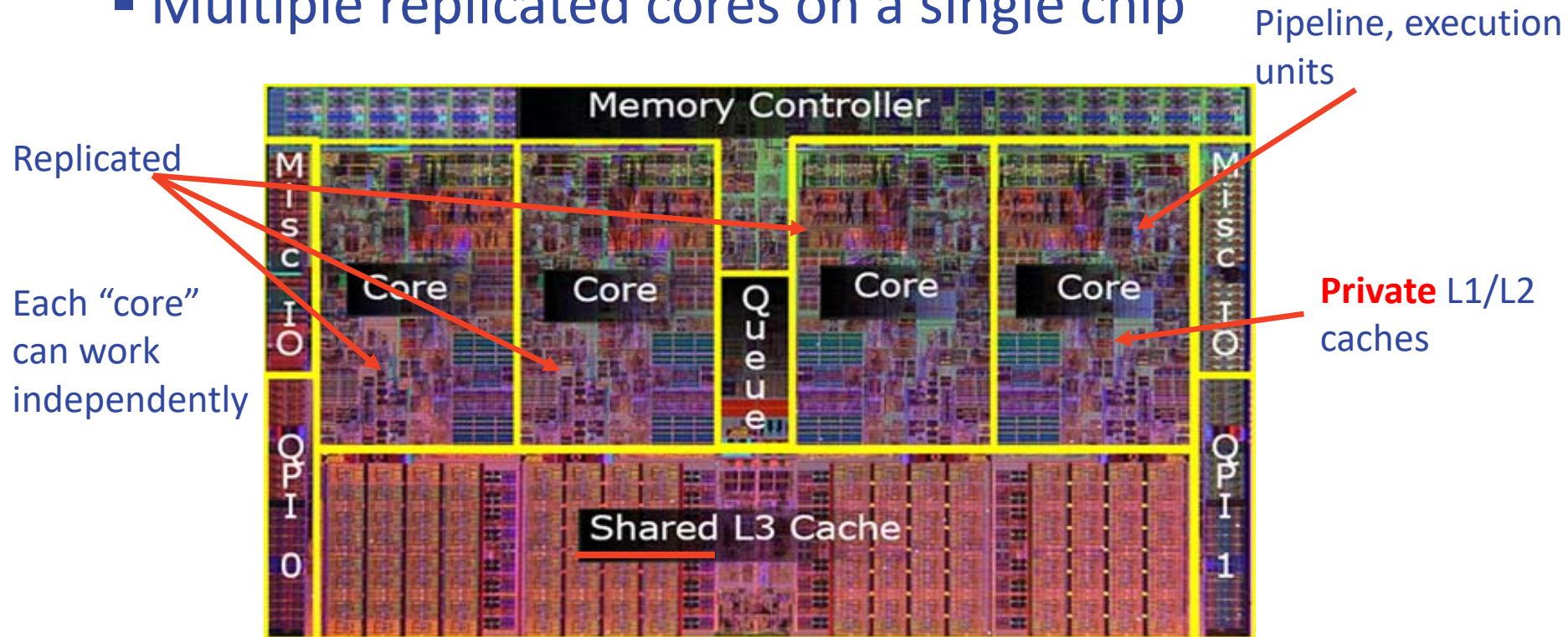- Multiple replicated cores on a single chip

# What are Multicore chips?

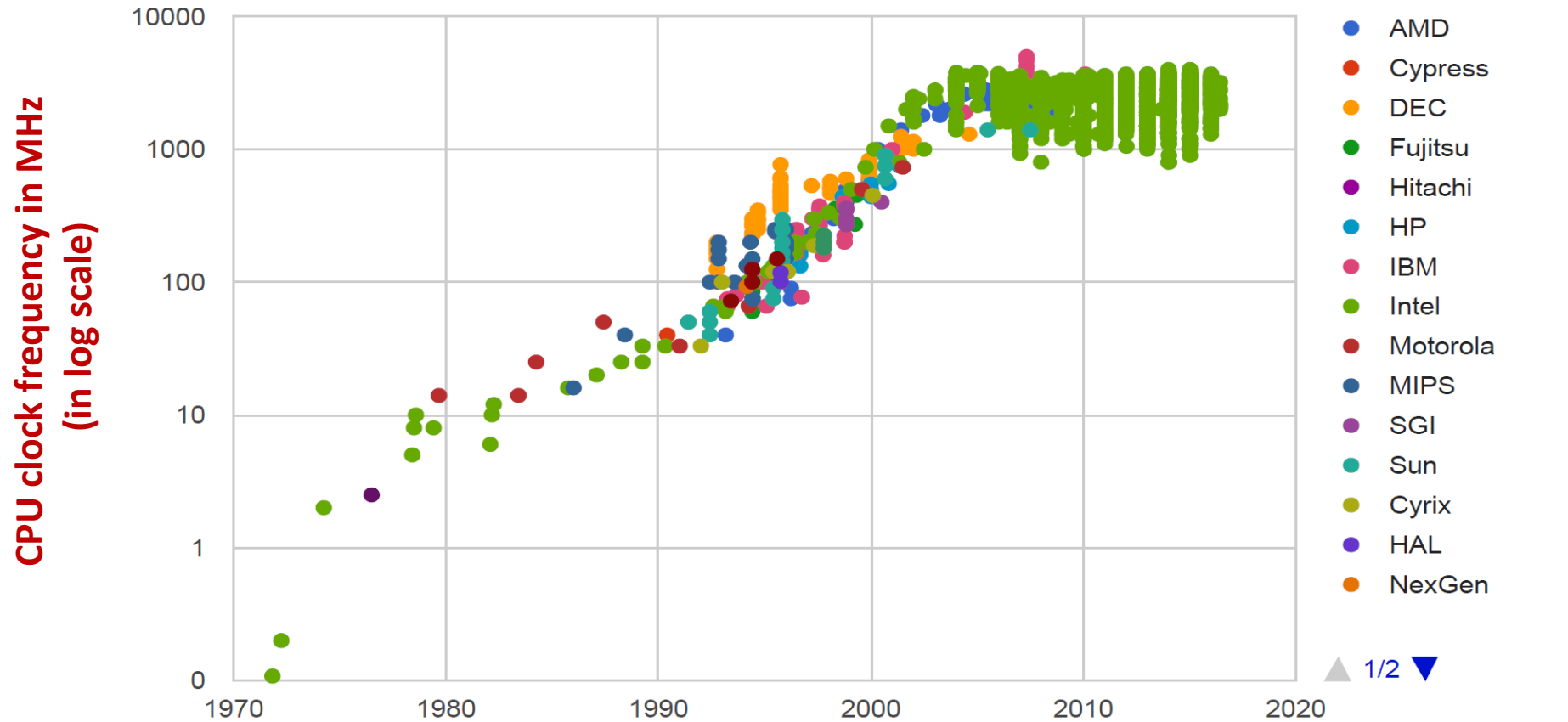- Multiple replicated cores on a single chip

Replicated

Each "core" can work independently

# What are Multicore chips?

- Multiple replicated cores on a single chip

Pipeline, execution units

Replicated

Each "core" can work independently

**Private** L1/L2 caches

# What are Multicore chips?

- Multiple replicated cores on a single chip

Pipeline, execution units

Replicated

Each "core" can work independently

**Private** L1/L2 caches

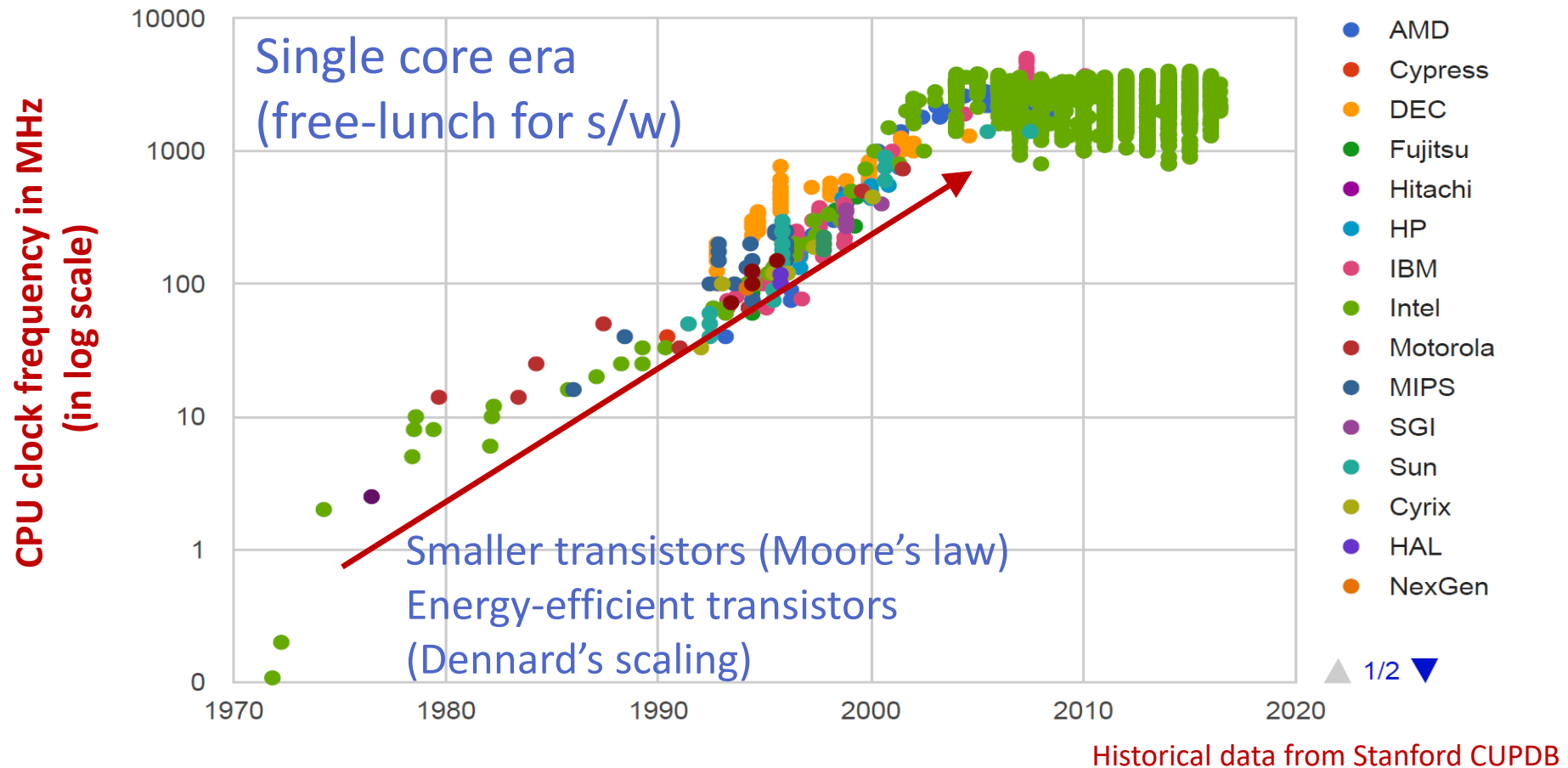Memory Controller | Misc IO | Core | Core | Queue | Core | Core | Misc IO | QPI 0 | QPI 1 | Shared L3 Cache

Examples:    Apple A14 (iPhone 12) 6 cores
Intel i7-10700T (desktop class) 8 cores
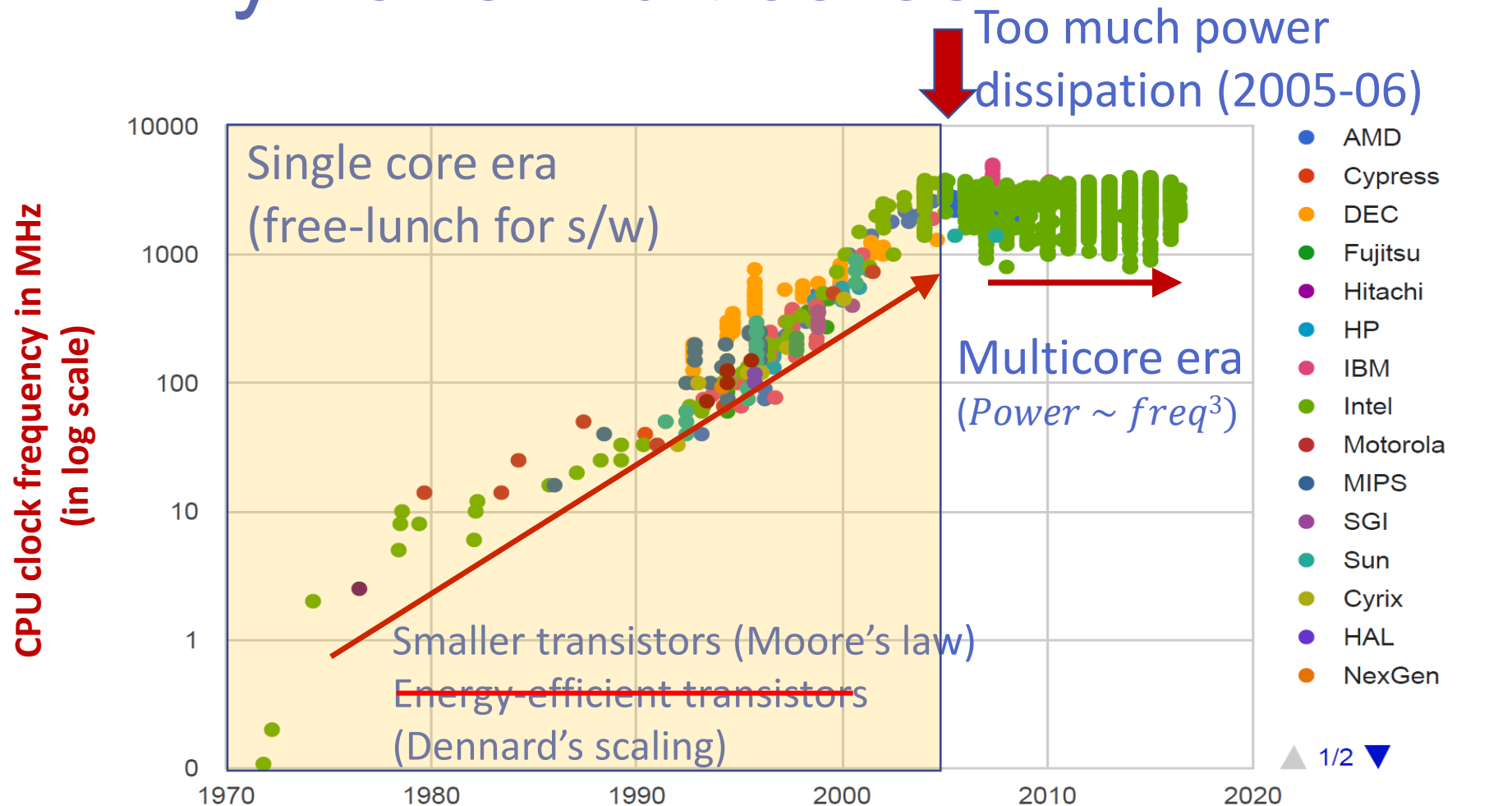AMD Epyc gen 2 (server class) 64 cores

# Why have Multicores?



CPU clock frequency in MHz (in log scale)

Legend: AMD, Cypress, DEC, Fujitsu, Hitachi, HP, IBM, Intel, Motorola, MIPS, SGI, Sun, Cyrix, HAL, NexGen

1/2

Historical data from Stanford CUPDB

# Why have Multicores?



Single core era (free-lunch for s/w)

CPU clock frequency in MHz (in log scale)

Smaller transistors (Moore's law)
Energy-efficient transistors (Dennard's scaling)

Legend: AMD, Cypress, DEC, Fujitsu, Hitachi, HP, IBM, Intel, Motorola, MIPS, SGI, Sun, Cyrix, HAL, NexGen

1/2

Historical data from Stanford CUPDB

# Why have Multicores?

Too much power dissipation (2005-06)

Single core era (free-lunch for s/w)

Multicore era
($Power \sim freq^3$)

Smaller transistors (Moore's law)

~~Energy-efficient transistors~~ (Dennard's scaling)

**CPU clock frequency in MHz (in log scale)**

Legend:
- AMD
- Cypress
- DEC
- Fujitsu
- Hitachi
- HP
- IBM
- Intel
- Motorola
- MIPS
- SGI
- Sun
- Cyrix
- HAL
- NexGen

△ 1/2 ▼

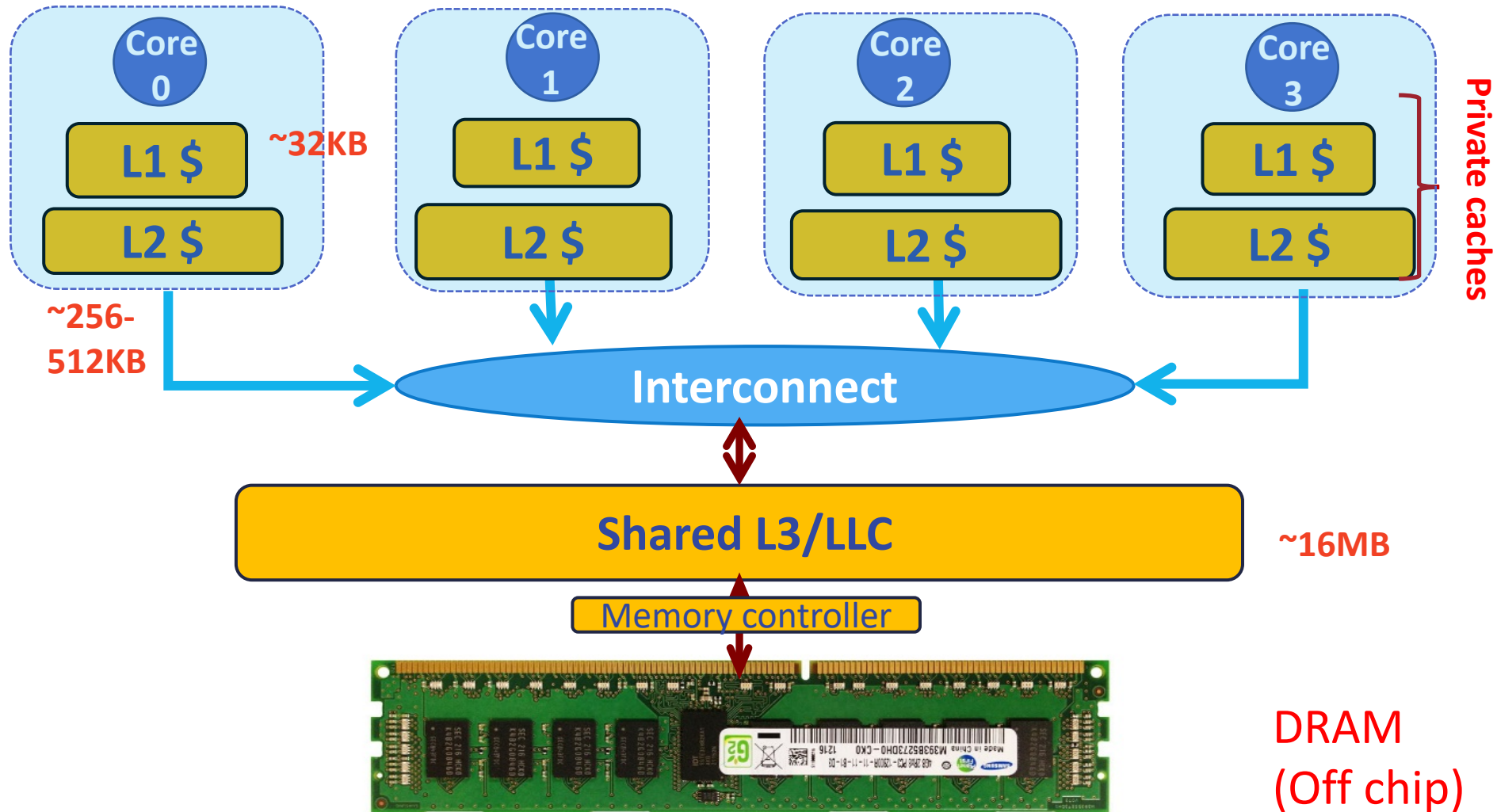Historical data from Stanford CUPDB

Multicores was a compulsion and not a choice
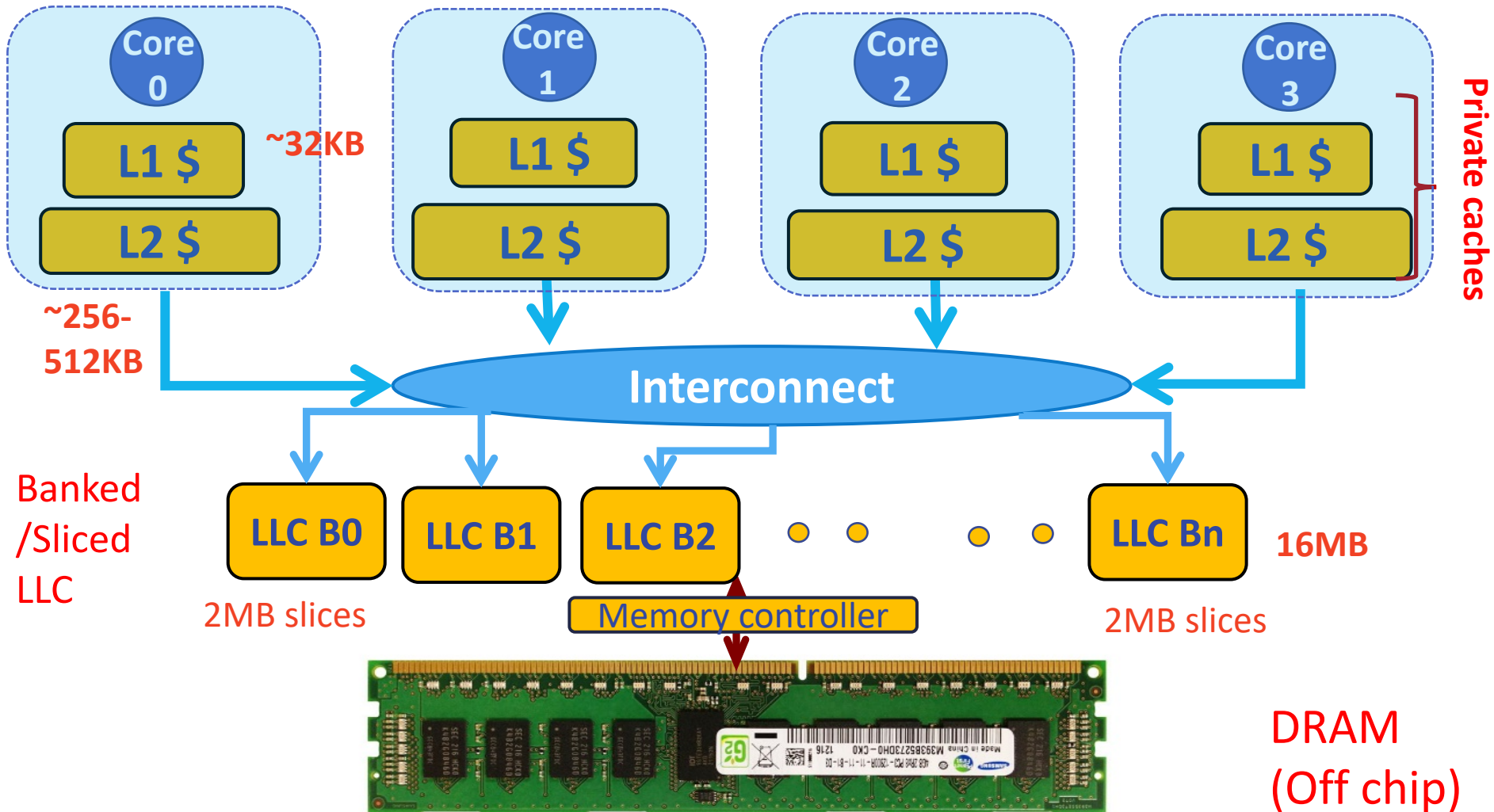
# Schematic: Physical connection



DRAM
(Off chip)

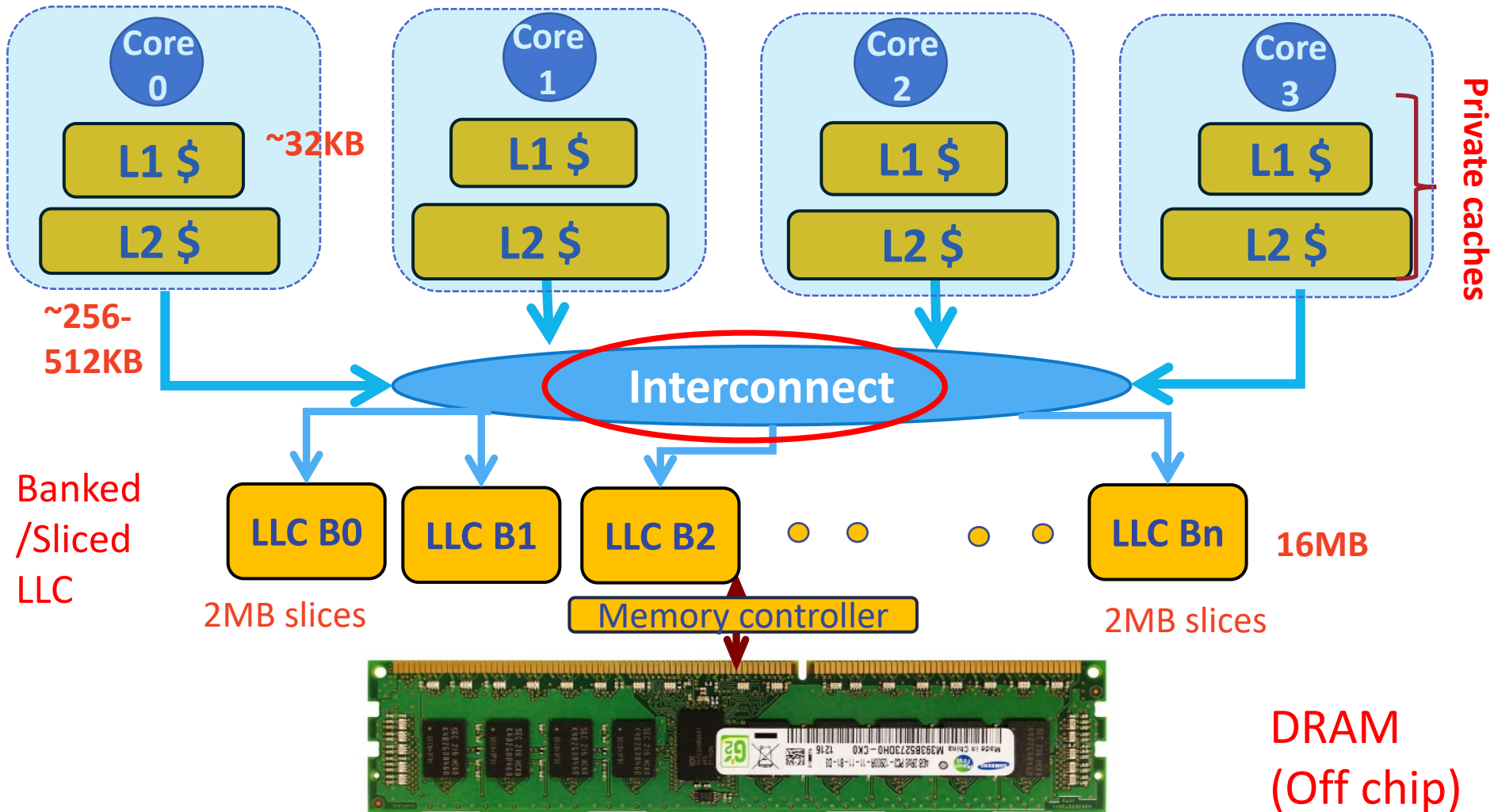# Schematic: Physical connection

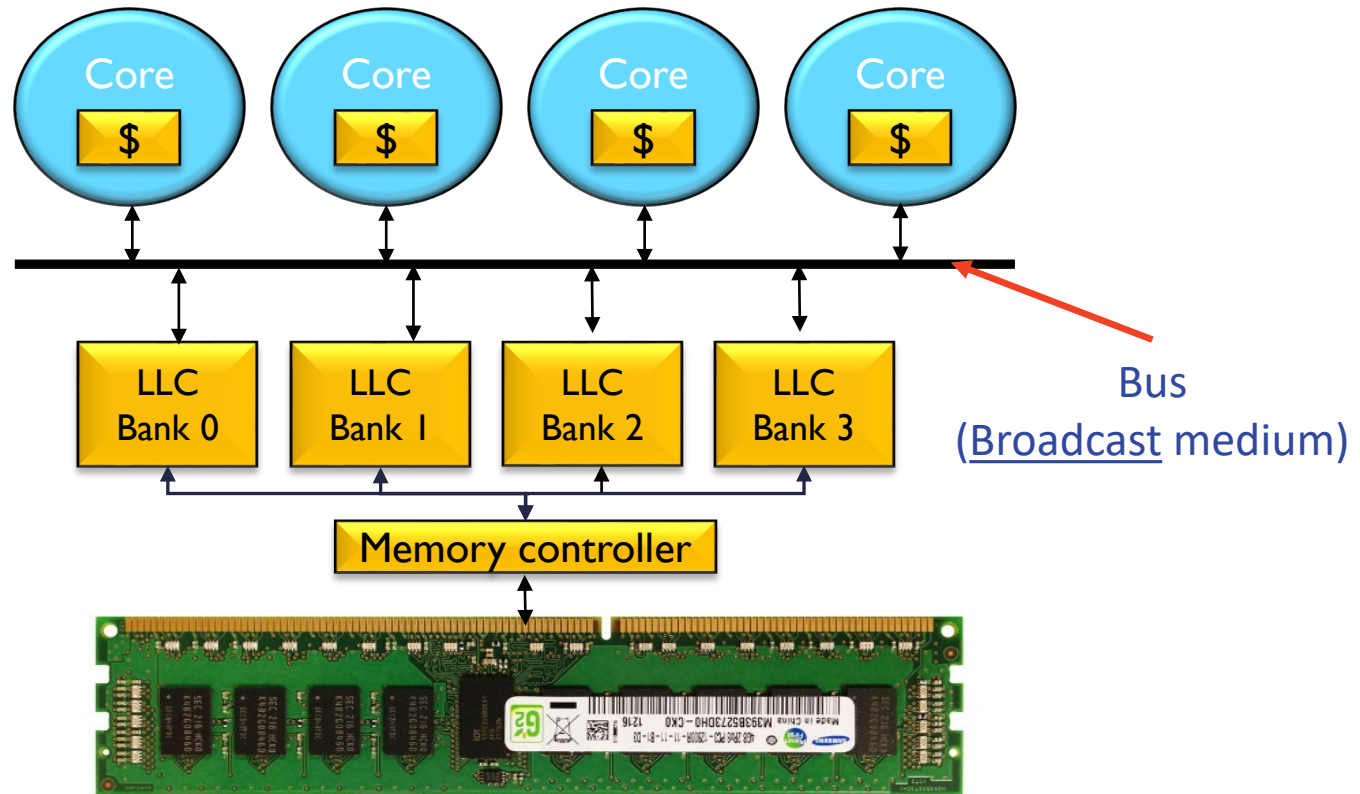# Schematic: Physical connection

# Schematic: Physical connection
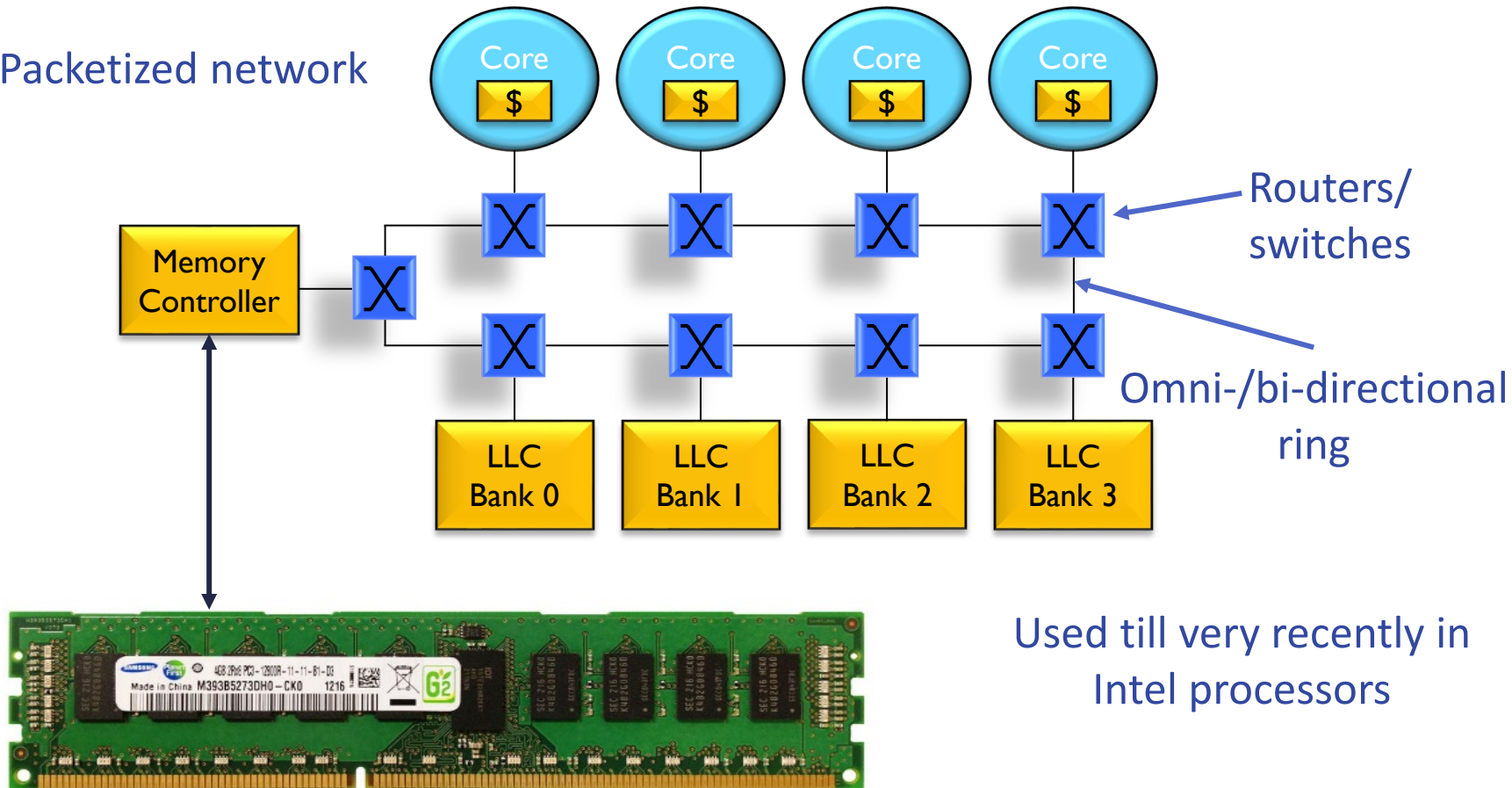
# Choice of Interconnects: Bus



Bus
(<u>Broadcast</u> medium)

Good for small number of cores

# Choice of Interconnects: Ring



Packetized network

Routers/ switches

Omni-/bi-directional ring

# Choice of Interconnects: Ring



Packetized network
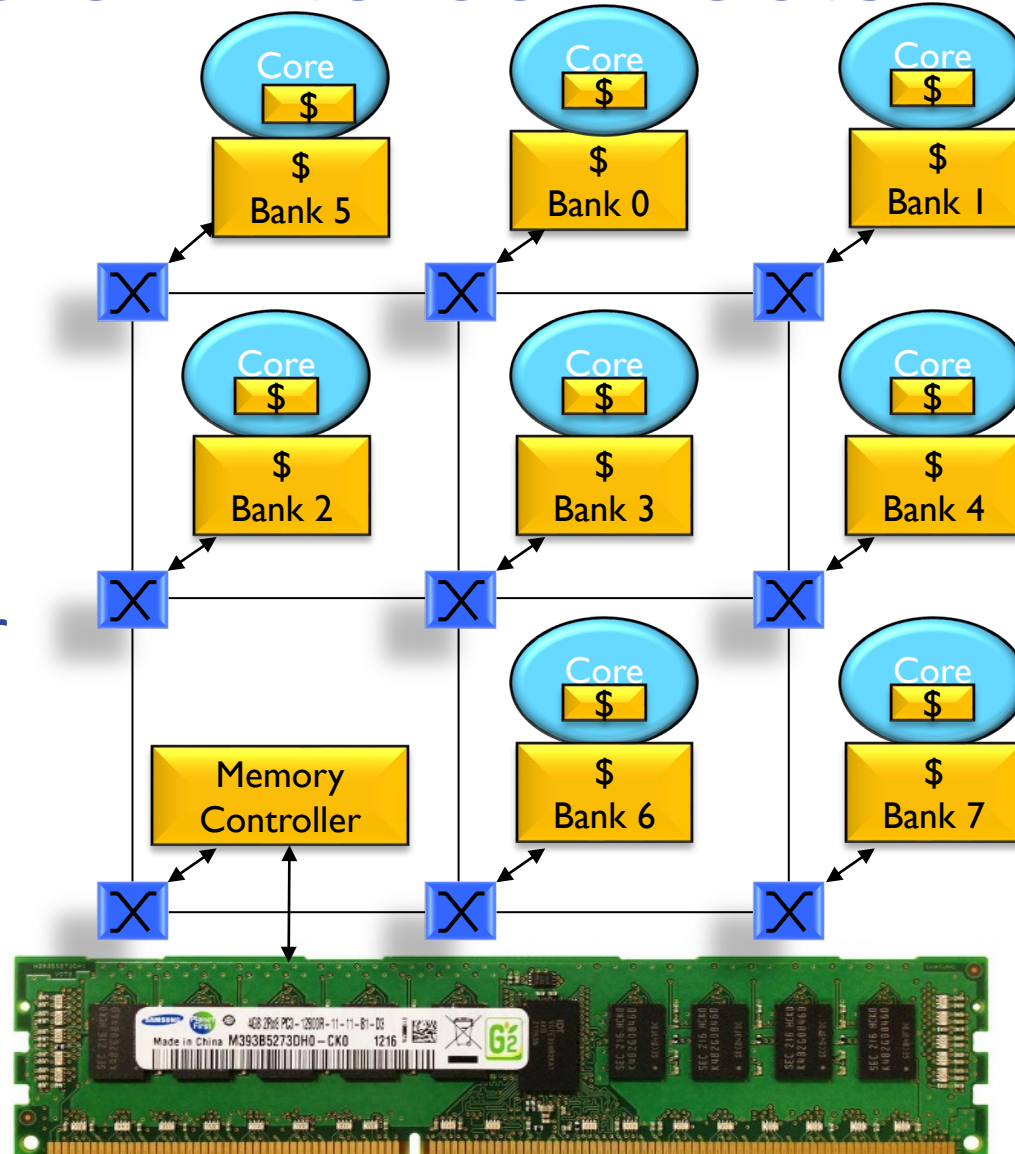
Routers/switches

Omni-/bi-directional ring

Used till very recently in Intel processors

Can be used as broadcast medium; good for moderate number of cores

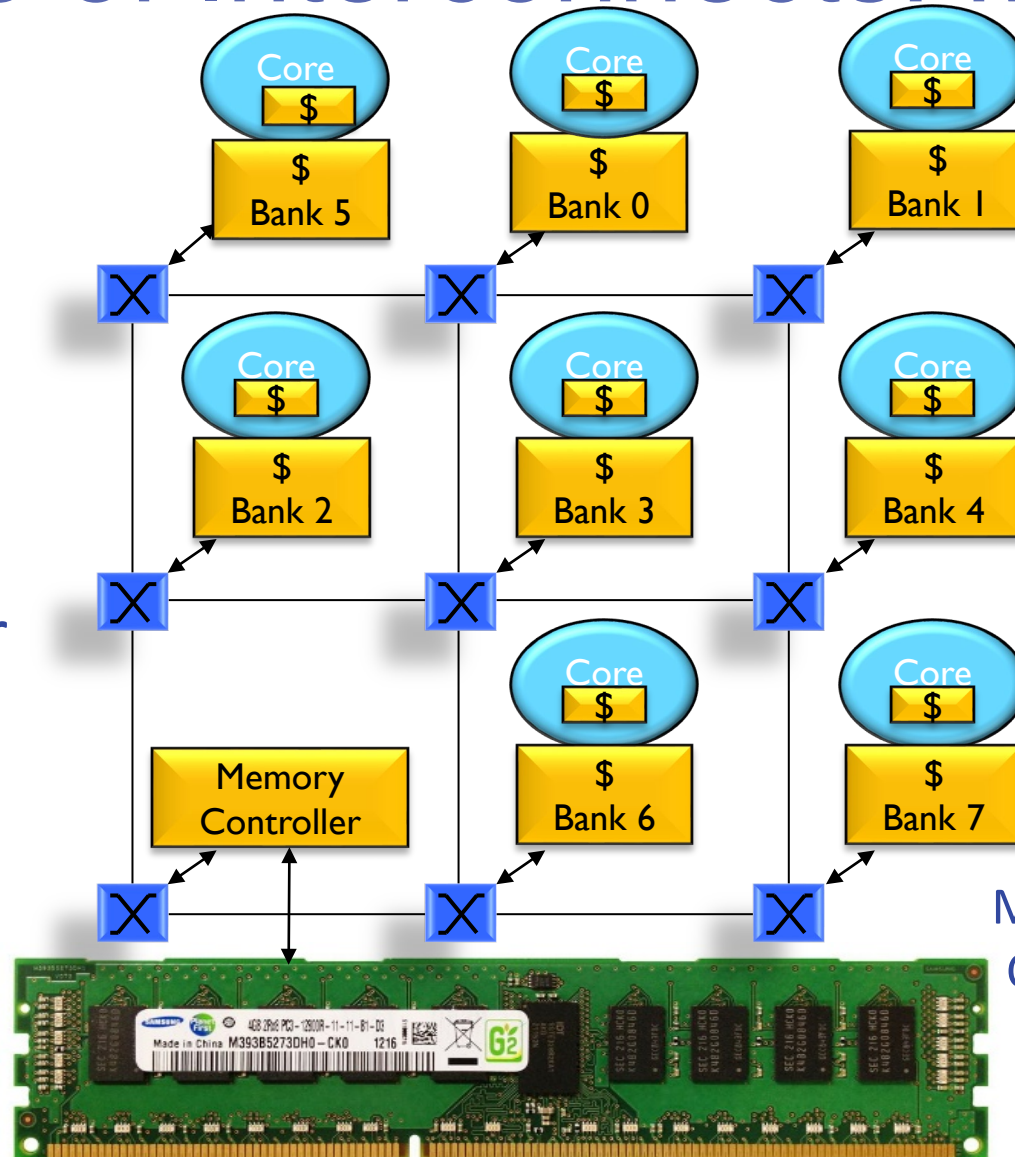# Choice of Interconnects: Mesh



Scalable,
point-to-point,

Better for larger
number f cores

Latest Intel
server processors
uses mesh

# Choice of Interconnects: Mesh

Scalable, point-to-point,

Better for larger number f cores
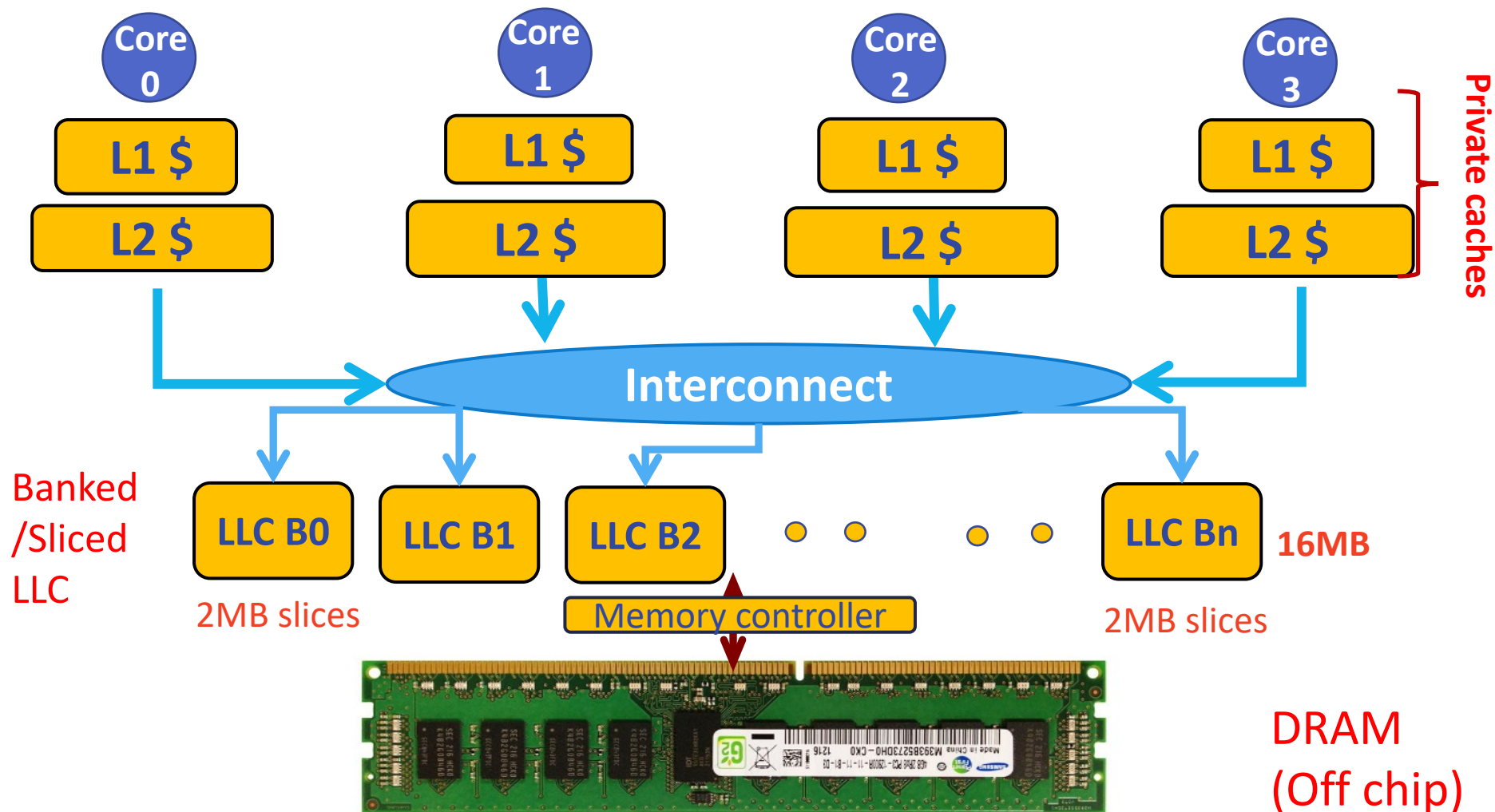
Latest Intel server processors uses mesh

Many other options: Crossbar, torus etc.

# Software's view of multicore

Need for cache coherence

# Software's view of multicore

# Software's view of multicore

**Private caches**

| Core 0 | Core 1 | Core 2 | Core 3 |
|--------|--------|--------|--------|
| L1 $ | L1 $ | L1 $ | L1 $ |
| L2 $ | L2 $ | L2 $ | L2 $ |

**Interconnect**

Banked /Sliced LLC

LLC B0    LLC B1    LLC B2    •  •    •  •    LLC Bn    **16MB**

2MB slices

Memory controller

2MB slices

DRAM (Off chip)

# Software's view of multicore



Core 0 Core 1 Core 2 Core 3

L1 $ L1 $ L1 $ L1 $

L2 $ L2 $ L2 $ L2 $

Private caches

Memory

Interconnect

Banked /Sliced LLC

LLC B0 LLC B1 LLC B2 LLC Bn 16MB

2MB slices 2MB slices

Memory controller

DRAM (Off chip)

# Cache In-coherence problem

- Variable A initially has value 0
- C1 stores value 1 into A
- C2 loads A from memory and sees old value 0



**Private caches**

C1 → A: 0 (L1)

C2 → (L1)

Bus

A: 0

Shared LLC/ Main Memory

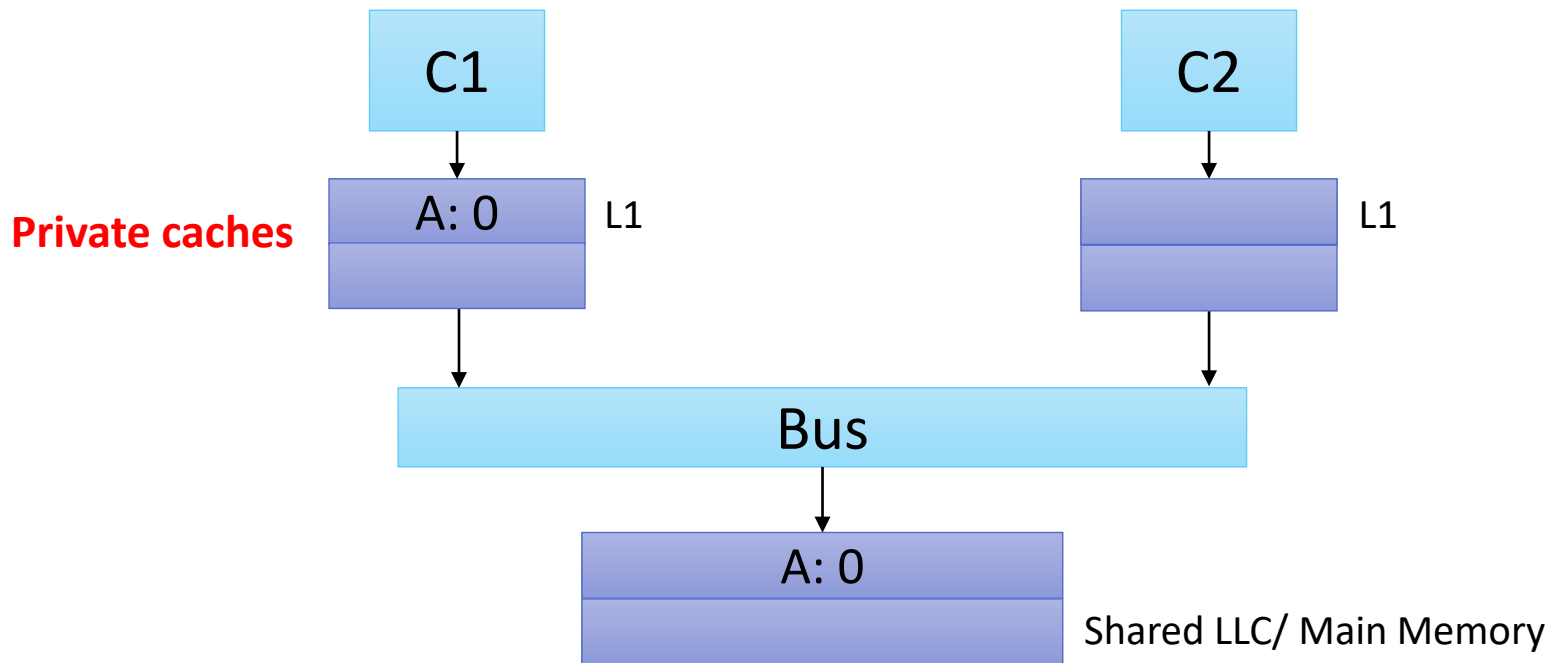Simplified view of a multicore

# Cache In-coherence problem

- Variable A initially has value 0
- C1 stores value 1 into A
- C2 loads A from memory and sees old value 0

t1: Store A=1

C1

C2

**Private caches**

A: 0 1    L1

L1

Bus

A: 0
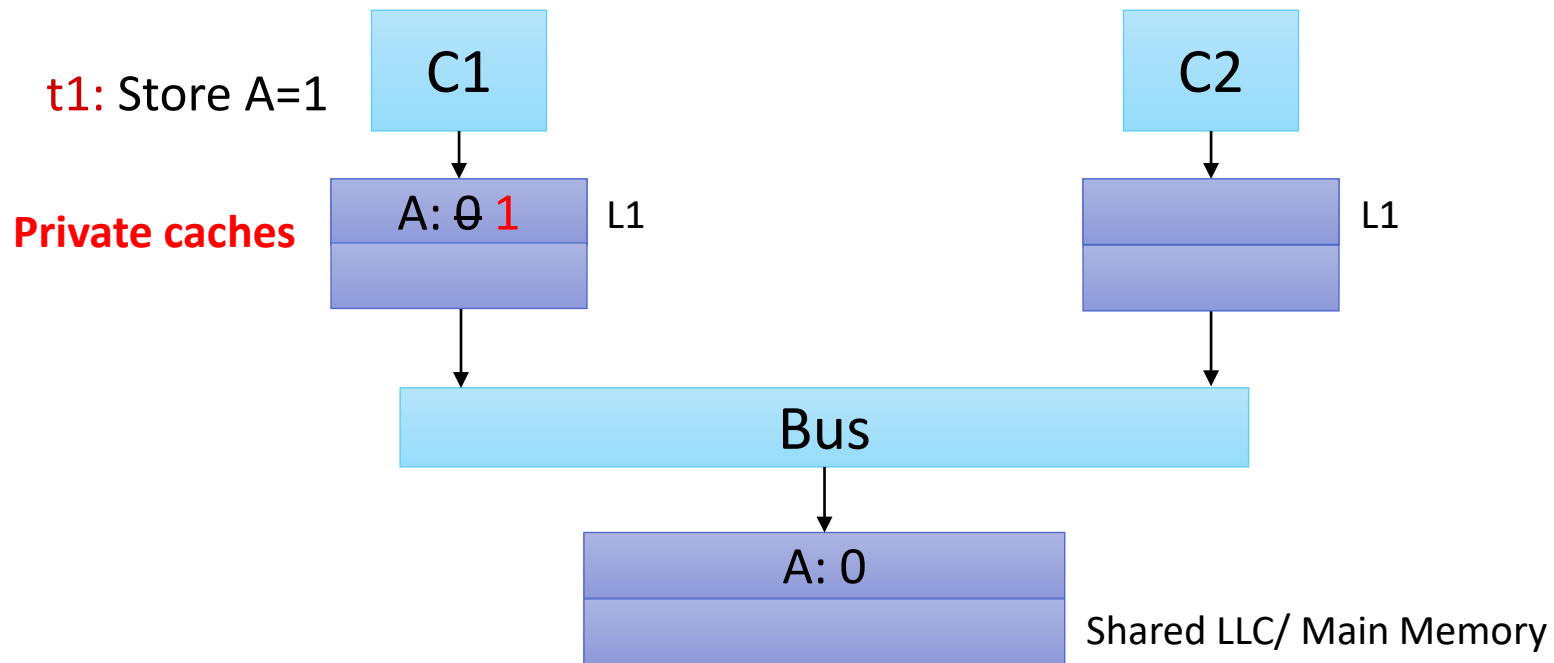
Shared LLC/ Main Memory

Simplified view of a multicore

# Cache In-coherence problem

- Variable A initially has value 0
- C1 stores value 1 into A
- C2 loads A from memory and sees old value 0

t1: Store A=1    C1          C2    t2: Load A?

**Private caches**    A: 0 1    L1              L1

Bus

A: 0
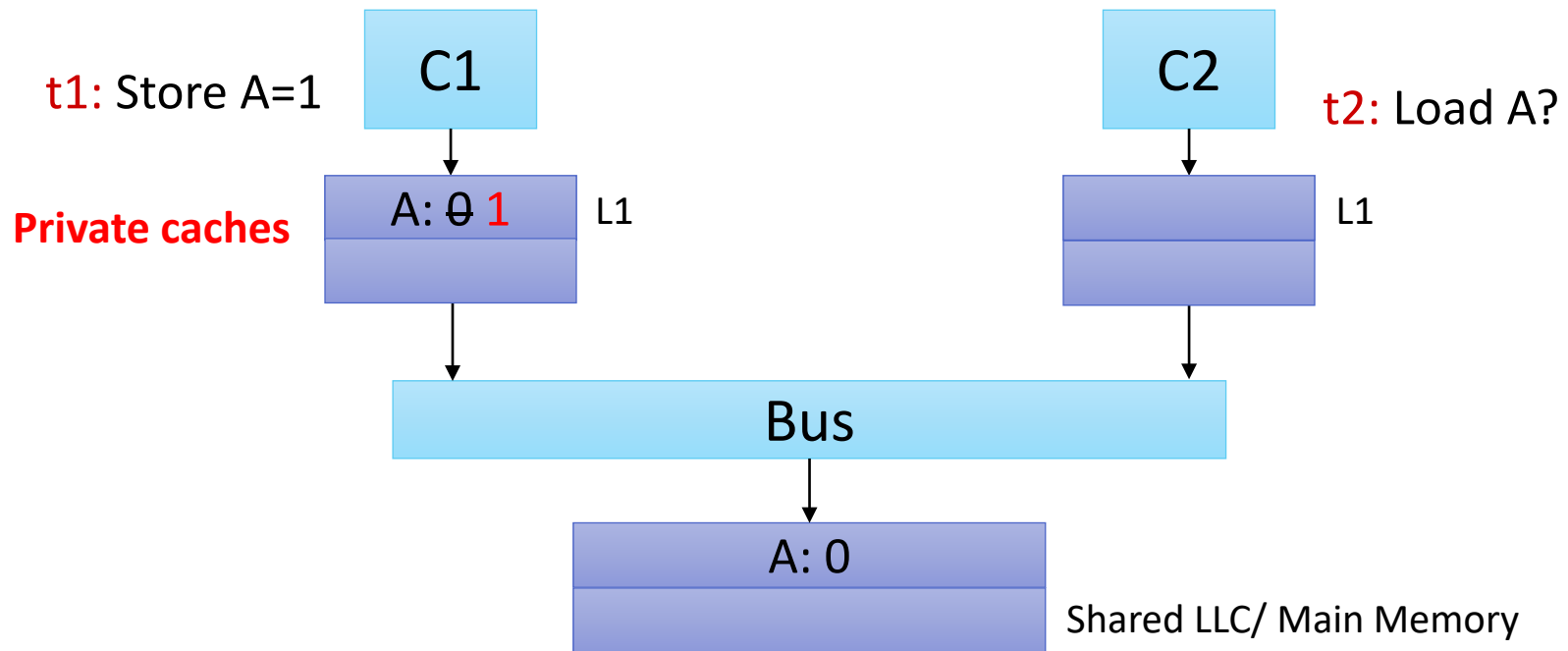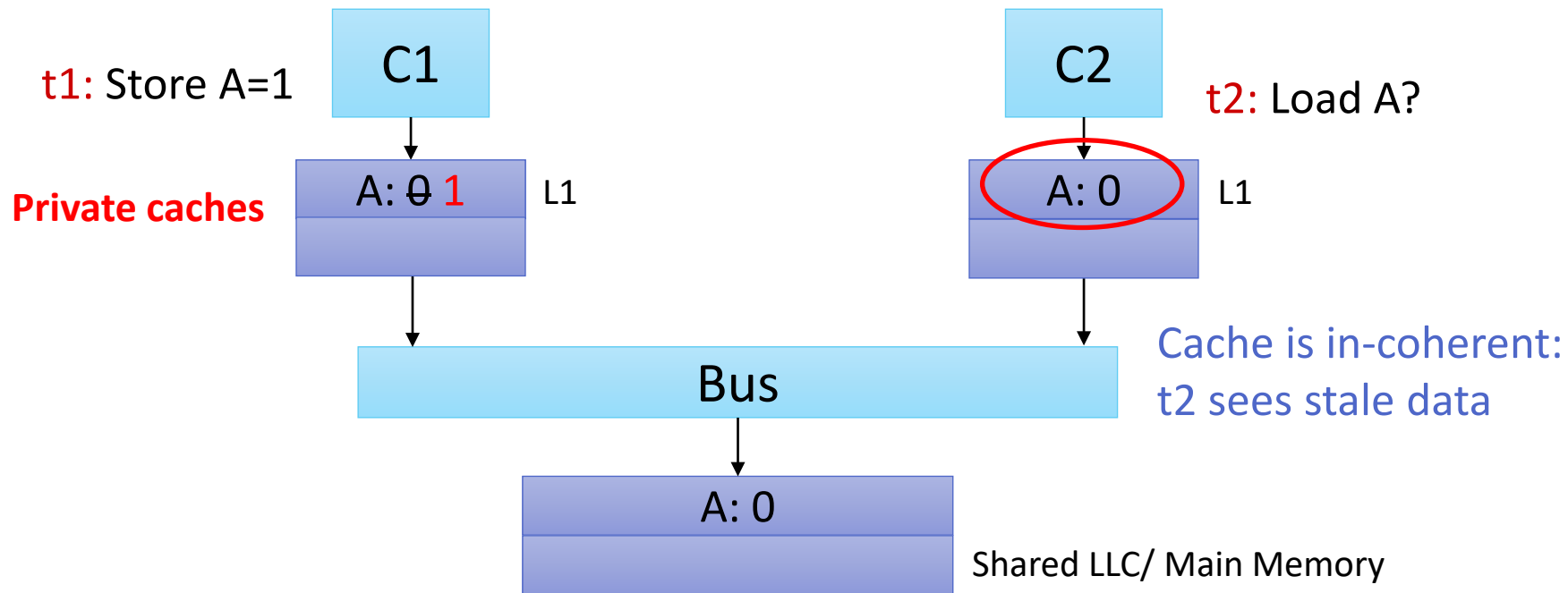
Shared LLC/ Main Memory

Simplified view of a multicore

# Cache In-coherence problem

- Variable A initially has value 0
- C1 stores value 1 into A
- C2 loads A from memory and sees old value 0

t1: Store A=1    C1                    C2    t2: Load A?

**Private caches**    A: 0̶ 1    L1        A: 0    L1

Bus                            Cache is in-coherent:
                               t2 sees stale data

A: 0

Shared LLC/ Main Memory

Simplified view of a multicore

# Goal of Cache coherence

- Keep contents of private caches coherent
  - ▸ Software should get the "latest" data

- Cache coherence subsystem orders (total order) all writes to a given memory address/location
  - ▸ It does it for all memory address/location

# Goal of Cache coherence

- Keep contents of private caches coherent
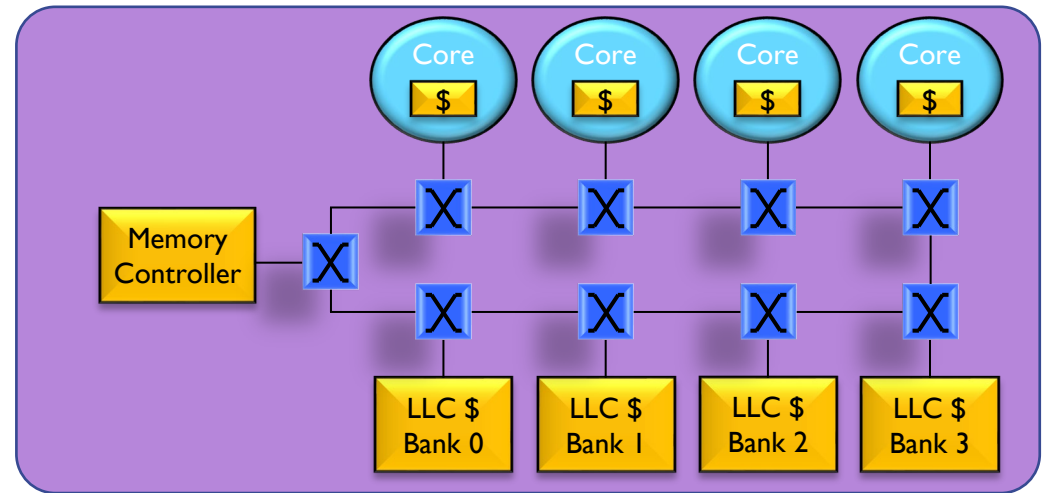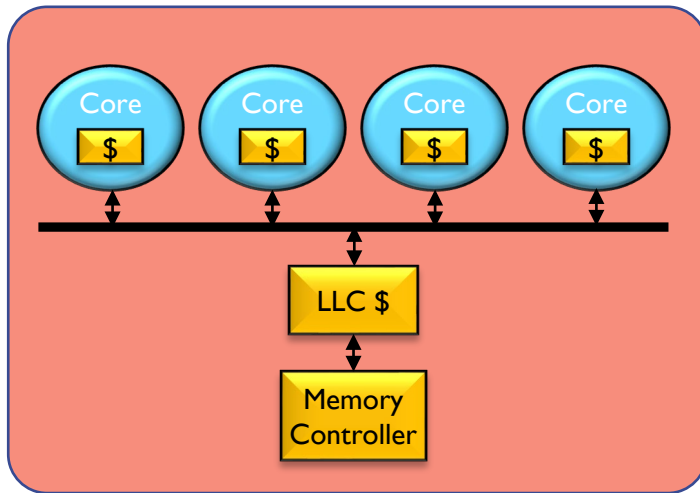  - ‣ Software should get the "latest" data

- Cache coherence subsystem orders (total order) all writes to a given memory address/location
  - ‣ It does it for all memory address/location

- Two design classes of cache coherence subsystems
  - ‣ Snoopy coherence protocol ⟵
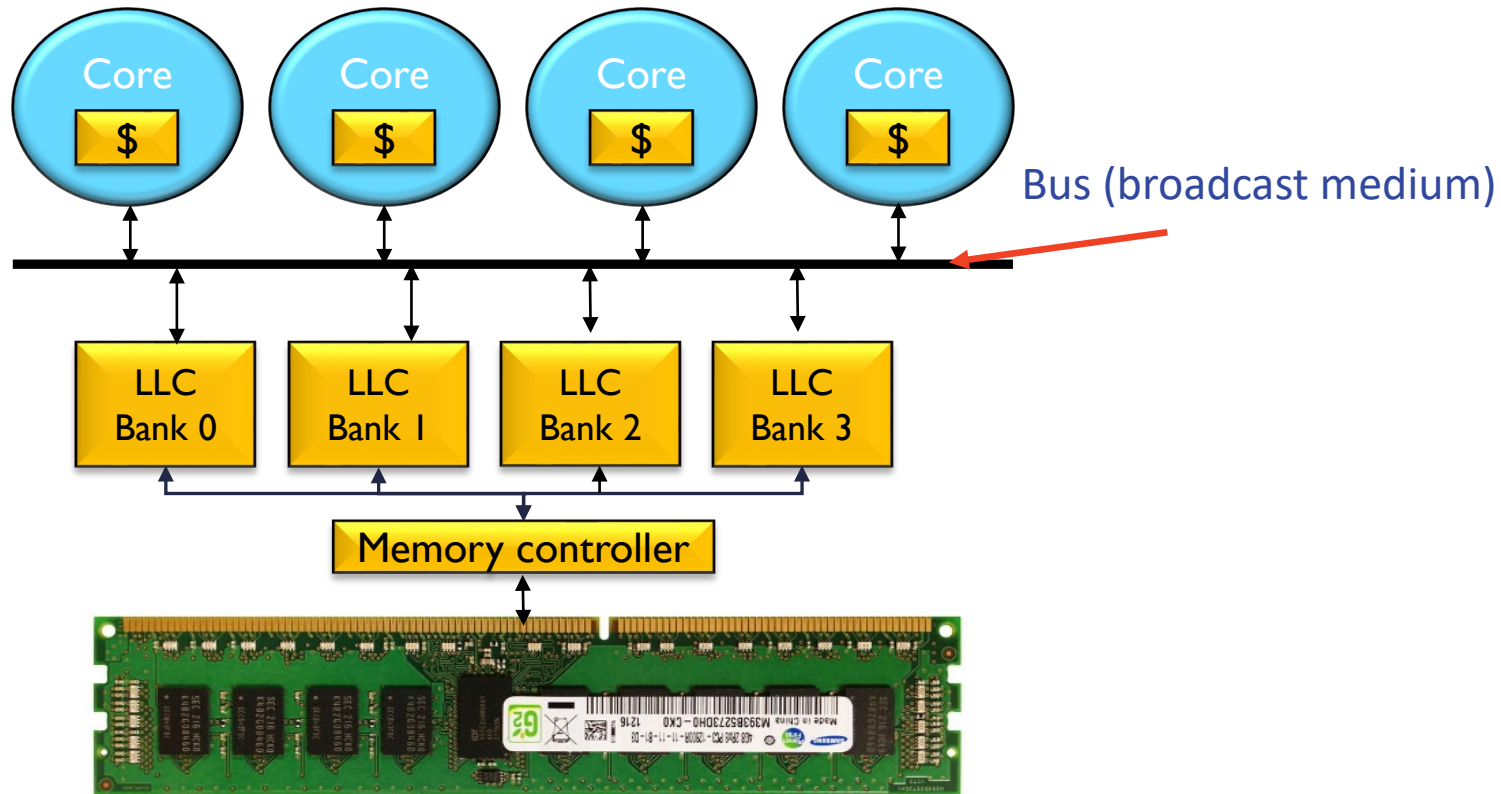  - ‣ Directory coherence protocol  (if time permits)

# Snoopy cache coherence

- Relies on broadcast-based interconnection network
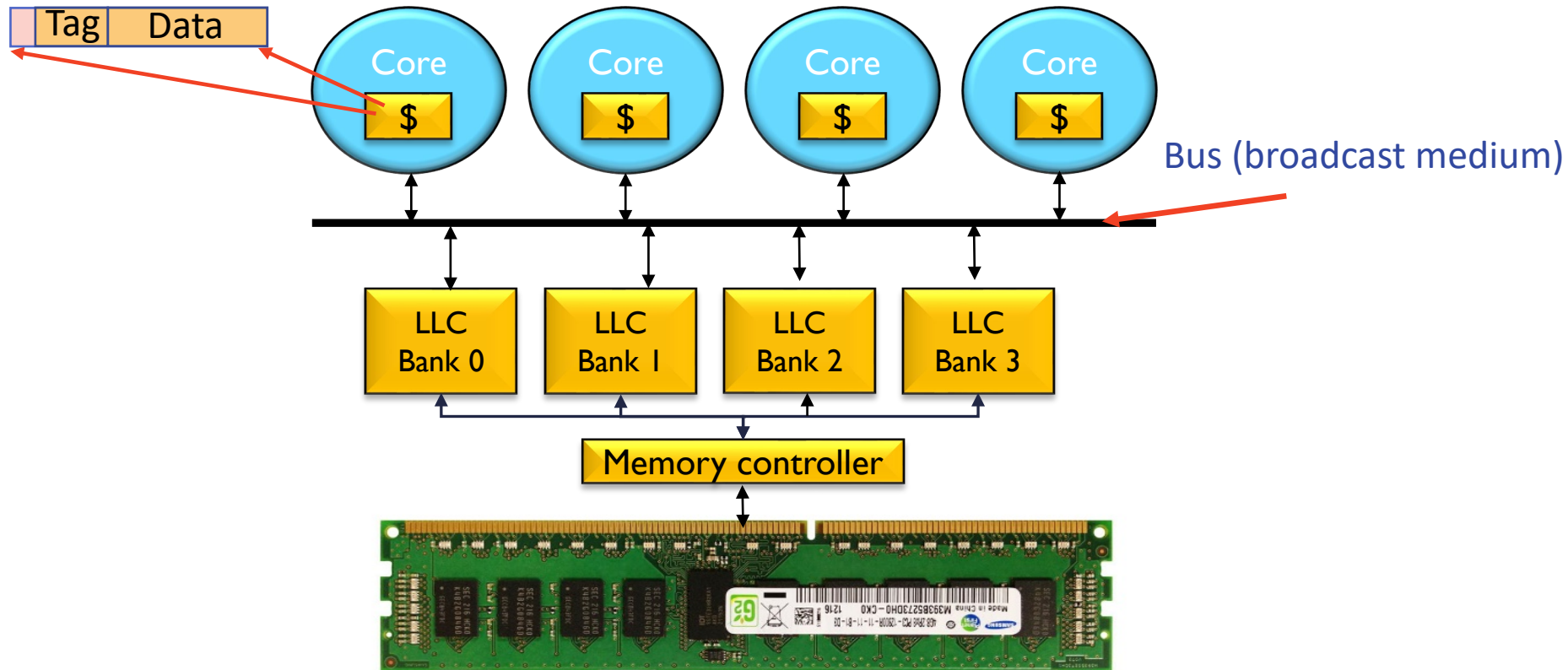  - Typically Bus or Ring



- All caches must monitor (aka "*snoop*") all traffic
  - And keep track of cache line *states* based on the observed traffic

# Snoopy cache coherence



Bus (broadcast medium)

Example: MSI coherence protocol

# Snoopy cache coherence



Example: MSI coherence protocol

# Snoopy cache coherence

Coherence state of the cache block

| Tag | Data |

**M** = Modified (Read/Write)

**S** = Shared (Read-only)

**I** = Invalid (Not-present/ No permission)

Core — $
Core — $
Core — $
Core — $

Bus (broadcast medium)

LLC Bank 0
LLC Bank 1
LLC Bank 2
LLC Bank 3

Memory controller

Example: MSI coherence protocol

# Snoopy cache coherence

Coherence state of the cache block

| Tag | Data |

Coherence controller (CC)

Core $

Core $

Core $

Core $

Bus (broadcast medium)

**M** = Modified (Read/Write)

**S** = Shared (Read-only)

**I** = Invalid (Not-present/ No permission)

| LLC Bank 0 | LLC Bank 1 | LLC Bank 2 | LLC Bank 3 |

Memory controller

Role of CC:
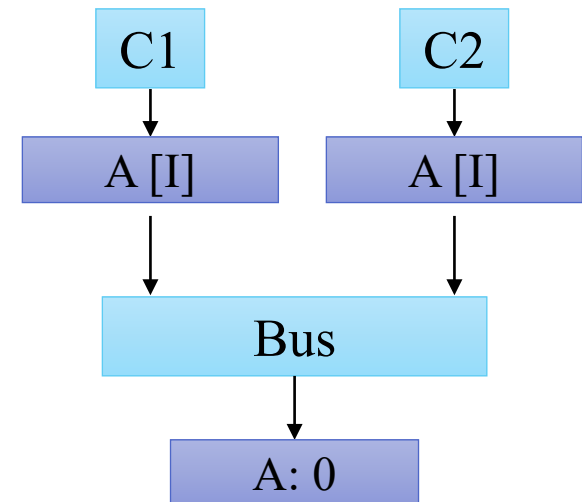(1) Snoop on every transaction on bus.
(2) Respond to bus transactions
(3) Change coherence state or invalidate local $ block
(4) Observe local core's ld/st and (possibly) initiate bus transactions [Implements an FSM]

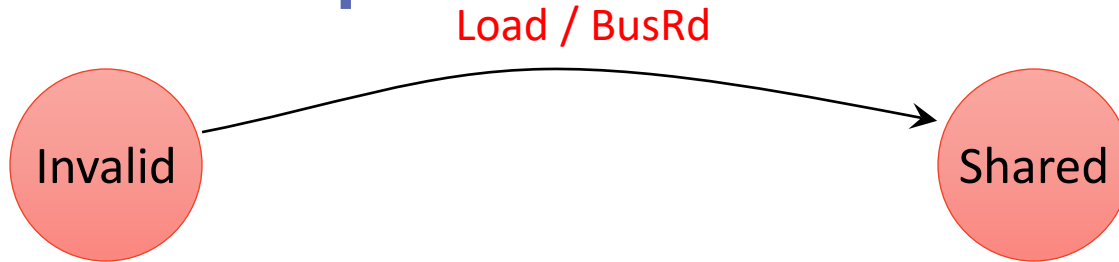Example: MSI coherence protocol

# Ex: MSI protocol in action (1)

Invalid

⟶ **Transition caused by local action**

C1      C2

A [I]      A [I]

Bus

A: 0

Invariant: Only one writer (M) per block; Many sharers (S) okay

# Ex: MSI protocol in action (1)

Load / BusRd

Invalid → Shared

→ **Transition caused by local action**

1: Load A

C1      C2

A [I]    A [I]

Bus

A: 0

Invariant: Only one writer (M) per block; Many sharers (S) okay

# Ex: MSI protocol in action (1)

Load / BusRd

Invalid → Shared

→ **Transition caused by local action**

1: Load A

C1    C2

A [I]    A [I]

2: BusRd A

Bus

A: 0

Invariant: Only one writer (M) per block; Many sharers (S) okay

# Ex: MSI protocol in action (1)

Load / BusRd

Invalid → Shared

⟶ **Transition caused by local action**

1: Load A

C1        C2

A [I̶ S]: 0        A [I]

2: BusRd A

Bus

3: BusReply A        A: 0

Invariant: Only one writer (M) per block; Many sharers (S) okay

# Ex: MSI protocol in action (2)

Load / BusRd

Invalid → Shared

→ **Transition caused by local action**

- - - → **Transition caused by bus message**

C1    C2

A [S]: 0    A [I]

Bus

A: 0

Invariant: Only one writer (M) per block; Many sharers (S) okay

# Ex: MSI protocol in action (2)

Load / BusRd

Invalid → Shared

Load / --

→ **Transition caused by local action**

- - → **Transition caused by bus message**

1: Load A

C1    C2

A [S]: 0    A [I]

Bus

A: 0

Invariant: Only one writer (M) per block; Many sharers (S) okay

# Ex: MSI protocol in action (2)



Invariant: Only one writer (M) per block; Many sharers (S) okay

# Ex: MSI protocol in action (2)



**Invariant: Only one writer (M) per block; Many sharers (S) okay**

# Ex: MSI protocol in action(3)



**Invariant: Only one writer (M) per block; Many sharers (S) okay**

# Ex: MSI protocol in action(3)



**Invariant: Only one writer (M) per block; Many sharers (S) okay**

# Ex: MSI protocol in action (4)

Load / BusRd

Invalid → Shared

BusRd / [BusReply]

Evict / --

Load / --

C1    C2

A [S]: 0    A [I]

Bus

A: 0

Invariant: Only one writer (M) per block; Many sharers (S) okay

# Ex: MSI protocol in action (4)



**Invariant: Only one writer (M) per block; Many sharers (S) okay**

# Ex: MSI protocol in action (4)



**Invariant: Only one writer (M) per block; Many sharers (S) okay**

# Ex: MSI protocol in action (4)



**Invariant: Only one writer (M) per block; Many sharers (S) okay**

# Ex: MSI protocol in action (4)



**Invariant: Only one writer (M) per block; Many sharers (S) okay**

# Ex: MSI protocol in action (4)



**Invariant: Only one writer (M) per block; Many sharers (S) okay**

# Ex: MSI protocol in action (4)



Invariant: Only one writer (M) per block; Many sharers (S) okay

# Ex: MSI protocol in action (4)



Invariant: Only one writer (M) per block; Many sharers (S) okay

# Ex: MSI protocol in action (5)



**Invariant: Only one writer (M) per block; Many sharers (S) okay**

# Ex: MSI protocol in action (5)

Load / BusRd

Invalid

BusRdX / [BusReply]

Evict / --

Shared

BusRd / [BusReply]

Load / --

Store / BusRdX

Modified

Load, Store / --

1: Load A

C1

C2

A [I]

A [M]: 1

Bus

A: 0

**Invariant: Only one writer (M) per block; Many sharers (S) okay**

# Ex: MSI protocol in action (5)



**Invariant: Only one writer (M) per block; Many sharers (S) okay**

# Ex: MSI protocol in action (5)



**Invariant: Only one writer (M) per block; Many sharers (S) okay**

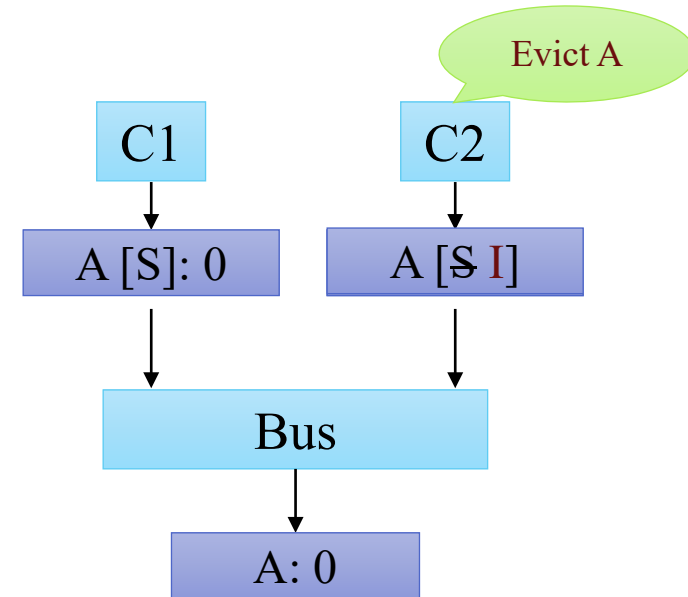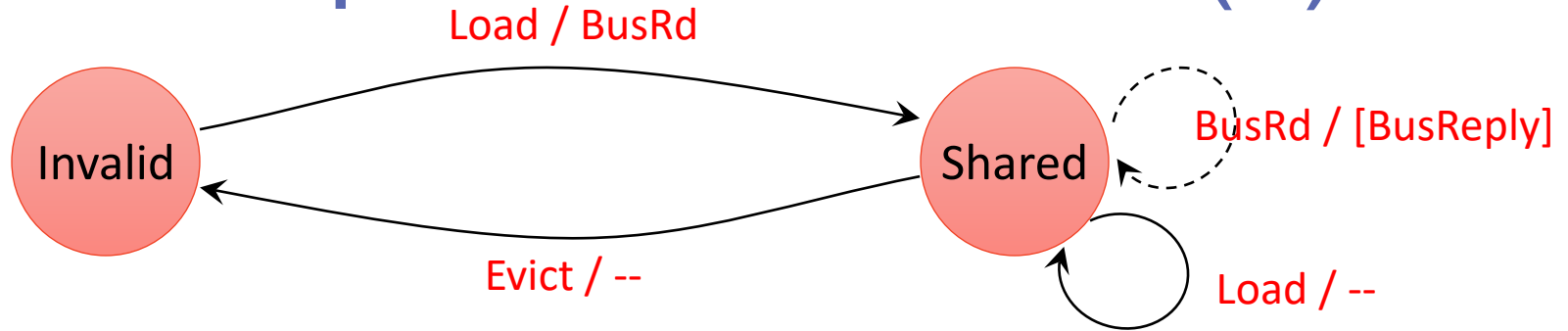# Ex: MSI protocol in action (5)



Invariant: Only one writer (M) per block; Many sharers (S) okay

# Ex: MSI protocol in action (6)

Load / BusRd

Invalid

Shared

BusRdX / [BusReply]

BusRd / [BusReply]

Evict / --

Load / --

Store / BusRdX

BusRd / BusReply

Modified

Load, Store / --

C1

C2

A [S]: 1

A [S]: 1

Bus

A: 1

Invariant: Only one writer (M) per block; Many sharers (S) okay

# Ex: MSI protocol in action (6)



Load / BusRd

BusRdX / [BusReply]

BusRd / [BusReply]

Evict / --

Load / --

Store / BusRdX

BusRd / BusReply

Store / BusInv

Load, Store / --

Invalid

Shared

Modified

1: Store A
aka "*Upgrade*"
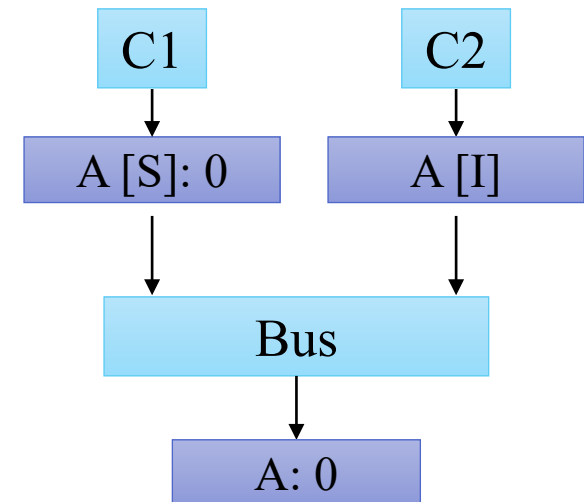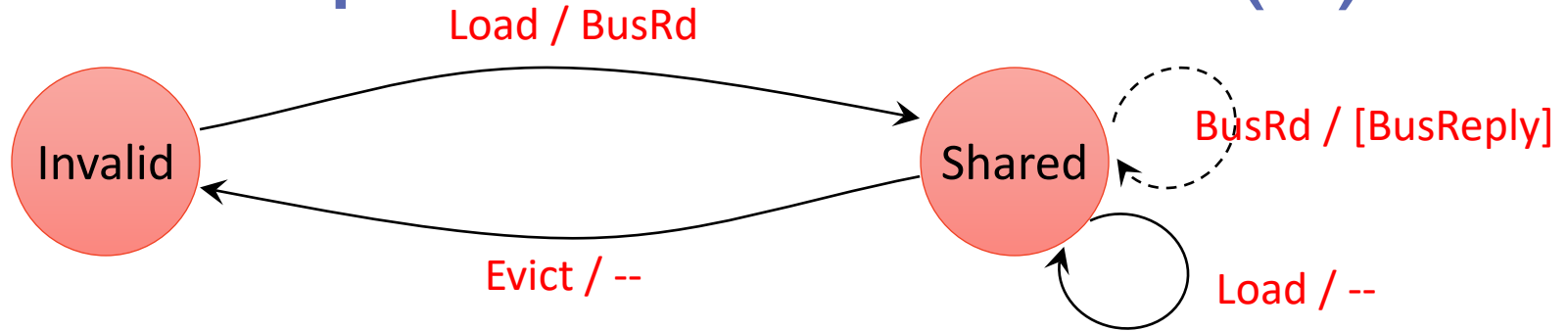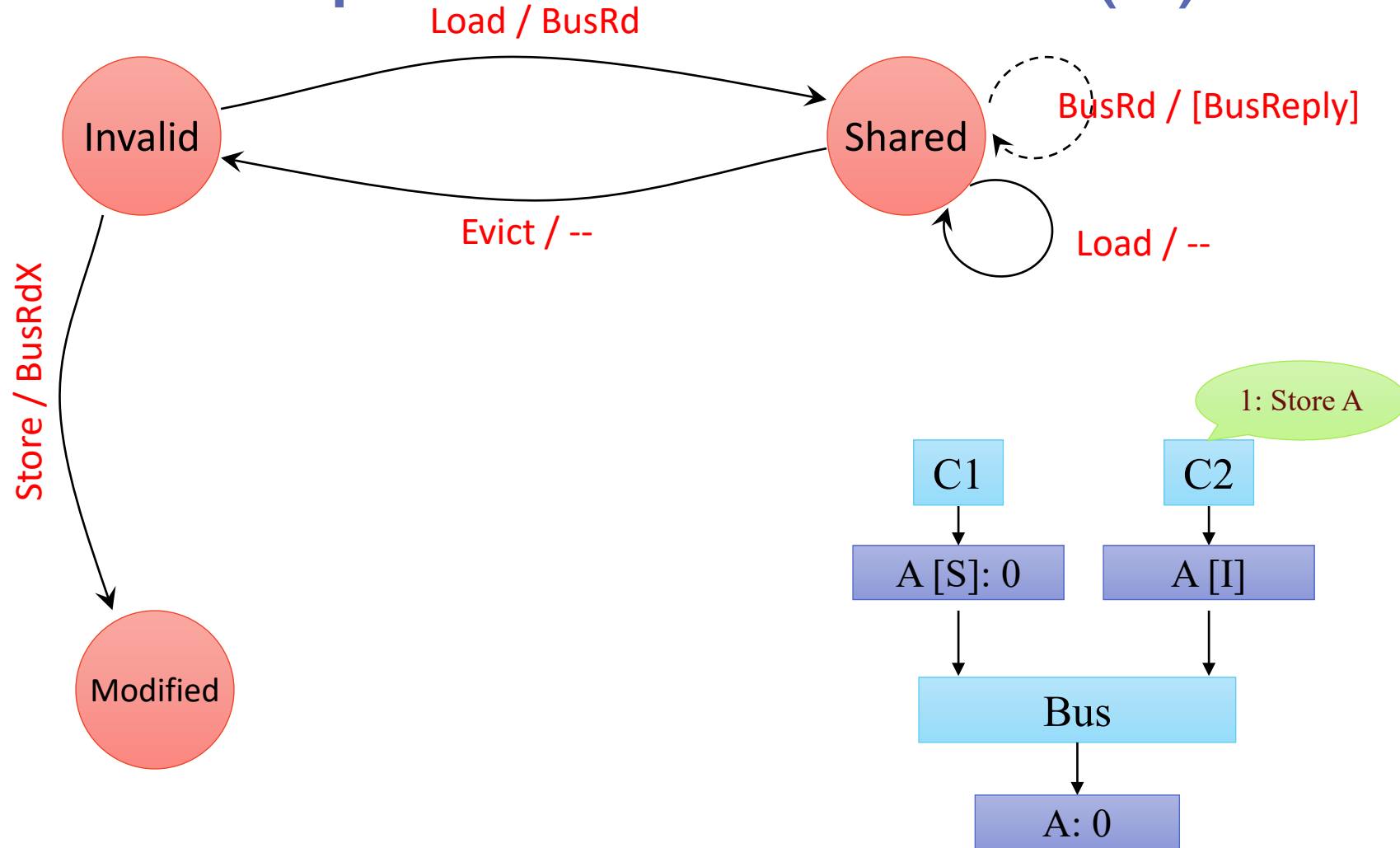
C1

C2

A [S]: 1

A [S]: 1

Bus

A: 1

Invariant: Only one writer (M) per block; Many sharers (S) okay

# Ex: MSI protocol in action (6)



Invariant: Only one writer (M) per block; Many sharers (S) okay
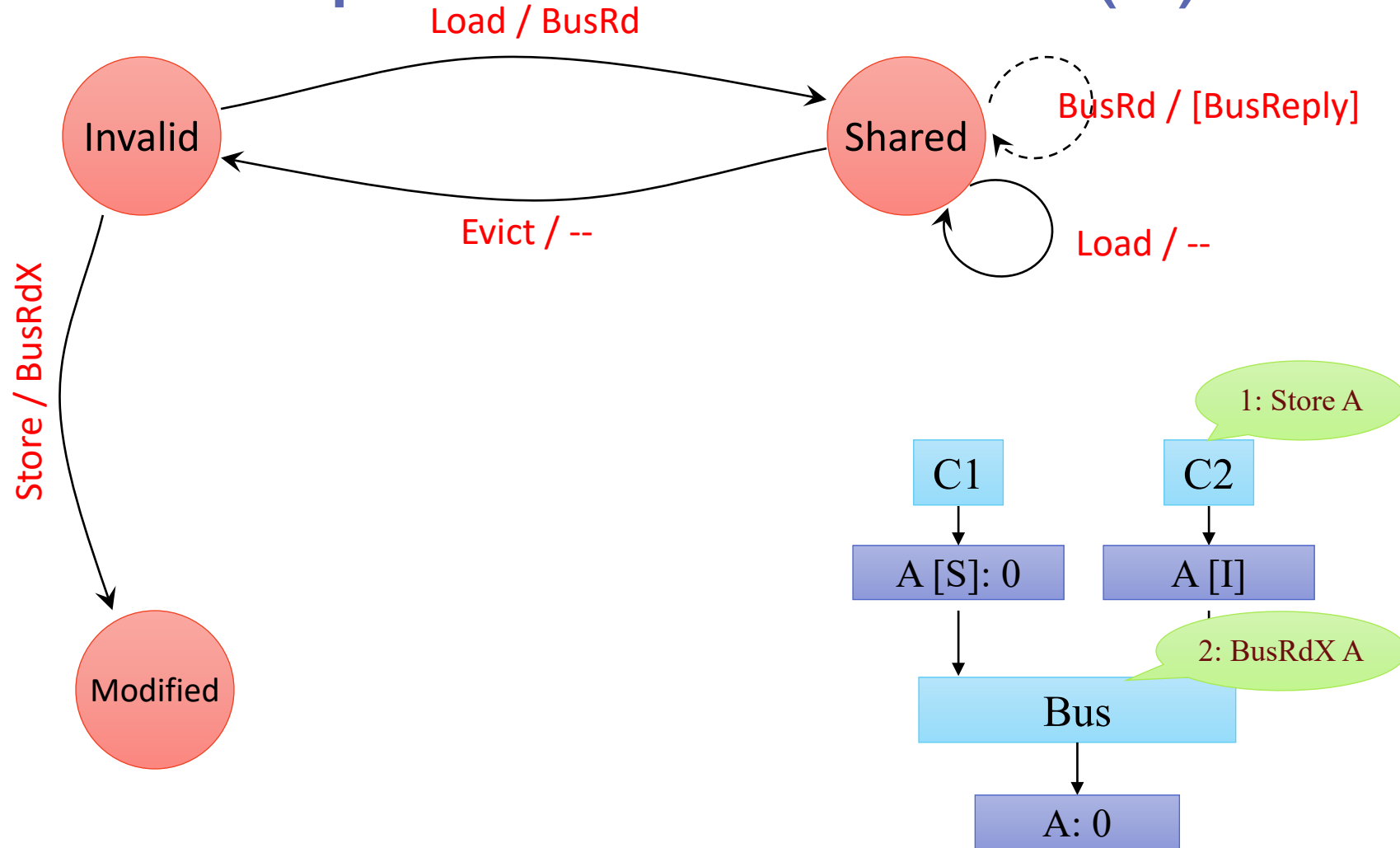
# Ex: MSI protocol in action (6)
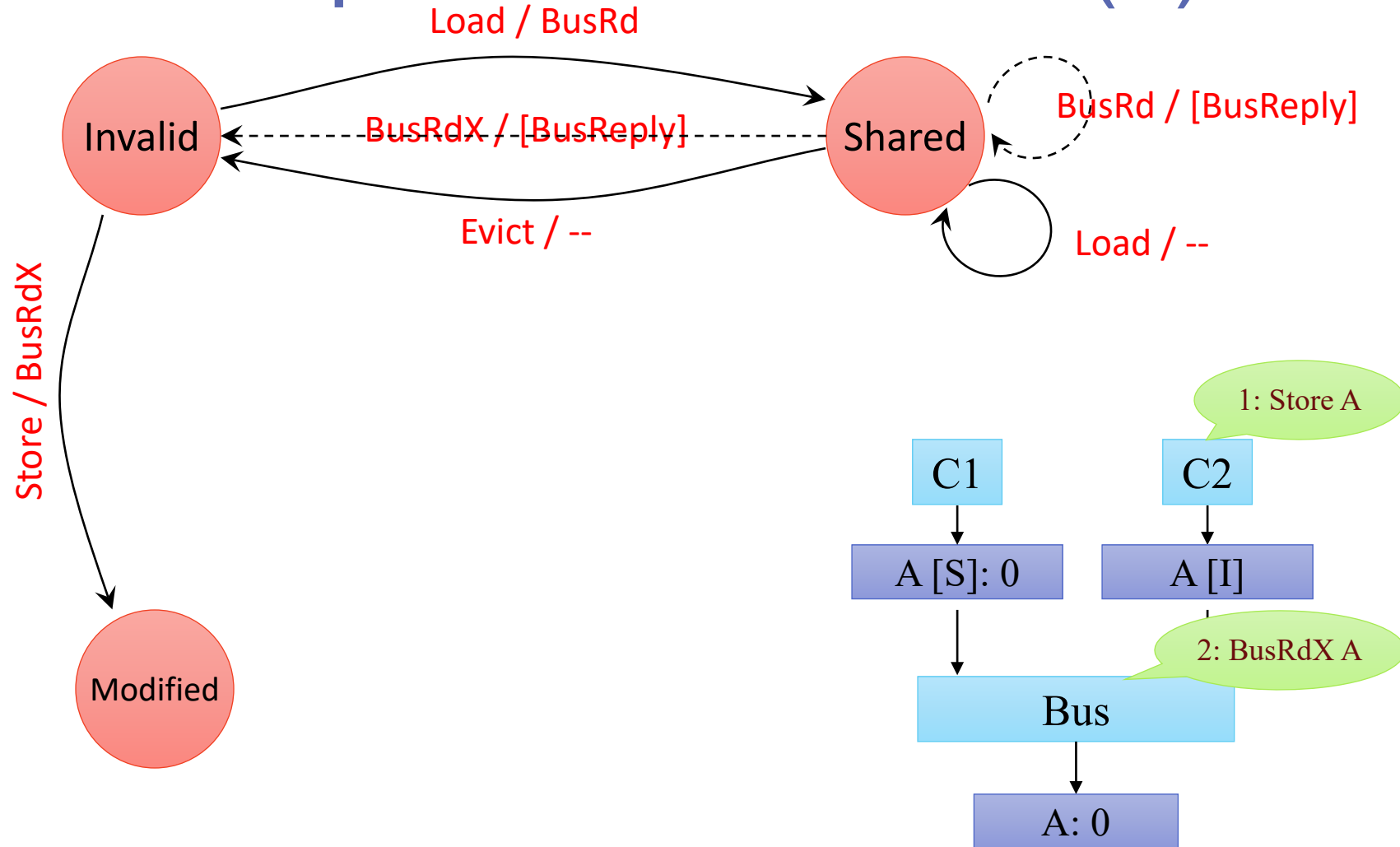


Invariant: Only one writer (M) per block; Many sharers (S) okay

# Ex: MSI protocol in action (6)



**Invariant: Only one writer (M) per block; Many sharers (S) okay**

# Ex: MSI protocol in action (7)



Invariant: Only one writer (M) per block; Many sharers (S) okay
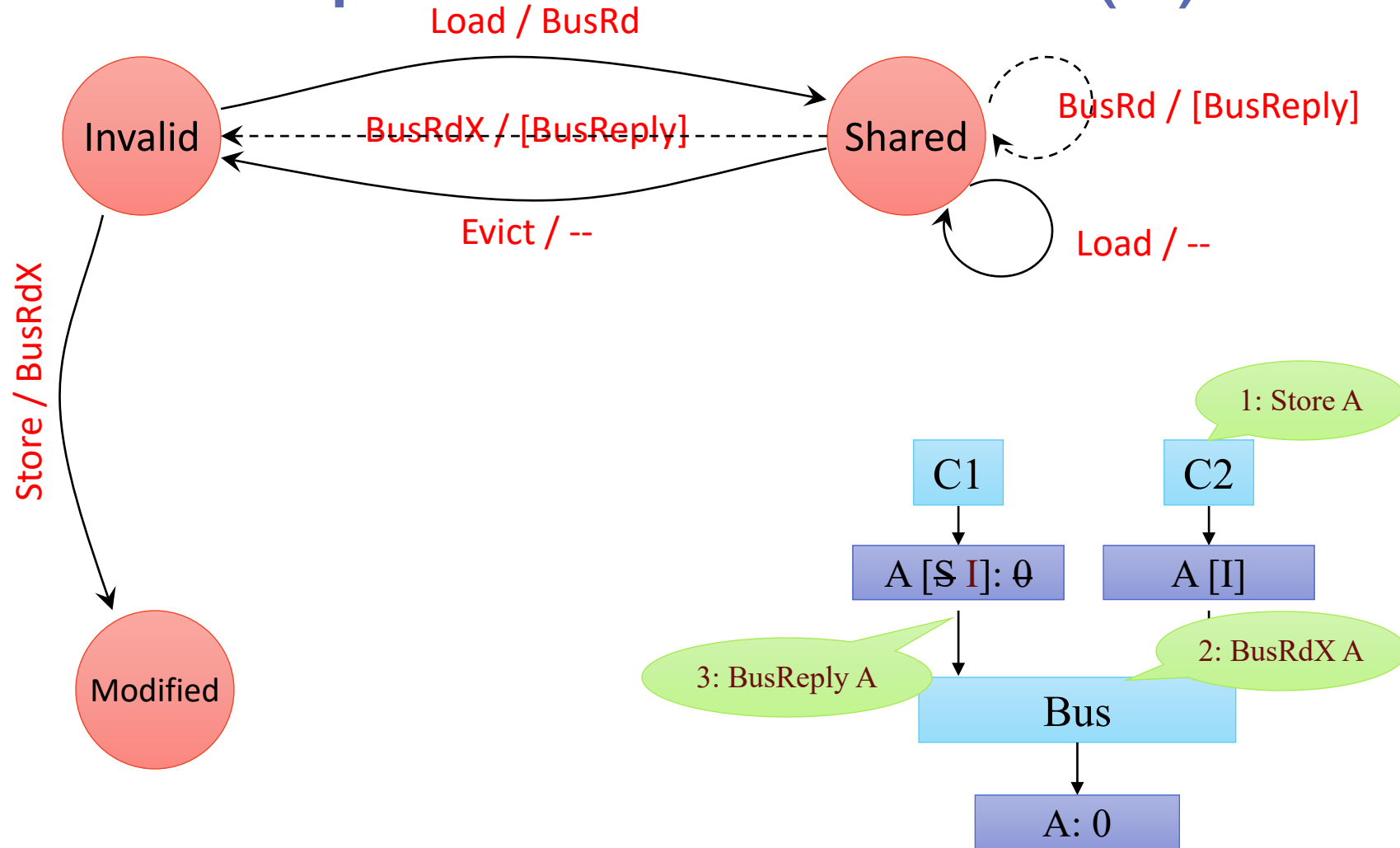
# Ex: MSI protocol in action (7)



Load / BusRd

BusInv, BusRdX / [BusReply]

BusRd / [BusReply]

Evict / --

Load / --

Store / BusRdX

BusRd / BusReply

Store / BusInv

Invalid

Shared

Modified

Load, Store / --

1: Store A

C1

C2

A [M]: 2

A [I]

Bus

A: 1

**Invariant: Only one writer (M) per block; Many sharers (S) okay**

# Ex: MSI protocol in action (7)

Load / BusRd

Invalid

BusInv, BusRdX / [BusReply]

Shared

BusRd / [BusReply]

Evict / --

Load / --

Store / BusRdX

BusRd / BusReply

Store / BusInv

Modified

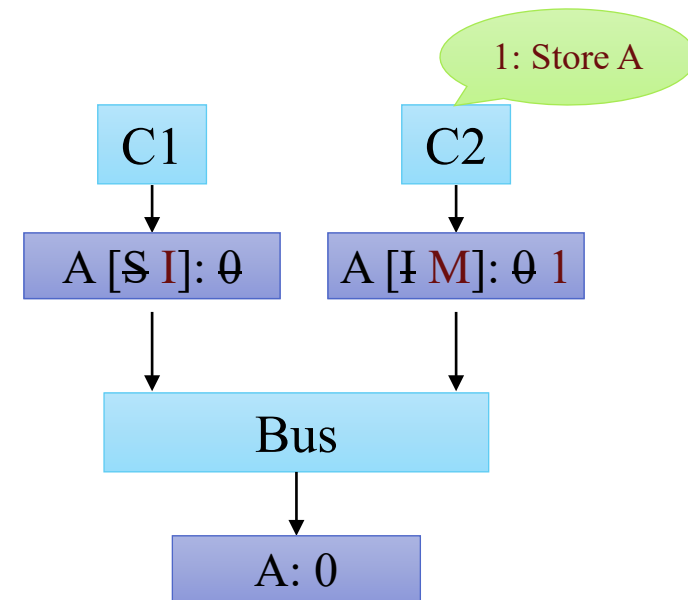Load, Store / --

C1

A [M]: 2

C2

A [I]

1: Store A

2: BusRdX A

Bus

A: 1

**Invariant: Only one writer (M) per block; Many sharers (S) okay**

# Ex: MSI protocol in action (7)



**Invariant: Only one writer (M) per block; Many sharers (S) okay**

# Ex: MSI protocol in action (7)

Load / BusRd

Invalid

BusInv, BusRdX / [BusReply]

Shared

BusRd / [BusReply]

Evict / --

Load / --

Store / BusRdX

BusRdX / BusReply

BusRd / BusReply

Store / BusInv

Modified

Load, Store / --

1: Store A

C1

C2

A [M I]: 2
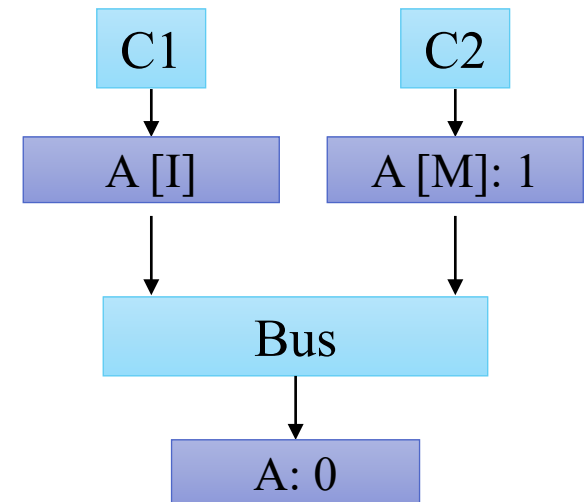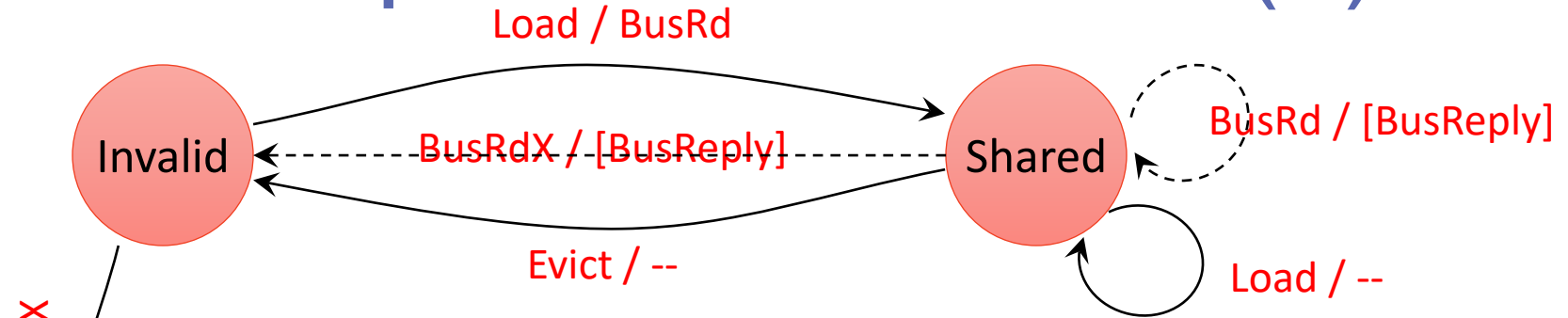
A [I]

3: BusReply A

2: BusRdX A

Bus

A: 1

**Invariant: Only one writer (M) per block; Many sharers (S) okay**

# Ex: MSI protocol in action (7)



**Invariant: Only one writer (M) per block; Many sharers (S) okay**
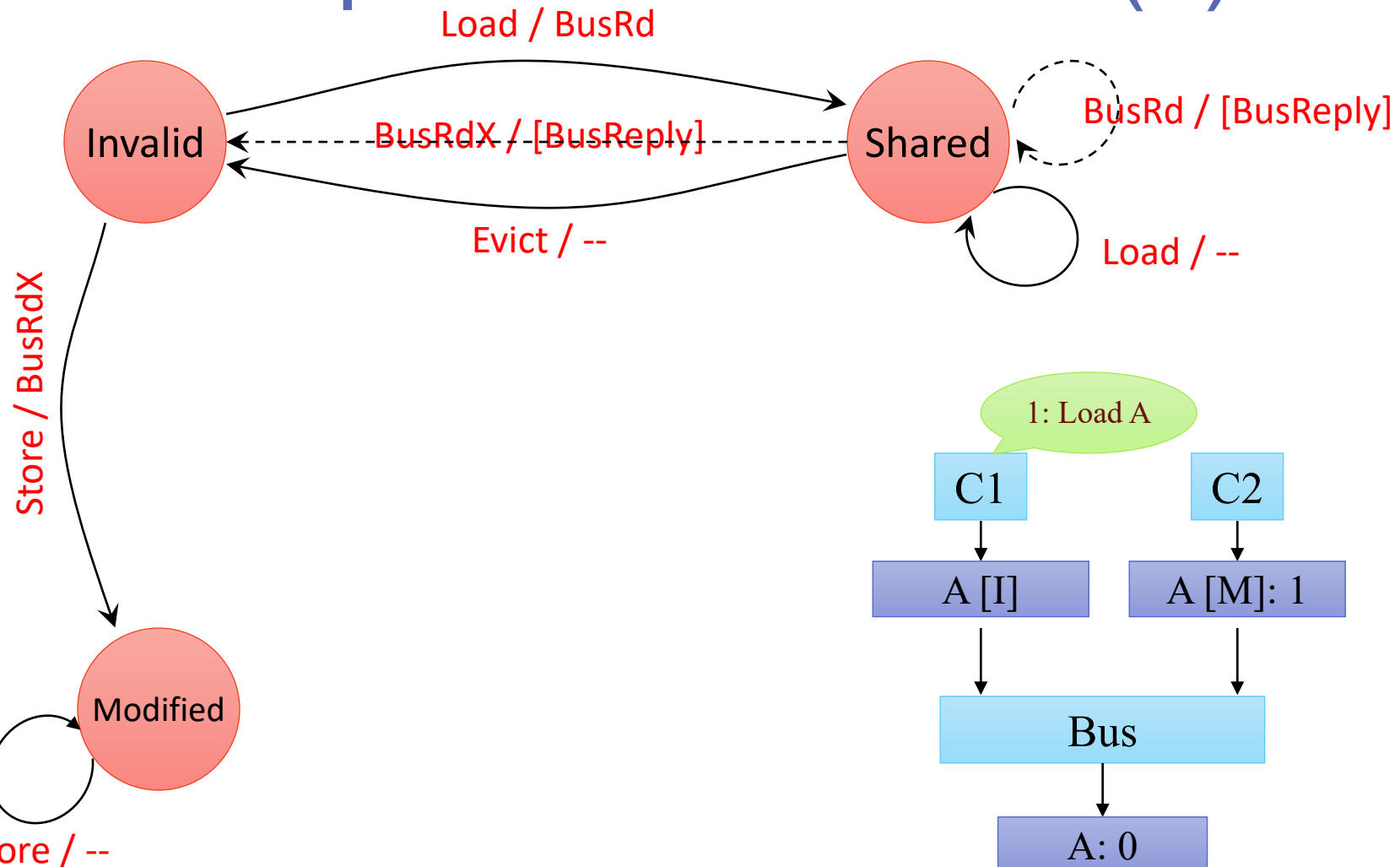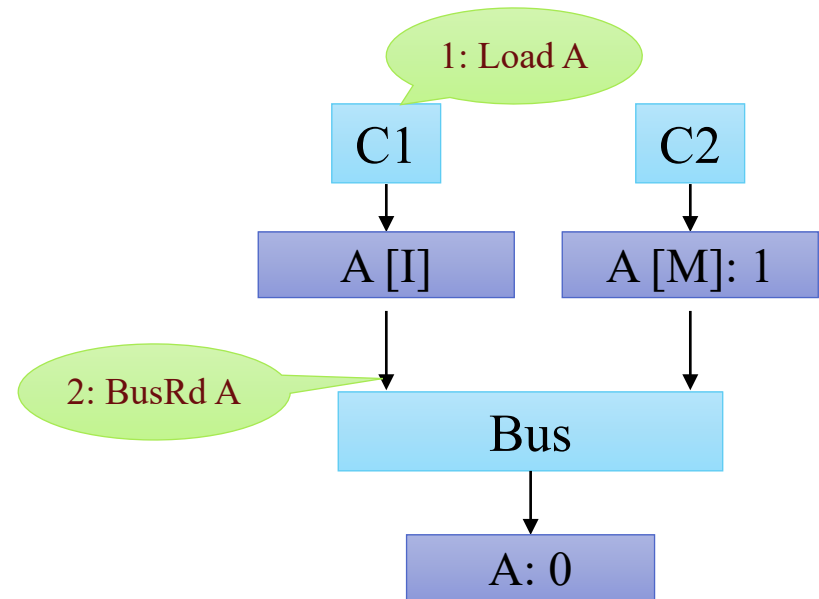
# Ex: MSI protocol in action (8)



**Invariant: Only one writer (M) per block; Many sharers (S) okay**

# Ex: MSI protocol in action (8)



**Invariant: Only one writer (M) per block; Many sharers (S) okay**

# Ex: MSI protocol in action (8)
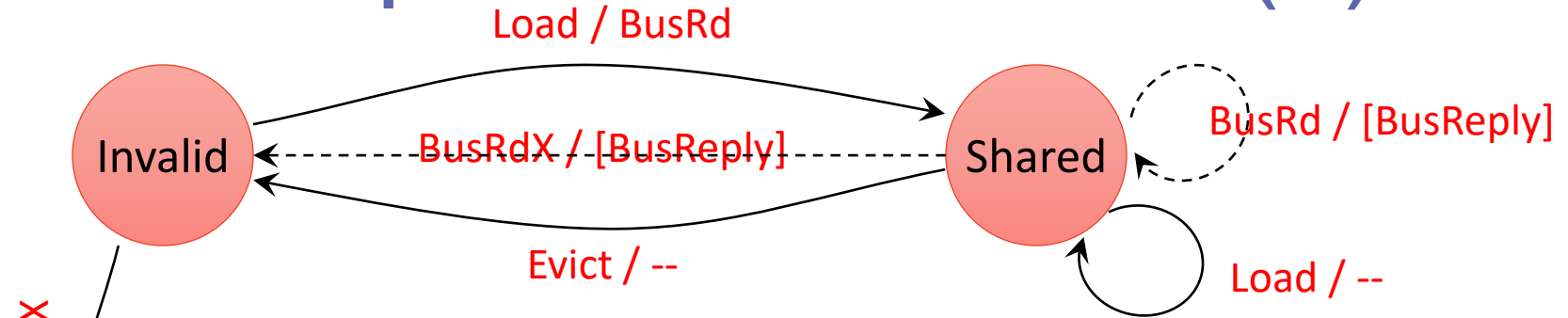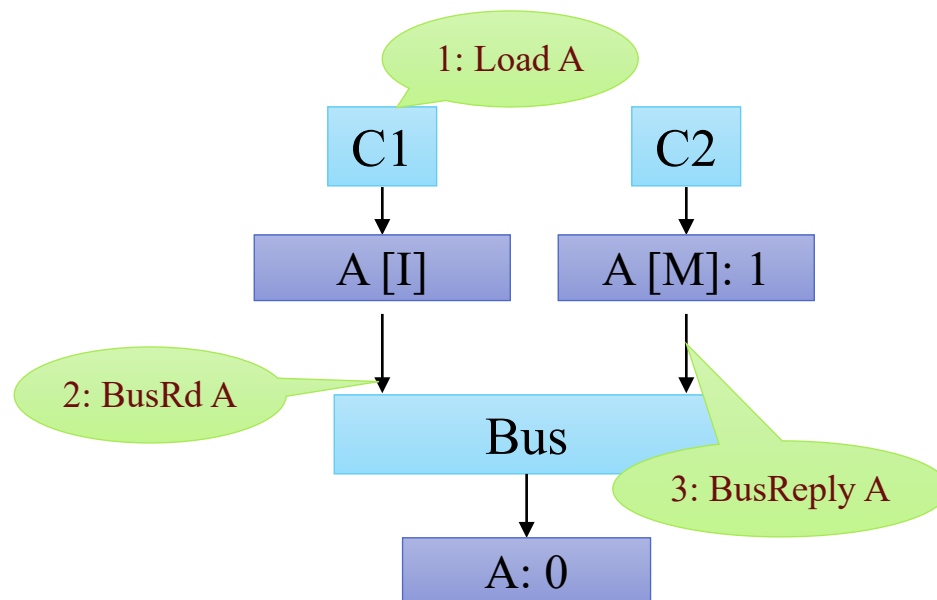


**Invariant: Only one writer (M) per block; Many sharers (S) okay**

# Ex: MSI protocol in action (8)



**Invariant: Only one writer (M) per block; Many sharers (S) okay**

# Ex: MSI protocol in action (10)



Load / BusRd

BusInv, BusRdX / [BusReply]

BusRd / [BusReply]

Evict / --

Load / --

Store / BusRdX

Evict / BusWB

BusRdX / BusReply

BusRd / BusReply

Store / BusInv

Load, Store / --

Invalid

Shared

Modified

- Cache cctions:
  - ► Load, Store, Evict

- Bus transaction /messages:
  - ► BusRd, BusRdX, BusInv, BusWB, BusReply

Invariant: Only one writer (M) per block; Many sharers (S) okay
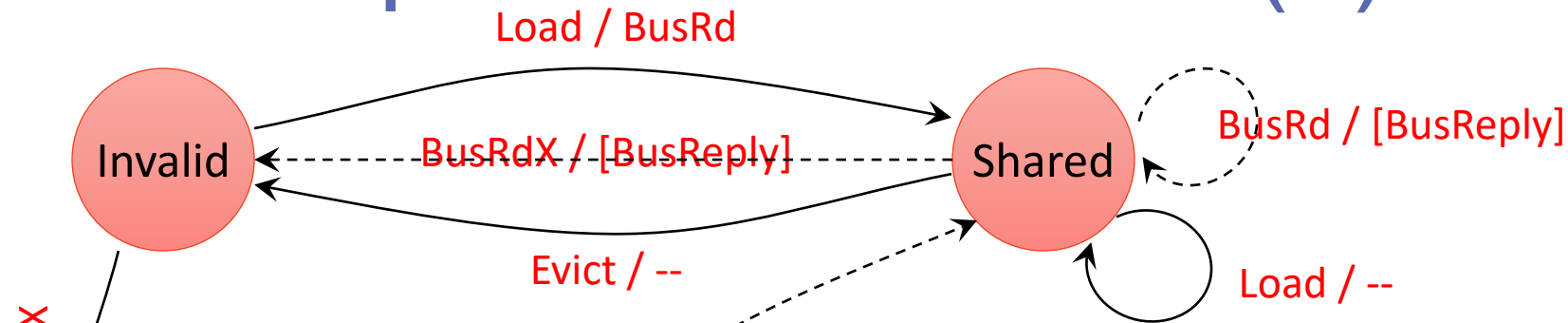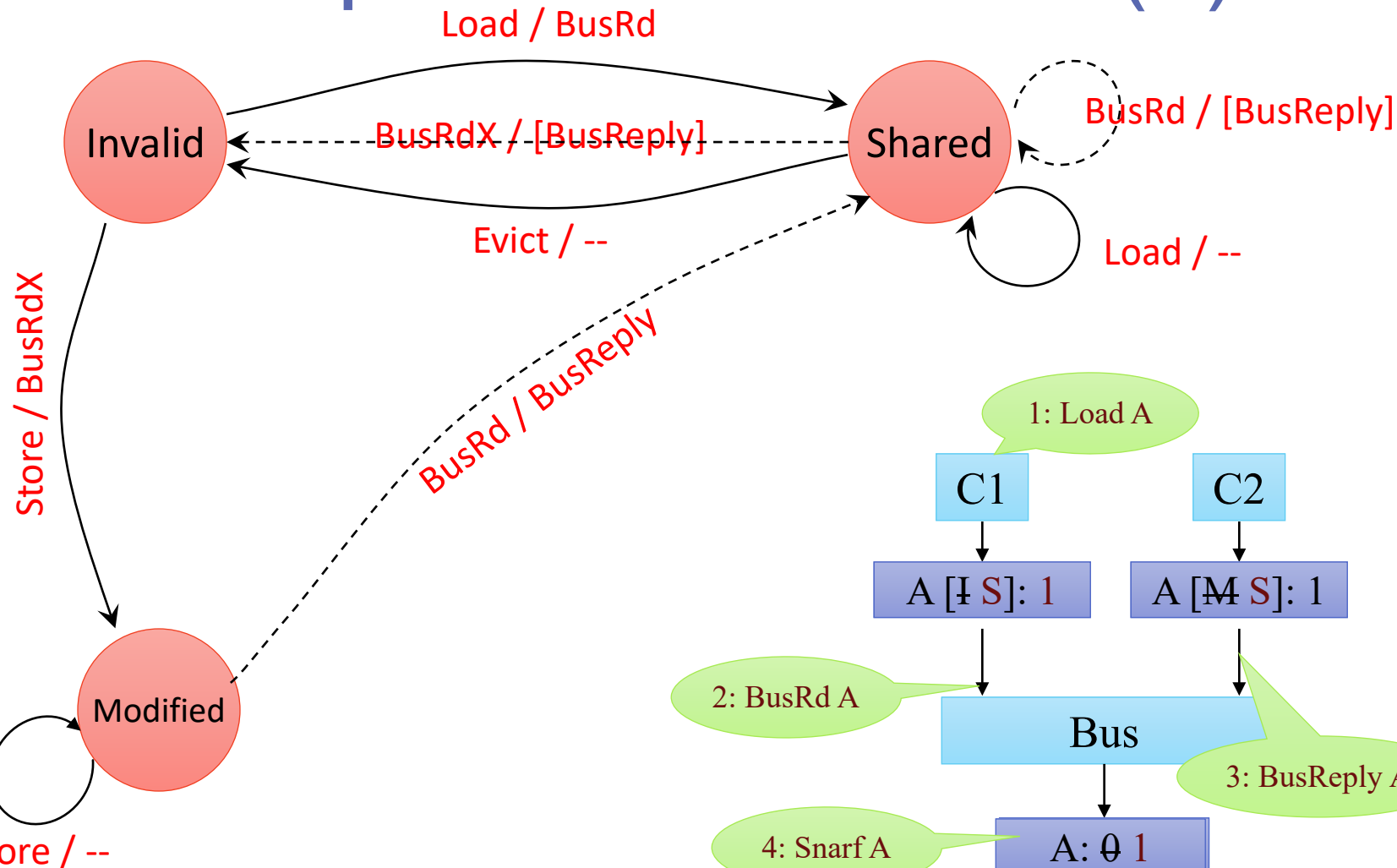
# Ex: MSI protocol in action (10)



- Cache cctions:
  - ▸ Load, Store, Evict

- Bus transaction /messages:
  - ▸ BusRd, BusRdX, BusInv, BusWB, BusReply

CCs in every private cache implements the same FSM for each $ block

Invariant: Only one writer (M) per block; Many sharers (S) okay

# Putting it together: Coherence



Example: MSI coherence protocol

# Putting it together: Coherence



Example: MSI coherence protocol

# Putting it together: Coherence



Example: MSI coherence protocol

# Putting it together: Coherence



**$ A (if present, "I" state otherwise)**

What is the state?

BusTxn A

Example: MSI coherence protocol

# Putting it together: Coherence

**$ A (if present, "I" state otherwise)**

| | Tag | Data |

**What is the state?**

Core — cc — $

Core — cc — $

Core — cc — $

Core — cc — $

**BusTxn A**

LLC Bank 0

LLC Bank 1

LLC Bank 2

LLC Bank 3

Memory controller

**Invalid** — Load / BusRd — **Shared** — BusRd / [BusReply]

Inv, BusRdX / [BusReply]

Load / --

Evict / --

Store / BusRdX

Evict / BusWB

Flush / BusReply

BusRd / BusReply

Store / BusRd

**Modified**

Load, Store / --

**Change state of "A", as per FSM**

Example: **MSI** coherence protocol

# Putting it together: Coherence



Example: MSI coherence protocol

# Putting it together: Coherence



Example: MSI coherence protocol

# Putting it together: Coherence



Example: MSI coherence protocol

# Putting it together: Coherence



**$ A (if present, "I" state otherwise)**

**What is the state?**

Example: MSI coherence protocol

# Putting it together: Coherence



Example: MSI coherence protocol

# Putting it together: Coherence

**$ A (if present, "I" state otherwise)**

| | Tag | Data |
|---|---|---|

**What is the state?**

**ld/st A**

**BusTxn A**

Load / BusRd

Invalid ← Inv, BusRdX / [BusReply] → Shared · BusRd / [BusReply]

Evict / --

Load / --

Store / BusRdX

Evict / BusWB

BusRd / BusReply

Store / BusInv

Store / BusInv

Modified

Load, Store / --

**Send Bus txn. for "A", as per FSM**

| Core | Core | Core | Core |
|---|---|---|---|
| cc $ | cc $ | cc $ | cc $ |

| LLC Bank 0 | LLC Bank 1 | LLC Bank 2 | LLC Bank 3 |
|---|---|---|---|

Memory controller

## Example: MSI coherence protocol

# Many possible protocols

- MSI protocol is <u>not</u> the only possibility
  - ▸ It's a basic protocol

- M<span style="color:red">E</span>SI protocol optimizes for access pattern where the same data is read and then immediately written by a thread

# FSM for MESI protocol



**Load / BusRd**
**(if someone else has it)**

BusInv, BusRdX / [BusReply]

BusRd / [BusReply]

Evict / --

Load / --

**Load / BusRd**
**(if no one else has it)**

Store / BusInv

Store / BusRdX

Evict / BusWB

BusRdX / BusReply

BusRd / BusReply

BusRdX / [BusReply]

BusRd / BusReply

Evict / --

Store / --

Load / --

Load, Store / --

**Invalid**

**Shared**

**Modified**

**Exclusive**

"E" state: Similar to "M" but not written to yet (clean). On a Read, "I" –>"E", if no other private cache have the block

# Many possible protocols

- MSI protocol is <u>not</u> the only possibility
  - ▶ It's a basic protocol

- MESI protocol optimizes for access pattern where the same data is read and then immediately written by a thread

- MOESI further optimizes for producer-consumer access pattern in applications

# Take away

- Multicores: Multiple replicated cores, along with own private caches that typically shares a larger LLC

- Cache coherence protocols keeps the private caches coherent

- Many possible ways to implement coherence

# Directory coherence protocols

Scalable cache coherence

# Problems w/ Snoopy Coherence

1) Interconnect bandwidth
   ▸ Problem: Bus and Ring are not scalable interconnects
     • Limited bandwidth
     • Cannot support more than a dozen or so processors
   ▸ Solution: Replace non-scalable interconnect (ring or bus) with a scalable one (e.g., mesh)

# Problems w/ Snoopy Coherence

1) Interconnect bandwidth
   - ▸ Problem: Bus and Ring are not scalable interconnects
     - Limited bandwidth
     - Cannot support more than a dozen or so processors
   - ▸ Solution: Replace non-scalable interconnect (ring or bus) with a scalable one (e.g., mesh)

2) Cache snooping bandwidth
   - ▸ Problem: All caches must monitor all bus traffic; most snoops result in no action
   - ▸ Solution: Replace non-scalable broadcast protocol (spam everyone) with scalable directory protocol (notify cores that care)
     - The "directory" keeps track of "sharers"

# Directory Coherence Protocols

- Typically use point-to-point scalable networks
  - ▸ Such as Crossbar or Mesh

# Issues with point-to-point networks



Advantages:
- better electrical behavior (shorter wires)
- coherence transaction parallelism

Problem: unordered network
Nodes may observe messages in different orders
Is this a problem?
(May break the write propagation constraint)

# Ordering – What is wrong here?

# Ordering – What is wrong here?



What caused the problem here?

# Ordering – What is wrong here?



What caused the problem here?

How do we fix it?

# Ordering – What is wrong here?



**What caused the problem here?**

**How do we fix it?**

**Need a single point of ordering for all coherence transactions to a given address → directory**

# Directory Coherence Protocols

- Each physical cache line has a ***home node/core/contoller***

- Extend memory controller (or LLC bank) to track caching information for cache lines for which it is home
  - ▸ Information kept in a hardware structure called ***Directory***

# Directory Coherence Protocols

- Each physical cache line has a ***home node/core/contoller***

- Extend memory controller (or LLC bank) to track caching information for cache lines for which it is home
  - Information kept in a hardware structure called ***Directory***

- For each physical cache line, a ***home directory*** tracks:
  - ***Owner***: core that has a dirty copy (i.e., M state)
  - ***Sharers***: cores that have clean copies (i.e., S state)

# Directory Coherence Protocols

- Each physical cache line has a ***home node/core/contoller***

- Extend memory controller (or LLC bank) to track caching information for cache lines for which it is home
  - ▸ Information kept in a hardware structure called ***Directory***

- For each physical cache line, a ***home directory*** tracks:
  - ▸ ***Owner***: core that has a dirty copy (i.e., M state)
  - ▸ ***Sharers***: cores that have clean copies (i.e., S state)

- Cores send coherence requests to home directory

- Home directory forwards messages only to cores that "care" (i.e., cores that might have a copy of the line)

# MSI Directory Example: Step #1



C₀   C₁   C₂

Load A

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| -- | -- | -- |
| -- | -- | -- |

Miss!

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | 500 | M |
| -- | -- | -- |

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| -- | -- | -- |
| -- | -- | -- |

## Point-to-Point Interconnect

**Directory at LLC**

| Addr | Data | State | Sharers |
|------|------|-------|---------|
| A | 1000 | Modified | C1 |
| B | 0 | Idle | -- |

62

# MSI Directory Example: Step #2



63

# MSI Directory Example: Step #2



63

# MSI Directory Example: Step #3

**C₀** **C₁** **C₂**

Load A

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| -- | -- | -- |
| -- | -- | -- |

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | 500 | S |
| -- | -- | -- |

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| -- | -- | -- |
| -- | -- | -- |

DataResp: Addr=A, Data=500

## Point-to-Point Interconnect

**Directory at LLC**

| Addr | Data | State | Sharers |
|------|------|-------|---------|
| A | 1000 | Blocked | C1 |
| B | 0 | Idle | -- |

64

# MSI Directory Example: Step #4



**C_0**

Load A

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | 500 | S |
| -- | -- | -- |

**C_1**

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | 500 | S |
| -- | -- | -- |

**C_2**

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| -- | -- | -- |
| -- | -- | -- |

DataResp: Addr=A, Data=500

## Point-to-Point Interconnect

**Directory at LLC**

| Addr | Data | State | Sharers |
|------|------|-------|---------|
| A | 1000 | Blocked | C1 |
| B | 0 | Idle | -- |

65

# MSI Directory Example: Step #5

# MSI Directory Example: Step #6

C<sub>0</sub>

C<sub>1</sub>

C<sub>2</sub>

Store 400 -> A

Miss!

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | 500 | S |
| -- | -- | -- |

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | 500 | S |
| -- | -- | -- |

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| -- | -- | -- |
| -- | -- | -- |

## Point-to-Point Interconnect

**LLC & Directory**

| Addr | Data | State | Sharers |
|------|------|-------|---------|
| A | 500 | Shared, Dirty | C0, C1 |
| B | 0 | Idle | -- |

# MSI Directory Example: Step #7



Store 400 -> A

**L1 Cache** (C0)

| Addr | Data | State |
|------|------|-------|
| A | 500 | S |
| -- | -- | -- |

**L1 Cache** (C1)

| Addr | Data | State |
|------|------|-------|
| A | 500 | S |
| -- | -- | -- |

**L1 Cache** (C2)

| Addr | Data | State |
|------|------|-------|
| -- | -- | -- |
| -- | -- | -- |

GetM: Addr=A

## Point-to-Point Interconnect

**LLC & Directory**

| Addr | Data | State | Sharers |
|------|------|-------|---------|
| A | 500 | Blocked | C0, C1 |
| B | 0 | Idle | -- |

68

# MSI Directory Example: Step #8



C₀  C₁  C₂

Store 400 -> A

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | 500 | S |
| -- | -- | -- |

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | 500 | S |
| -- | -- | -- |

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | 500 | Blk |
| -- | -- | -- |

DataResp: Addr=A, 500, Ack=2

Point-to-Point Interconnect

**LLC & Directory**

| Addr | Data | State | Sharers |
|------|------|-------|---------|
| A | 500 | Blocked | C0, C1 |
| B | 0 | Idle | -- |

# MSI Directory Example: Step #9

# MSI Directory Example: Step #10

C_0

C_1

C_2

Store 400 -> A

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | -- | I |
| -- | -- | -- |

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | -- | I |
| -- | -- | -- |

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | 500 | Blk |
| -- | -- | -- |

## Point-to-Point Interconnect

**Directory & LLC**

| Addr | Data | State | Sharers |
|------|------|-------|---------|
| A | 500 | Blocked | C0, C1 |
| B | 0 | Idle | -- |

# MSI Directory Example: Step #10



C₀ gives C_0:

**L1 Cache** (C0)

| Addr | Data | State |
|------|------|-------|
| A | -- | I |
| -- | -- | -- |

**L1 Cache** (C1)

| Addr | Data | State |
|------|------|-------|
| A | -- | I |
| -- | -- | -- |

**L1 Cache** (C2)

Store 400 -> A

| Addr | Data | State |
|------|------|-------|
| A | 500 | Blk |
| -- | -- | -- |

Ack: Addr=A

**Point-to-Point Interconnect**

**Directory & LLC**

| Addr | Data | State | Sharers |
|------|------|-------|---------|
| A | 500 | Blocked | C0, C1 |
| B | 0 | Idle | -- |

# MSI Directory Example: Step #10

C0

C1

C2

Store 400 -> A

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | -- | I |
| -- | -- | -- |

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | -- | I |
| -- | -- | -- |

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | 500 | Blk |
| -- | -- | -- |

Ack: Addr=A,

Ack: Addr=A

Point-to-Point Interconnect

**Directory & LLC**

| Addr | Data | State | Sharers |
|------|------|-------|---------|
| A | 500 | Blocked | C0, C1 |
| B | 0 | Idle | -- |

# MSI Directory Example: Step #11



C_0     C_1     C_2

Store 400 -> A

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | -- | I |
| -- | -- | -- |

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | -- | I |
| -- | -- | -- |

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | 500 | M |
| -- | -- | -- |

Point-to-Point Interconnect

**Directory & LLC**

| Addr | Data | State | Sharers |
|------|------|-------|---------|
| A | 500 | Blocked | C0, C1 |
| B | 0 | Idle | -- |

# MSI Directory Example: Step #12

C0

C1

C2

Store 400 -> A

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | -- | I |
| -- | -- | -- |

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | -- | I |
| -- | -- | -- |

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | 400 | M |
| -- | -- | -- |

## Point-to-Point Interconnect

**Directory & LLC**

| Addr | Data | State | Sharers |
|------|------|-------|---------|
| A | 500 | Blocked | C0, C1 |
| B | 0 | Idle | -- |

# MSI Directory Example: Step #12

C$_0$

C$_1$

C$_2$

Store 400 -> A

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | -- | I |
| -- | -- | -- |

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | -- | I |
| -- | -- | -- |

**L1 Cache**

| Addr | Data | State |
|------|------|-------|
| A | 400 | M |
| -- | -- | -- |

Unblock: Addr=A

## Point-to-Point Interconnect

**Directory & LLC**

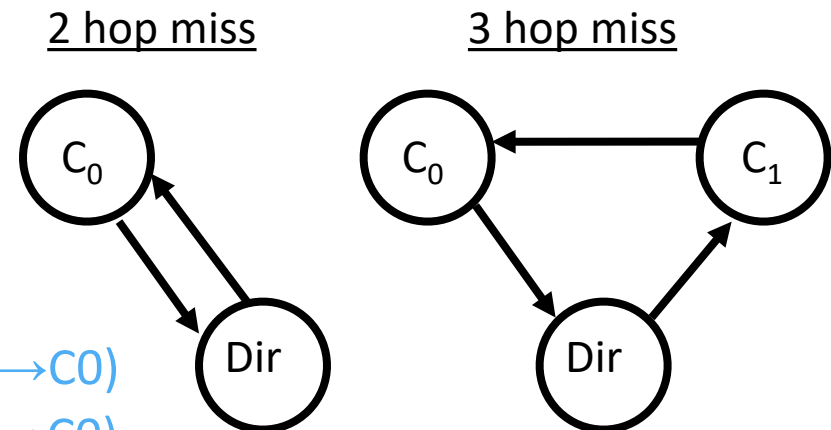| Addr | Data | State | Sharers |
|------|------|-------|---------|
| A | 500 | Blocked | C0, C1 |
| B | 0 | Idle | -- |

# MSI Directory Example: Step #13

# Directory Flip Side: Latency

- Directory protocols
  - + Lower bandwidth consumption → more scalable
  - ▶ Longer latencies

- Two read miss situations

- Unshared: get data from memory
  - ▶ Snooping: 2 hops (C0→LLC/memory→C0)
  - ▶ Directory: 2 hops (C0→LLC/memory→C0)

- Shared or exclusive: get data from other processor (P1)
  - ▶ Assume cache-to-cache transfer optimization
  - ▶ Snooping: 2 hops (C0→C1→C0)
  - ▶ Directory: **3 hops** (C0→Dir/LLC→C1→C0)
  - ▶ Common, with many processors high probability someone has it

2 hop miss          3 hop miss

# Directory Flip Side: Complexity

- Latency not the only issue for directories
  - ▸ Subtle correctness issues as well
  - ▸ Stem from unordered nature of underlying inter-connect
- Individual requests to single cache must be ordered
  - ▸ Point-to-point network: requests may arrive in different orders
    - Directory has to enforce ordering explicitly
    - Cannot initiate actions on request B…
    - Until all relevant processors have completed actions on request A
    - Requires directory to collect acks, queue requests, etc.

- Directory protocols
  - ▸ Obvious in principle
  - ▸ Complicated in practice
  - ▸ **State space explosion due to unordered network**
  - ▸ **Need to consider various possible coherence message races**