# Multiprocessors and Coherent Memory

Erik Hagersten

Uppsala University

# Goal for this course

- Understand **how and why** modern computer systems are designed the way the are:
  - pipelines
  - ✓ memory organization
  - ✓ virtual/physical memory …

- Understand **how and why** multiprocessors are built
  - Cache coherence
  - Memory models
  - Synchronization…

  *This batch of lectures*

- Understand **how and why** parallelism is created
  - Instruction-level parallelism
  - √ Memory-level parallelism
  - Thread-level parallelism…

- Understand **how and why** multiprocessors of combined SIMD/MIMD type are built
  - GPU
  - Vector processing…

- Understand **how** computer systems are adopted to different usage areas
  - General-purpose processors
  - Embedded/network processors…

- Understand the physical limitation of modern computers
  - ✓ Bandwidth
  - Energy
  - Cooling…

UPPSALA
UNIVERSITET

AVDARK
2013

# The era of the "Rocket Science Supercomputers" 1980-1995

- The one with the most blinking lights wins
- The one with the niftiest language wins
- The more different the better!

UPPSALA
UNIVERSITET

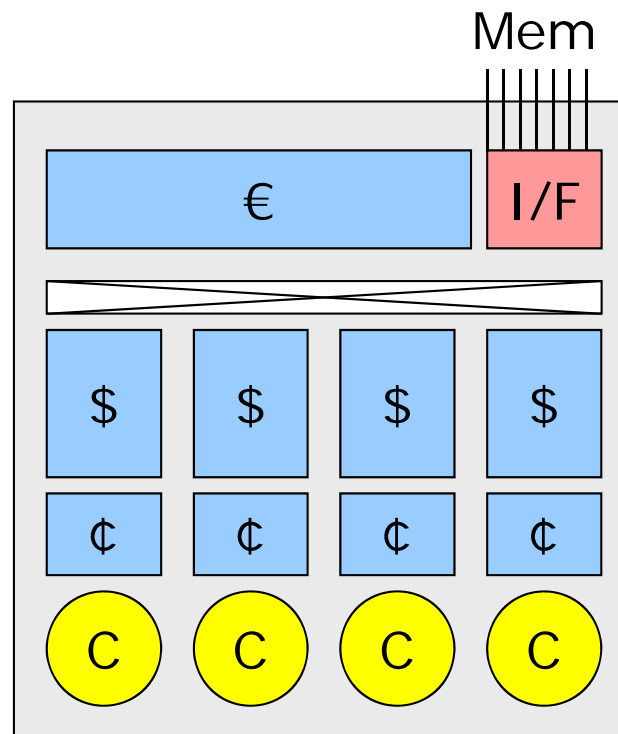AVDARK
2013

# The server market 1995

| Server Size | High-Perf. Computing | Commercial Computing |
|---|---|---|
| <$10k | 1% | 19% |
| <$50k | 5% | 25% |
| <$250k | 5% | 24% |
| <$1M | 2% | 9% |
| >$1M | 3% | 8% |

UNIX shared-mem servers

The target of the rocket science supercomputers

AVDARK
2013

UPPSALA
UNIVERSITET

# Multicore: Who has <u>not</u> got one?

# *MP Taxonomy (more later…)*

```
                        ●
                       / \
                      /   \
                     /     \
                  SIMD     MIMD
                            / \
                           /   \
                          /     \
                  Message-      Shared
                  passing       Memory
                   / \           /|\
                  /   \         / | \
                 /     \       /  |  \
              Fine-  Coarse-  UMA NUMA COMA
            grained  grained
```

UPPSALA UNIVERSITET

AVDARK 2013

# Models of parallelism

- Processes (**fork** or **&** in UNIX)
  - A parallel execution, where each process has its own process state, e.g., its own VA→PA mapping
- Threads (**thread_create** in POSIX)
  - Parallel threads of control inside a process
  - There are some thread-shared state, e.g., VA→PA mapping.
- More: OpenMP, OpenACC, OpenCL, SILC, ...
- Common property: Each thread has its own PC (i.e. executes its own code independently)
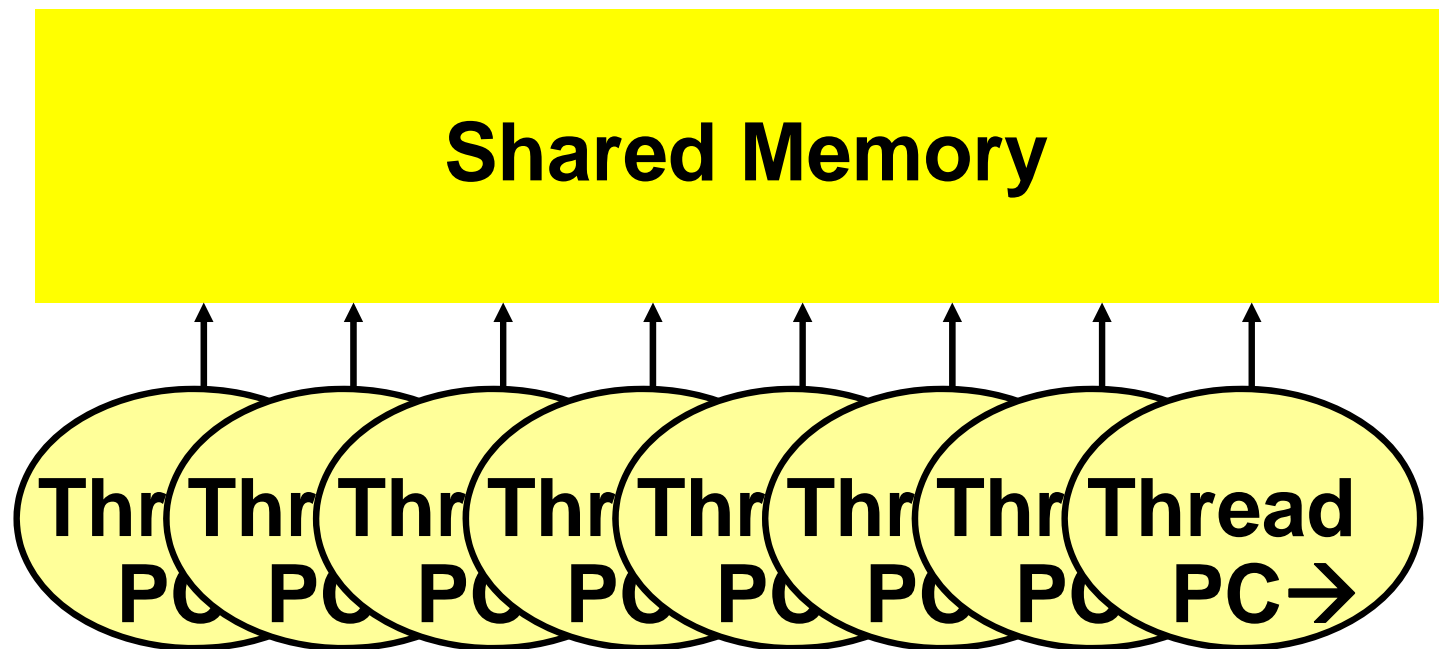- More during lab lecture...

AVDARK
2013

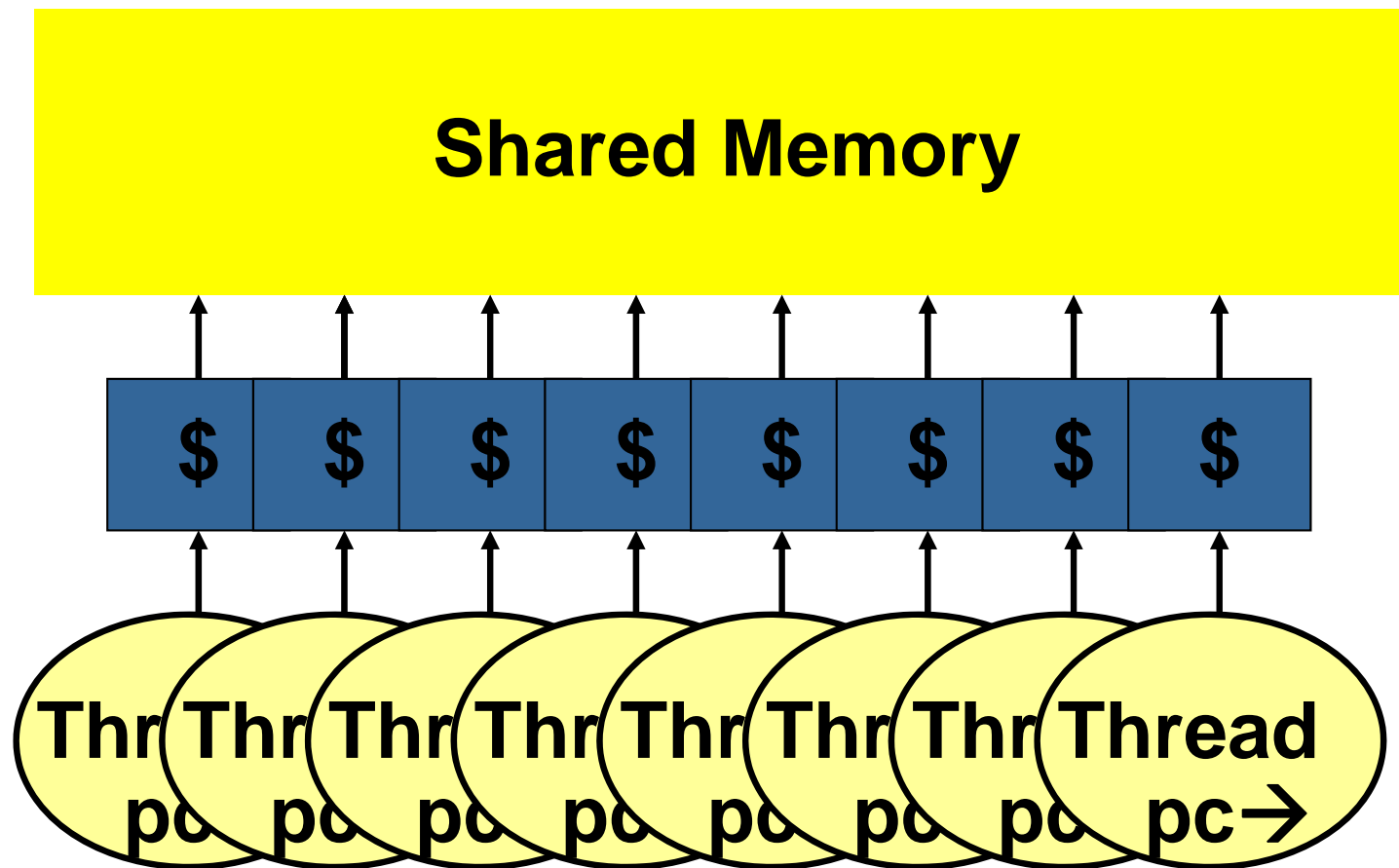# What is:
# Coherent Shared Memory?

Erik Hagersten

Uppsala University

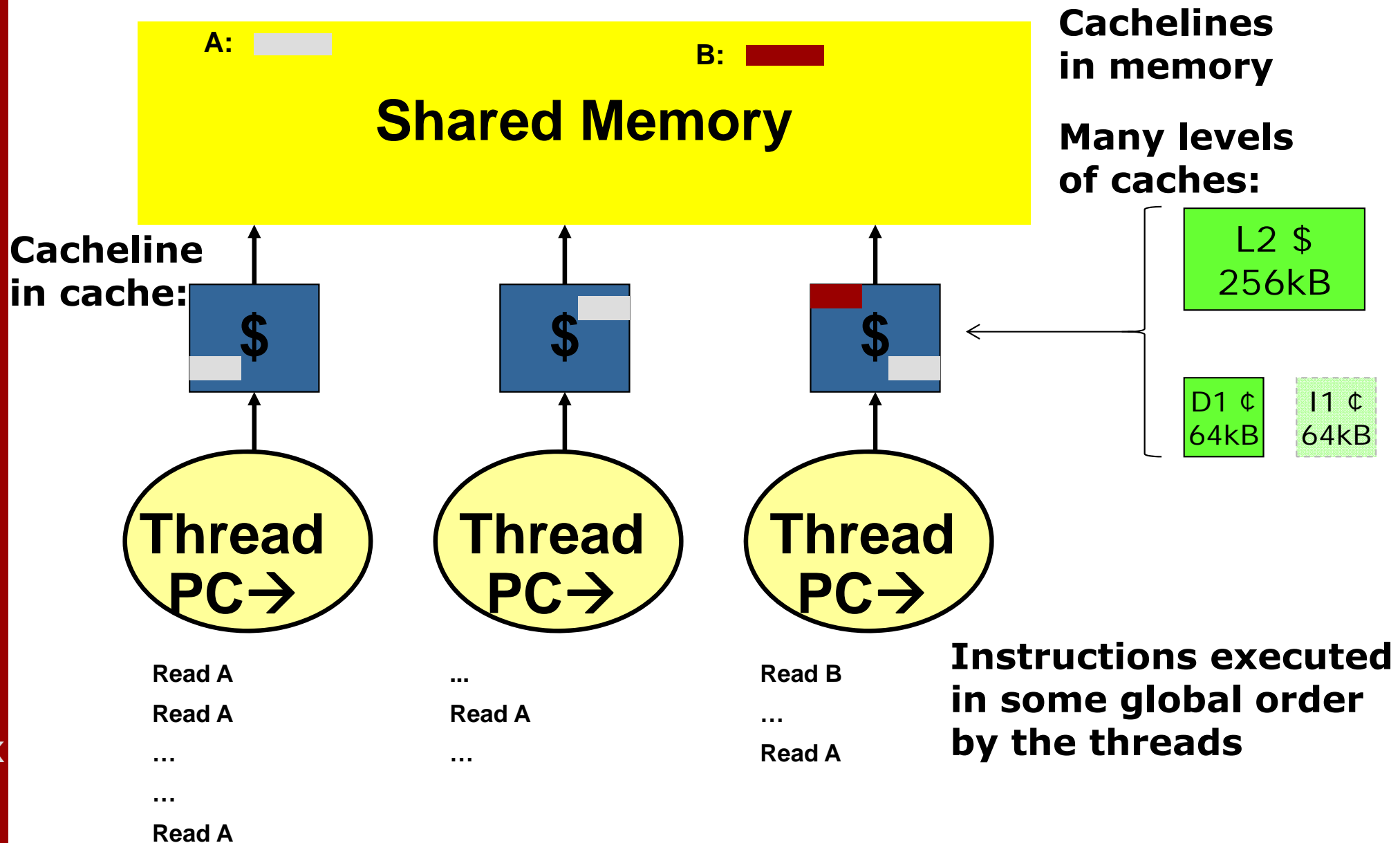# Programming Model: Coherent shared memory



**Shared Memory**

Thr Thr Thr Thr Thr Thr Thr Thread
PC PC PC PC PC PC PC PC→

UPPSALA UNIVERSITET

AVDARK 2013

# Adding Caches. Gives the illusion of:
- **Shorter Memory Latency**
- **Higher Memory Bandwidth**

# Our Generic Shared Memory Arch.

A:

B:

**Shared Memory**

**Cachelines in memory**

**Many levels of caches:**

L2 $
256kB

D1 ¢
64kB

I1 ¢
64kB

**Cacheline in cache:**

$

$

$

**Thread PC→**

**Thread PC→**

**Thread PC→**

Read A

Read A

...

...

Read A

...

Read A

...

Read B

...

Read A

**Instructions executed in some global order by the threads**

AVDARK
2013

UPPSALA
UNIVERSITET

# Automatic Replication of Data

A:

B:

**Shared Memory**

$

$

$

**Thread**

**Thread**

**Thread**

Read A

Read A

...

...

...

Read A

...

Read B

...

Read A

UPPSALA
UNIVERSITET

AVDARK
2013

# The Cache Coherent Memory System
## Coherent Write (Here: Write invalidate)

A: [ ]          B: [ ]

**Shared Memory**

$         INV          $          INV          $

**Thread**          **Thread**          **Thread**

Read A          ...          Read B
Read A          Read A          ...
...          ...          Read A
...          **Write A**

AVDARK
2013

# The Cache Coherent Memory System Coherent Read & Write-back
## (Here: Cache to Cache Transfer)

A:

B:

**Shared Memory**

$   $   $

**Thread**     **Thread**     **Thread**

Read A          ...              Read B

Read A          Read A           …

…               …                Read A

…               …                …

Read A          Write A          Replace B

UPPSALA UNIVERSITET

AVDARK 2013

# The Cache Coherent Memory System
# Coherent Read & Write-back
**(Here: Write Back)**



Thread:
Read A
Read A
…
…
Read A

Thread:
…
Read A
…
Write A
…
Replace A

Thread:
Read B
…
Read A

MP 15

AVDARK
2013

# The Cache Coherent Memory System
# Coherent Read & Write-back
## (Here: Cache-to-Cache Transfer and Write Back)

A:

B:

**Shared Memory**

$ 

$ 

$ 

**Thread**

**Thread**

**Thread**

| | | |
|---|---|---|
| Read A | ... | Read B |
| Read A | Read A | … |
| … | … | Read A |
| … | Write A | … |
| Read A | … | Replace B |
| | Replace A | |

MP 16

UPPSALA
UNIVERSITET

AVDARK
2013

# Summing up Coherence

*Good intuition, but too strong definition!*

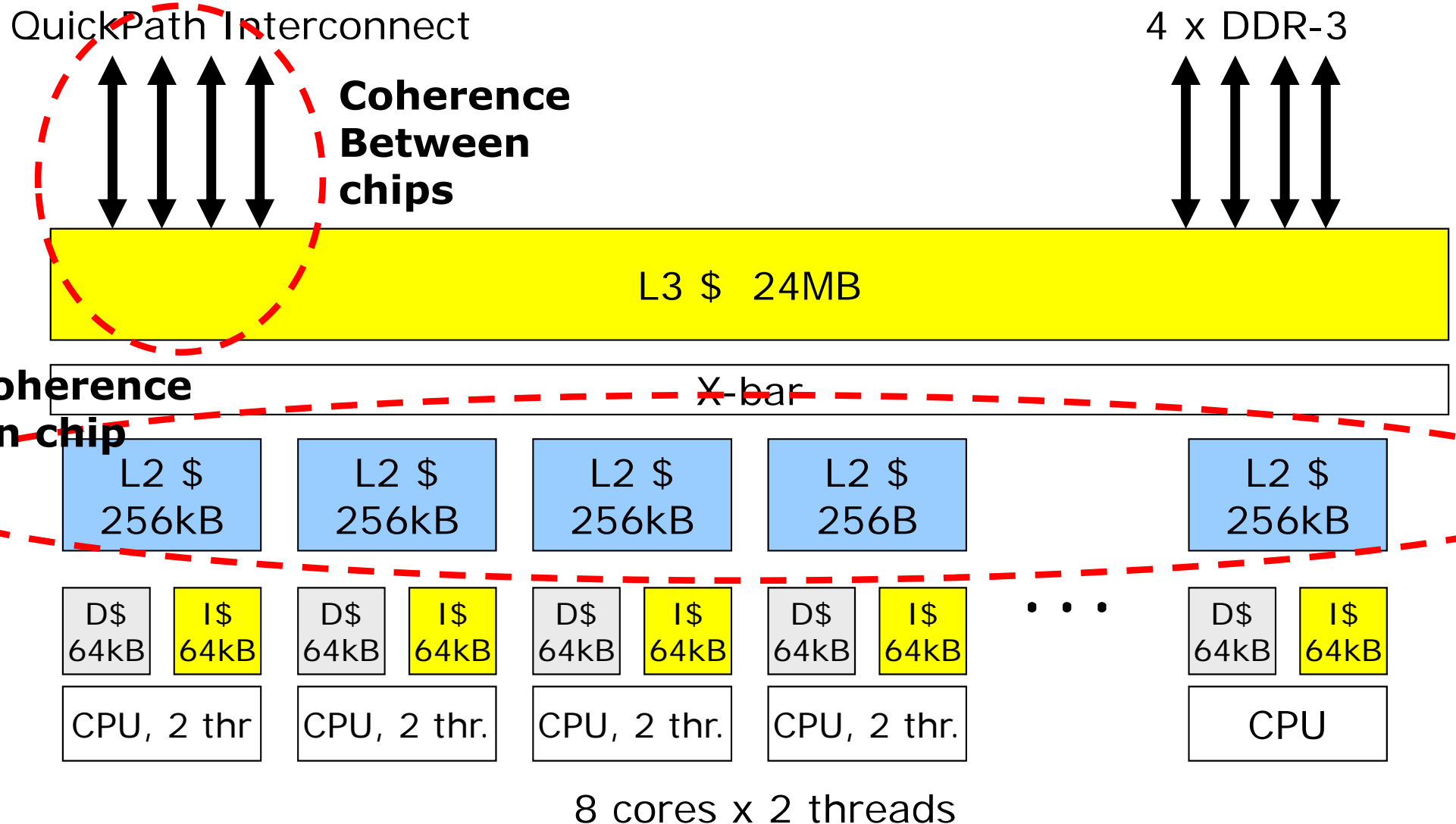*There can be many copies of a datum, but only one value*

*There is a single global order of value changes to each datum*

*After the computer stops, all copies should have the same value*

Thread1={1,2,3,4,5,6,7…}   Thread2={1,4,7…}   Thread3={1,8,7…}

Dept of Information Technology| www.it.uu.se

AVDARK
2013

# Where does coherence matter?

QuickPath Interconnect

4 x DDR-3

**Coherence Between chips**

L3 $  24MB

**Coherence On chip**

X-bar

| L2 $ 256kB | L2 $ 256kB | L2 $ 256kB | L2 $ 256B | | L2 $ 256kB |
|---|---|---|---|---|---|

| D$ 64kB | I$ 64kB | D$ 64kB | I$ 64kB | D$ 64kB | I$ 64kB | D$ 64kB | I$ 64kB | · · · | D$ 64kB | I$ 64kB |
|---|---|---|---|---|---|---|---|---|---|---|

| CPU, 2 thr | CPU, 2 thr. | CPU, 2 thr. | CPU, 2 thr. | | CPU |
|---|---|---|---|---|---|

8 cores x 2 threads

# Summary Coherence

- Coherent shared-memory programming model requires coherence

- (There is also non-coherent shared memory, e.g. single-sided MPI, PGAS)

- All threads can read and write shared data.

- Coherent view of the value of <u>a</u> datum

- Often: Coherence is kept per cache line.

UPPSALA
UNIVERSITET

AVDARK
2013

# Snooping Coherence

Erik Hagersten
Uppsala University

# Implementation options for coherence

- Two coherence options
  - Snoop-based ("broadcast protocol")
  - Directory-based ("point to point protocol")
- Different scalability
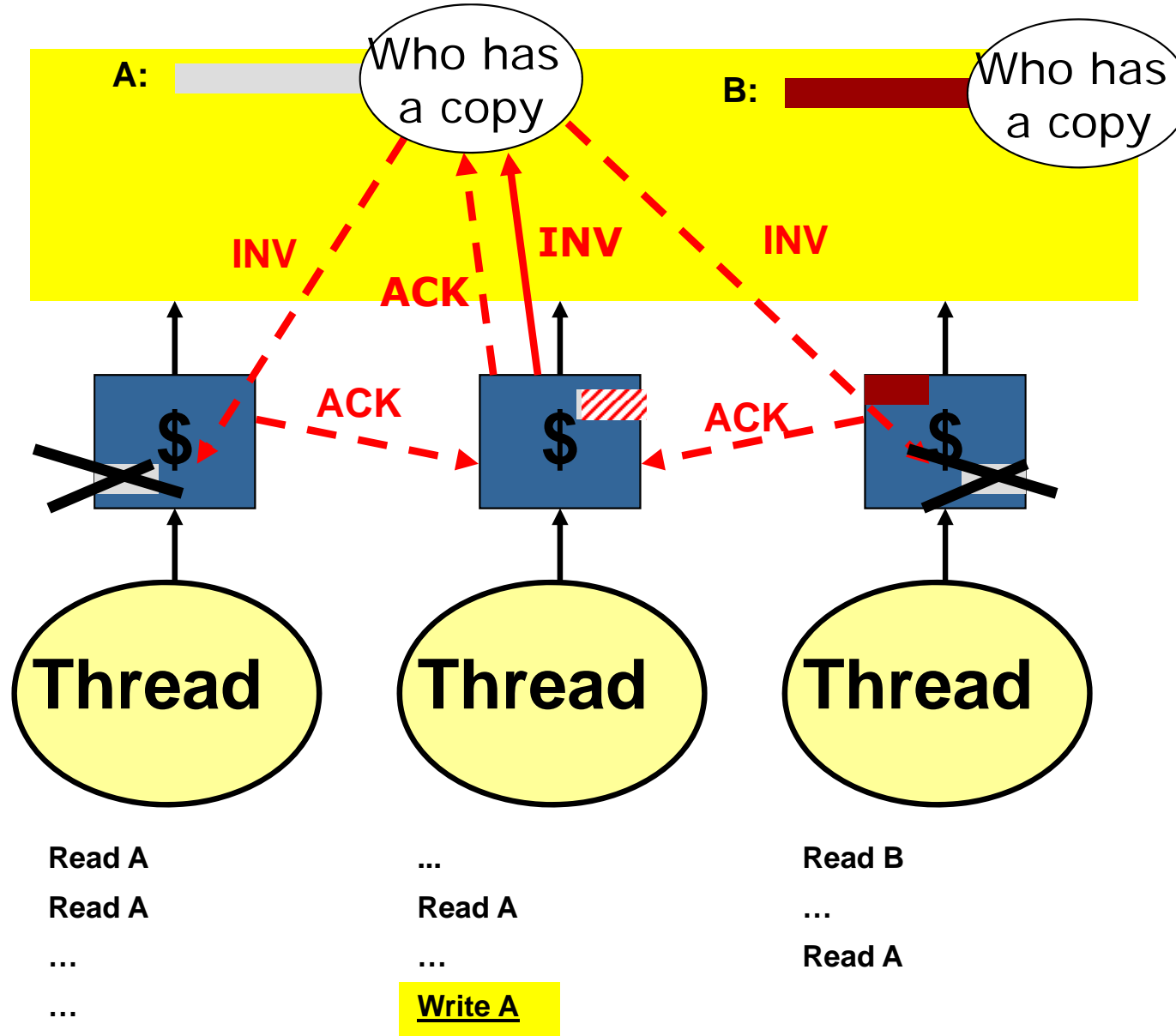- Different latency

# "Upgrade" in snoop-based

A: [ ]   B: [ ]

**BusINV**

**Have to INV**   **My INV**   **Have to INV**

$ $   $ $   $ $

Thread   Thread   Thread

| | | |
|---|---|---|
| Read A | ... | Read B |
| Read A | Read A | ... |
| ... | ... | Read A |
| ... | **Write A** | |

**How can we implement coherence on a write?**
- ❑ Invalidate A in mem.
- ❑ Invalidate A in left $
- ❑ Invalidate A in right $
- ❑ Invalidate A in both left and right $

AVDARK 2013

# "Upgrade" in dir-based
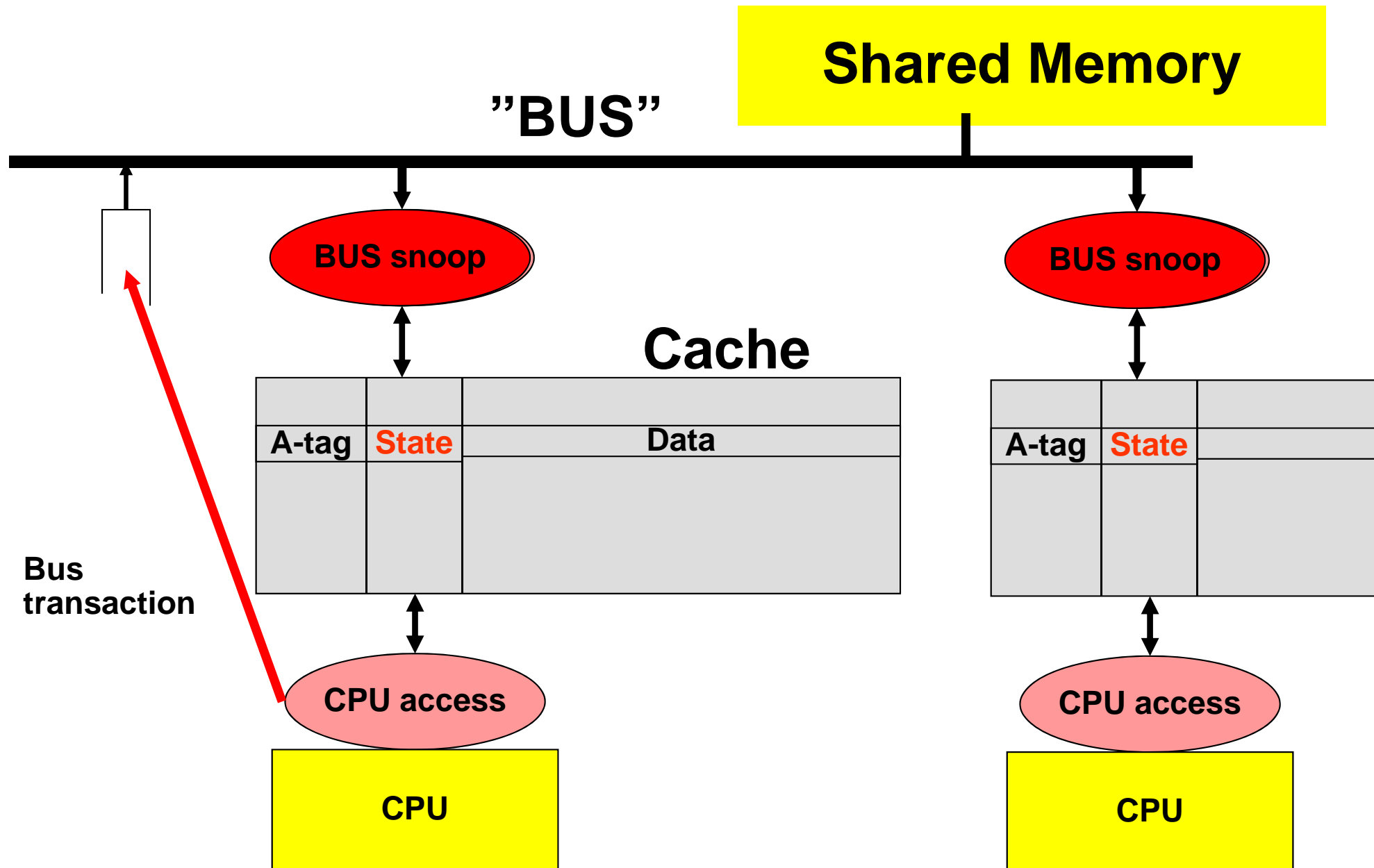
# Snoop-based Protocol Implementation

**Shared Memory**

**BUS**

**BUS snoop**

Per-cache-line "state" info

**Cache**

**Bus transaction**

| A-tag | S | Data |
|-------|---|------|
|       |   |      |

**CPU access**

"State machines"

**CPU**

UPPSALA UNIVERSITET

AVDARK 2013

# Cache implementation

State!

Cacheline here 64B:

| AT | S | Data = 64B |

Generic Cache:

MSB                    LSB        SRAM:

Addr [63..0]

index

Hit
Sel way "6"

...

mux

Data = 64B

**Which statement is true about this cache?**
❑ Coherence needs to be implemented on a cache line granularity
❑ Coherence with a word granularity can be implemented

MP 25

# Snoop-based Protocol Implementation

**Shared Memory**

**"BUS"**

**BUS snoop**

**BUS snoop**

**Cache**

| A-tag | State | Data |
|-------|-------|------|
|       |       |      |

| A-tag | State | |
|-------|-------|--|
|       |       |  |

**Bus transaction**

**CPU access**

**CPU access**

**CPU**

**CPU**

UPPSALA UNIVERSITET

AVDARK 2013

# Example: MOSI Bus Snoop

**Why is there no BUSinv arrow from M?**
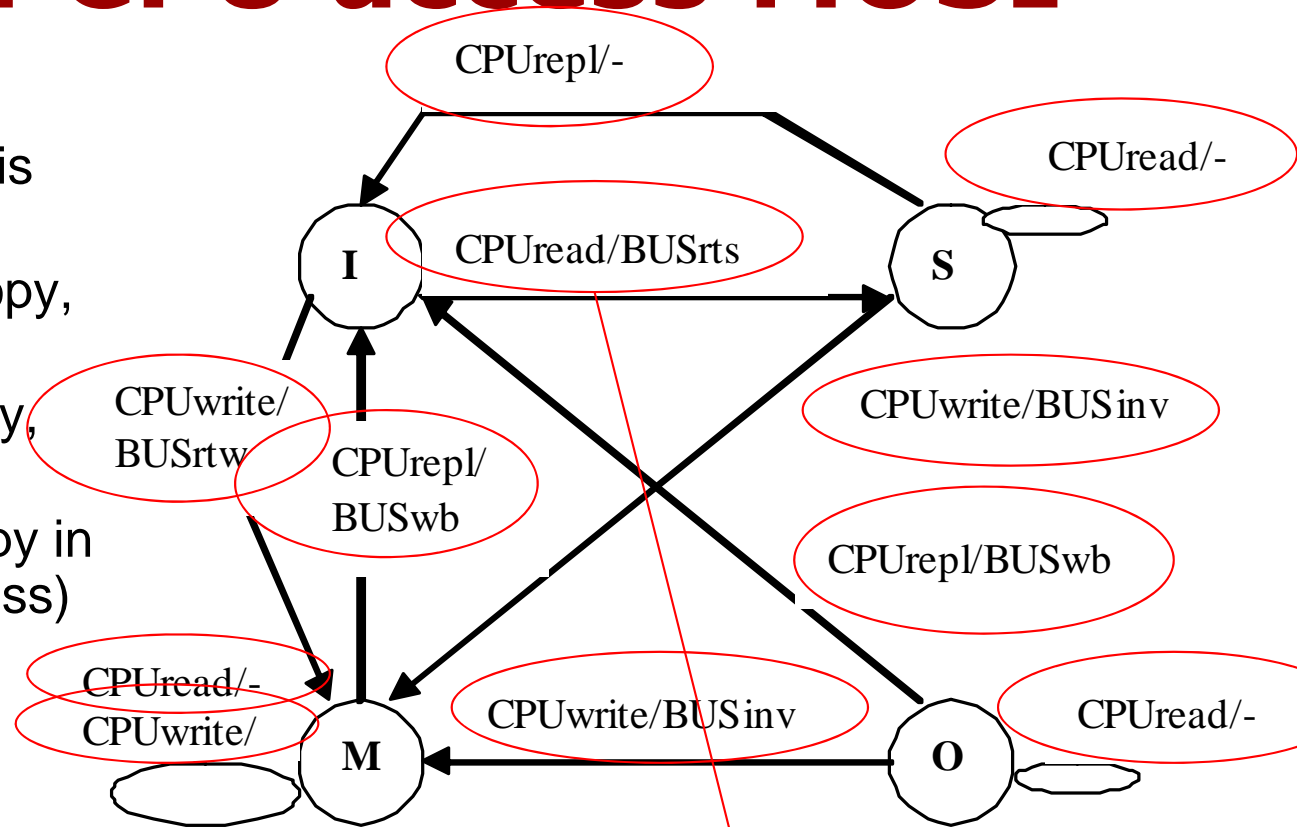- ❑ No other can have a cached copy
- ❑ BUSinv is only used for writes

## STATES:

**M – Modified**: My dirty* copy is the only cached copy

**S – Shared:** I have a clean copy, others may also have a copy

**O – Owner:** I have a dirty copy, others may also have a copy

**I – Invalid:** I have no valid copy in my cache (including cache miss)

## BUS TRANSACTIONS FROM OTHERS:

**BUSrts**  ReadtoShare. Reading the data

**BUSrtw:** ReadToWrite. Reading the data with the intention to modify it right away

**BUSinv**: Invalidating other caches copies

**BUSwb**: Writing data back to memory

*Dirty: my value differs from the old value in mem*

BUSrts
BUSrtw
BUSinv
BUSwb

BUSrts
BUSwb

BUSrtw   BUSinv

I          S

BUSrtw/Data

BUSrtw/Data
BUSinv

BUSrts/Data          O   BUSrts/Data

M

**Input-signal/Reply-signal**
Meaning: If you are in state M and  see BUSrts, goto state O and reply with Data

# Snoop-based Protocol Implementation

**Shared Memory**

**BUS**

**Cache**

**BUS snoop**

**BUS snoop**

| A-tag | State | Data |
|-------|-------|------|
|       |       |      |

| A-tag | State | |
|-------|-------|--|
|       |       |  |

**Bus transaction**

**CPU access**

**CPU access**

**CPU**

**CPU**

# Example: CPU access MOSI

## STATES:

**M – Modified**: My dirty* copy is the only cached copy

**S – Shared:** I have a clean copy, others may also have a copy

**O – Owner:** I have a dirty copy, others may also have a copy

**I – Invalid:** I have no valid copy in my cache (may be a cache miss)

CPUrepl/-

CPUread/-

CPUread/BUSrts

CPUwrite/BUSinv

CPUwrite/
BUSrtw

CPUrepl/
BUSwb

CPUrepl/BUSwb

I

S

CPUread/-
CPUwrite/

CPUwrite/BUSinv

CPUread/-

M

O

**Input-signal/Reply-signal**
Meaning:  If you are in state **I** and see CPUread, send a BUSrts and goto  S

## FROM MY CPU:

**CPUread**  Caused by a  Load instruction

**CPUwrite**: Caused by a Store or Atomic instruction

**CPUrepl**: Caused by a replacement of this cachline (caused by murphy ☺)

AVDARK
2013

# "Upgrade" in snoop-based



A:

B:

BUSinv

Have to INV

Have to INV

$

$

$

Thread

Thread

Thread

Read A
Read A
…
…

...
Read A
…
Write A

Read B
…
Read A

AVDARK 2013

# Summary

- Snooping was first suggested by Jim Goodman in ISCA, Stockholm 1984

- Effectively implements coherence through broadcast of "read and write misses"

- Best suited for on-chip coherence between a small number of caches
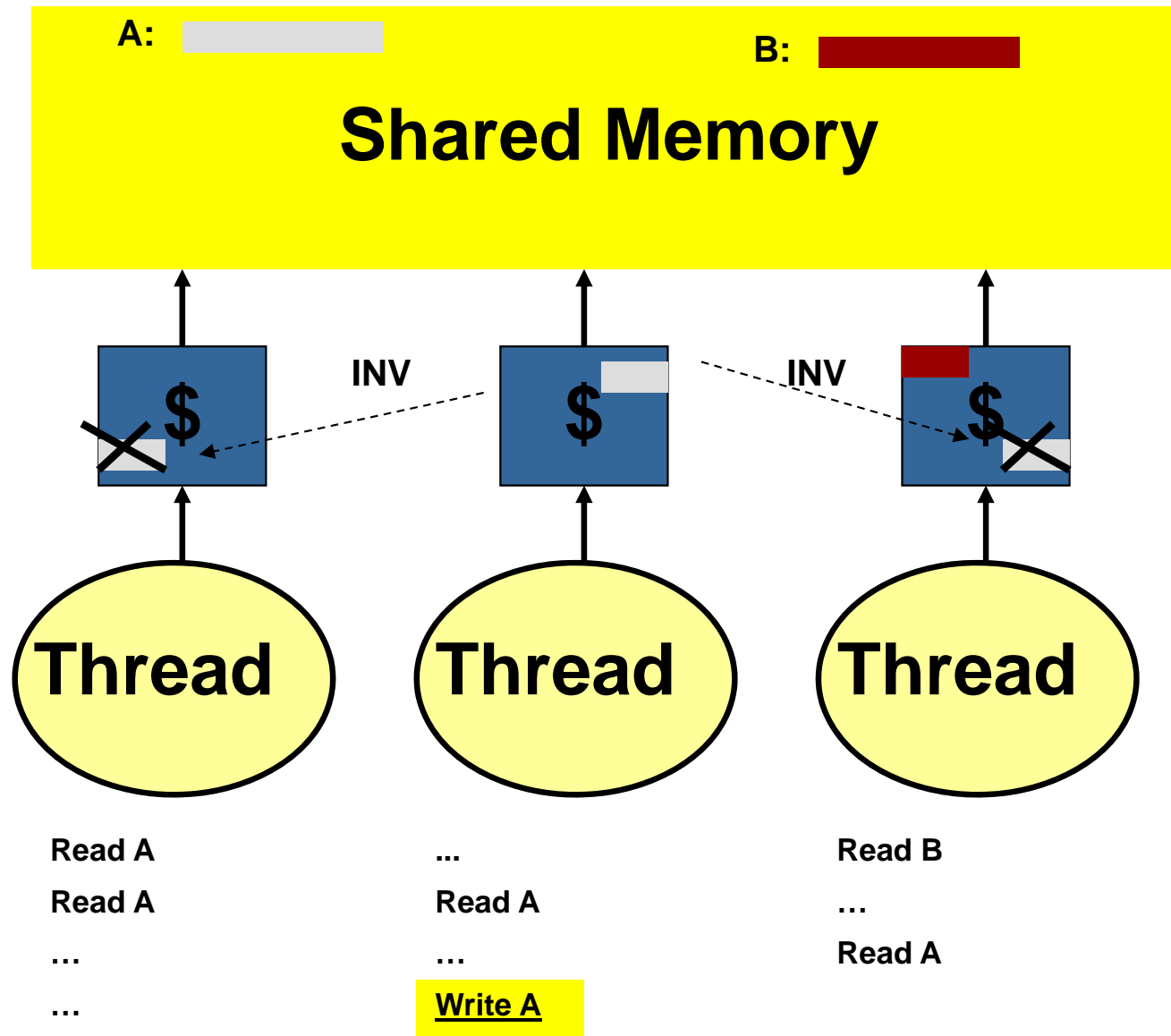
# MOSI Snooping Coherence Protocol Implementation

Erik Hagersten

Uppsala University

# Snoop-based Protocol Implementation



**UPPSALA UNIVERSITET**

**AVDARK 2013**

"BUS"

Shared Memory

BUS snoop

Cache

| A-tag | State | Data |
|-------|-------|------|
|       |       |      |
|       |       |      |

Bus transaction

CPU access

CPU

# Upgrade: Readable →Writable

AVDARK
2013

# Upgrade – the requesting CPU

**Defines action to CPU events:**

**CPUwrite**: Caused by a store miss

**CPUread** Caused by a load miss

**CPUrepl**: Caused by a replacement

**BUSinv**

snoop

BUSinv

Cache

A-tag | S→M | Data

access

Write

CPUrepl/-

CPUread/-

CPUread/BUSrts

I → S

CPUwrite/
BUSrtw

CPUrepl/
BUSwb

**CPUwrite/BUSinv**

CPUrepl/BUSwb

CPUread/-
CPUwrite/
-

M

CPUwrite/BUSinv

O

CPUread/-

AVDARK
2013

# Upgrade – the other CPUs

**Cache**

A-ta | S→I | Data

snoop

access

**CPU**

BUSrts
BUSrtw
BUSinv
BUSwb

**Defines action to Bus snoops:**

**BUSrts: ReadtoShare** (reading the cacheline)

**BUSrtw, ReadToWrite** (reading the cacheline with the intention to modify it right away)

**BUSwb**: Writing a cacheline back to memory

**BUSinv**: Invalidating other caches copies of the cacheline

BUSrts
BUSwb

BUSrtw  **BUSinv**

**I**                                **S**

BUSrtw

BUSrtw/Data
BUSinv

BUSrts/Data

BUSrts/Data

**M**                    **O**

AVDARK
2013

# **More Cache Lingo**

- **Capacity miss** – too small cache
- **Conflict miss** – limited associativity
- **Compulsory miss** – accessing data the first time
- **Coherence miss** – The cache would have had the data unless it had been invalidated by someone else
- **Upgrade miss:** (only for writes) – The cache would have had a writable copy, but answered a read request and "downgraded" itself to read-only state
- **False sharing:** Coherence/downgrade is caused by a shared <u>cacheline</u> and not by shared data:

| **False sharing example:** | **Read A** | **...** | cacheline: |
|---|---|---|---|
| | **...** | **Read D** | A, B, C, D |
| | **Write A** | **...** | |
| | **...** | **Write D** | |
| | **Read A** | | |

AVDARK
2013

# Implementing Snooping.
# One example
# (Sun E6000, ≈Intel P6,)

Erik Hagersten
Uppsala University

# Snoop-based architecture: Dual tags

**Shared Memory**

**BUS**

**BUS snoop**

| A-tag | State |
|-------|-------|

**Snoop Tag** (Obligation state)
(possibly time-sliced access
   to cache tags)

**Bus
transaction**

**Access Tag** (Permission sate)
(possibly time-sliced access
   to cache tags)

**Cache**

**CPU**

**Cache access**

| A-tag | State | Data |
|-------|-------|------|

# The Cache Coherent Memory System Coherent Write (Here: Write invalidate)



A:

B:

**Shared Memory**

INV

$ INV $

$ $

Thread

Thread

Thread

Read A

...

Read B

Read A

Read A

...

...

...

Read A

...

**Write A**

AVDARK 2013

# "Upgrade" in snooped-based

Shared Memory

"BusINV"

BUS snoop — A-tag S **M**

BUS snoop — A-tag S **I**

"right or left CPU"

"INV"

From earlier trans-actions

CPU: Store "middle CPU"

Cache access

A-tag **M**

A-tag **S**

AVDARK
2013

UPPSALA
UNIVERSITET

# The Cache Coherent Cache-to-cache

A: [          ]                    B: [████████]

## Shared Memory



**Thread**          **Thread**          **Thread**

Read A              ...                 Read B

Read A              Read A              …

…                   …                   Read A

…                   Write A

**Read A**

AVDARK
2013

UPPSALA
UNIVERSITET

# Cache2cache – the requesting CPU

**CPUwrite**: Caused by a store miss

**CPUread**  Caused by a loadmiss

**CPUrepl**: Caused by a replacement

**BUSrts**

snoop

Cache

A-ta | I→S | Data

Load

CPU

access

CPUrepl/-

CPUread/-

**I**    **CPUread/BUSrts**    **S**

CPUwrite/
BUSrtw

CPUrepl/
BUSwb

CPUwrite/BUSinv

CPUrepl/BUSwb

CPUread/-
CPUwrite/

**M**    CPUwrite/BUSinv    **O**    CPUread/-

# Cache-to-cache – the other CPU

**BUSrts**

**snoop**

**Cache**

A-tags | M→O | Data

**Data**
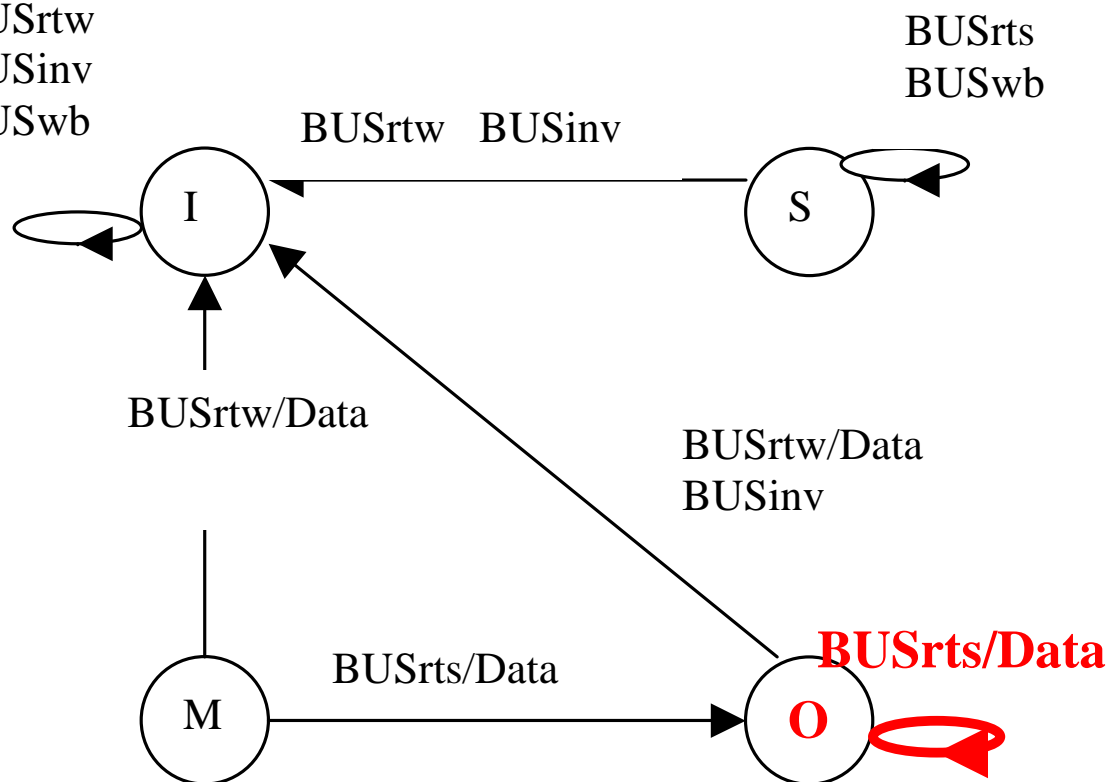
**access**

**CPU**

BUSrts
BUSrtw
BUSinv
BUSwb

**BUSrts: ReadToShare** (reading the data with the intention to read it)

**BUSrtw, ReadToWrite** (reading the data with the intention to modify it)

**BUSwb**: Writing data back to memory

**BUSinv**: Invalidating other caches copies

BUSrts
BUSwb

BUSrtw    BUSinv

I          S

BUSrtw/Data

BUSrtw/Data
BUSinv

M    **BUSrts/Data**    O    BUSrts/Data

**AVDARK 2013**

# Cache-to-cache in snoope-based



Shared Memory

BusRTS

BUS snoop   A-tag  I S

BUS snoop   A-tag  I O

Gotta' wait here for data

CPU: load

Cache access

A-tag  I S

A-tag  I O

"Left $"

"Middle $"

UPPSALA
UNIVERSITET

AVDARK
2013

# Yet Another Cache-to-cache

BUSrts

snoop

Cache

A-ta  O  Data

Data

access

CPU

**BUSrts: ReadtoShare** (reading the data with the intention to read it)

**BUSrtw, ReadToWrite** (reading the data with the intention to modify it)

**BUSwb**: Writing data back to memory

**BUSinv**: Invalidating other caches copies

BUSrts
BUSrtw
BUSinv
BUSwb

BUSrts
BUSwb

BUSrtw    BUSinv

I                    S

BUSrtw/Data

BUSrtw/Data
BUSinv

BUSrts/Data

M                    O    **BUSrts/Data**

UPPSALA UNIVERSITET

AVDARK 2013

# Summary

- Dual tags enable bus snoops and CPU lookups in parallel

- A datum may actually have several values "at the same wall-clock time"

- ... but not in "logic time": No software can detect that there are different values

- The value-change order maintained

AVDARK
2013

# Other Coherence Alternatives

Erik Hagersten
Uppsala University

# Common Cache States

- ## M – Modified
  My <u>dirty</u> (i.e. modified) copy is the only cached copy

- ## E – Exclusive
  My clean copy is the only cached copy

- ## O – Owner
  I have a dirty copy, others may also have a copy

- ## S – Shared
  I have a clean copy, others may also have a copy

- ## I – Invalid
  I have no valid copy in my cache

# Some Coherence Alternatives

Our first target

- **MOSI**
  - Leave one dirty copy in a cache on a cache2cache transfer

- **MSI**
  - Writeback to memory on a cache2cache.

- **MOESI**
  - The first reader will go to E and can later become a writer cheaply

AVDARK
2013

# Example A = A + 1
## Initially A is only in mem

**MOSI:**

| CPU | BUS | State |
|-----|-----|-------|
| LD A... | RTS(A) | S |
| ADD 1... | - | |
| ST A... | INV(A) | M |
| LD B | RTS(B) | S |
| ADD 1 | - | |
| ST B | INV(B) | M |
| ... | | |

**MOESI:**

| CPU | BUS | State |
|-----|-----|-------|
| LD A | RTS(A) | E |
| ADD 1 | - | |
| ST A | - | M |
| LD B | RTS(B) | E |
| ADD 1 | - | |
| ST B | - | M |
| ... | | |

AVDARK
2013

# Update-based MOSI protocol



A:

B:

**BusUpdate**

**Have to Update**

**My Update**

**Have to Update**

$

$

$

Thread

Thread

Thread

Read A

...

Read B

Read A

Read A

…

...

…

Read A

...

**Write A**

**Read A** ➔ HIT

AVDARK
2013

# Update-based MOSI protocol: Next write

A: [     ]                    B: [████]

BusUpdate

Have to Update                    My Update          Have to Update

$ | $ | $

Thread | Thread | Thread

| Read A | ... | Read B |
| Read A | Read A | ... |
| ... | ... | Read A |
| ... | Write A | |
| Read A | Write A | |

AVDARK
2013

# Update-based Coherence

- Write the new value to the other caches holding a shared copy (instead of invalidating…)

- Can avoid coherence misses

- May consume a large amount of snoop bandwidth

- Hard to implement some "memory models"

- Few commercial implementations: SPARCCenter2000, Xerox Dragon

AVDARK
2013

# Preparing for IRL

# Example during next IRL Class:

All the three RISC CPUs in a **MOSI** shared-memory (sequentially consistent) multiprocessor executes the following code almost at the same time:

```
while(A != my_id){};    /* this is a primitive kind of lock */
B = B + A;
A = A + 1;              /* this is a primitive kind of unlock */
while (A != 4) {};      /* this is a primitive kind of barrier sync */
<after a long time>
<some other execution replaces A and B from the caches, if still
present>
```

**Initially, CPU1 has its local variable my_id=1, CPU has my_id=2 and CPU3 has my_id=3 and the globally shared variables A is equal to 1 and B is equal to 0**.
Assume that CPU3, 2 and 1 first make one memory reference (i.e, a load or a store) each and then repeats that interleaving.

The following four bus transaction types can be seen on the snooping bus connecting the CPUs:
- **RTS:** ReadtoShare (reading the data with the intention to read it)
- **RTW,** ReadToWrite (reading the data with the intention to modify it)
- **WB:** Writing data back to memory
- **INV:** Invalidating other caches copies

Show every state <u>change</u> and/or value <u>change</u> of A and B in each CPU's cache according to one possible interleaving of the memory accesses. After the parallel execution is done for all of the CPUs, the cache lines still in the caches will be replaced. These actions should also be shown. For each line, also state what bus transaction occurs on the bus (if any) as well as which device is providing the corresponding data (if any).

**AVDARK 2013**

**UPPSALA UNIVERSITET**

# Example of a state transition sheet:

| CPU action | Bus Transaction (if any) | State/value after the CPU action | | | | | | Data is provided by [Cache 1, 2, 3 or Mem] (if any) |
|---|---|---|---|---|---|---|---|---|
| | | CPU1 A | B | CPU2 A | B | CPU3 A | B | |
| Initially | | I | I | I | I | I | I | |
| CPU3: LD A | RTS(A) | | | | | S/1 | | Mem |
| CPU2: LD A | RTS(A) | | | S/1 | | | | Mem |
| CPU1: LD A | RTS(A) | S/1 | | | | | | Mem |
| CPU3: LDA | — | | | | | | | — |
| | | | | | | | | |

AVDARK
2013

UPPSALA
UNIVERSITET

# What are Memory Models?

Erik Hagersten

Uppsala University

Sweden

# Where Memory Models Matter

- Flag synchronization

(initially flag = 0 and A = 0 )

```
…                          …
A = 1;                     while (flag != 1)  { };
flag = 1;                  X = A;
                           print(X);
```

- Causality (Causal correctness)

**(Initially A = 0 and  B = 0)**

```
…              …              Read A
A = 1;         …              …
…              while (A==0) {};   …
               B = 1;         …
                              while (B==0) {};
                              X = A;
                              print (X);
```

AVDARK
2013

# Dekker's Algorithm (mutual exclusion)

**Initially A = B = 0**

**"fork"**

**A := 1**

**if (B== 0) print("A won")**

**B := 1**

**if (A == 0)  print("B won")**

**Trick question**
**Is it possible that both threads "win"?**
- ❑ Yes
- ❑ No
- ❑ Undefined

UPPSALA
UNIVERSITET

AVDARK
2013

# Memory Ordering

- Coherence defines a per-datum valuechange order

- Memory model defines the valuechange order for all the data.

# Memory Ordering

- Defines the guaranteed memory ordering

- Is a "contract" between the HW and SW guys

- Without it, you may not be able to say much about the result of a parallel execution

# Observing order in SW
## In which order did A and B change value?

### Value order

LD newA; LD oldB → *newA before newB*

ST newA; LD oldB →*newA before newB*

ST newA; ST newB → *newA before newB*

### Program order

The order of program statements of each thread

AVDARK
2013

Dept of Information Technology| www.it.uu.se

© Erik Hagersten| user.it.uu.se/~eh

# One possible observed order

# Another possible observed order



**Thread 1 Thread 2**

**Thread 1 Thread 2**

# "The intuitive memory order" Sequential Consistency (Lamport)

**Shared Memory**

loads, stores

T T T T T T T T $T_{hread}$

* Global order achieved by *interleaving* <u>all</u> memory accesses from different threads

* SW should not be able to detect contradictive orders

* "Programmer's intuition is maintained"

* Unnecessarily restrictive ==> performance penalty

UPPSALA UNIVERSITET

# Where Memory Models Matters

- ## Flag synchronization

(initially flag = 0 and A = 0 )

```
…                              …
A = 1;                         while (flag != 1)  { };
flag = 1;                      X = A;
                               print(X);
```

- ## Causality (Causal correctness)

**(Initially A = 0 and  B = 0)**

```
…                   …                   Read A
A = 1;              …                   …
…                   while (A==0) {};    …
                    B = 1;              …
                                        while (B==0) {};
                                        X = A;
                                        print (X);
```

AVDARK
2013

# Dekker's Algorithm (mutual exclusion)

Initially A = B = 0

"fork"

A := 1

if (B== 0) print("A won")

B := 1

if (A == 0)  print("B won")

Given Sequential Consistecy:
Is it possible that both threads "win"?
- ❑ Yes
- ❑ No
- ❑ Undefined

AVDARK
2013

UPPSALA
UNIVERSITET

# Proving Dekker's Algorithm under SC

Initially A = B = 0

"fork"

A := 1

if (B== 0) print("A won")

B := 1

if (A == 0)  print("B won")

# Can the case "both win" happen under SC?

## Acess graph

c ┊ = VO: Value
d    order: c < d
(i.e., c happened before
d in the global order)

a ↓ = PO: Program
b    order: a < b
(the order specified
by the program)

A := B := 0

A:= 1          B:= 1

If (B == 0)      If (A == 0)
   print "Left wins"     print "Right wins"

A := B := 0

ST A, 1          ST B, 1

LD B ➔ 0      LD A ➔ 0

## Cyclic access graph ➔ Not SC (there is no global order)

AVDARK
2013

# One thread wins under SC

Only Right wins ➔  SC is OK

A := B := 0

A:= 1

If (B == 0)
   print "Left wins"

B:= 1

If (A == 0)
   print "Right wins"

A := B := 0

ST A, 1

ST B, 1

LD B ➔ 1

LD A ➔ 0

**Not cyclic graph ➔ SC**

**One global order:**
**STB < LDA < STA <LDB**

# No thread wins under SC

No thread wins ➔ SC is OK

A := B := 0

A:= 1

If (B == 0)
print "Left wins"

B:= 1

If (A == 0)
print "Right wins"

A := B := 0

ST A, 1                    ST B, 1

LD B ➔ 1                  LD A ➔ 1

**Not cyclic graph ➔ SC**

**Four Partial Orders, still SC**
**STB < LDA ;   STA < LDA;   STB < LDB ;   STA < LDA**

# Summary

- Gives the "illusion" of **one global order** between all memory accesses

- If two threads can observe two contratictive [partial] orders, SC is broken.

- Maintains human intuition

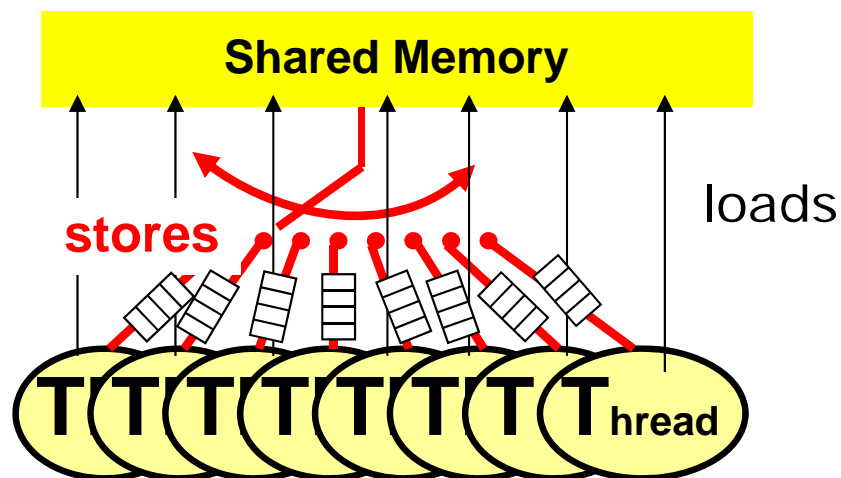- ... at the cost of performance (or complexity)

AVDARK
2013

# Other Memory Models

Erik Hagersten
Uppsala University

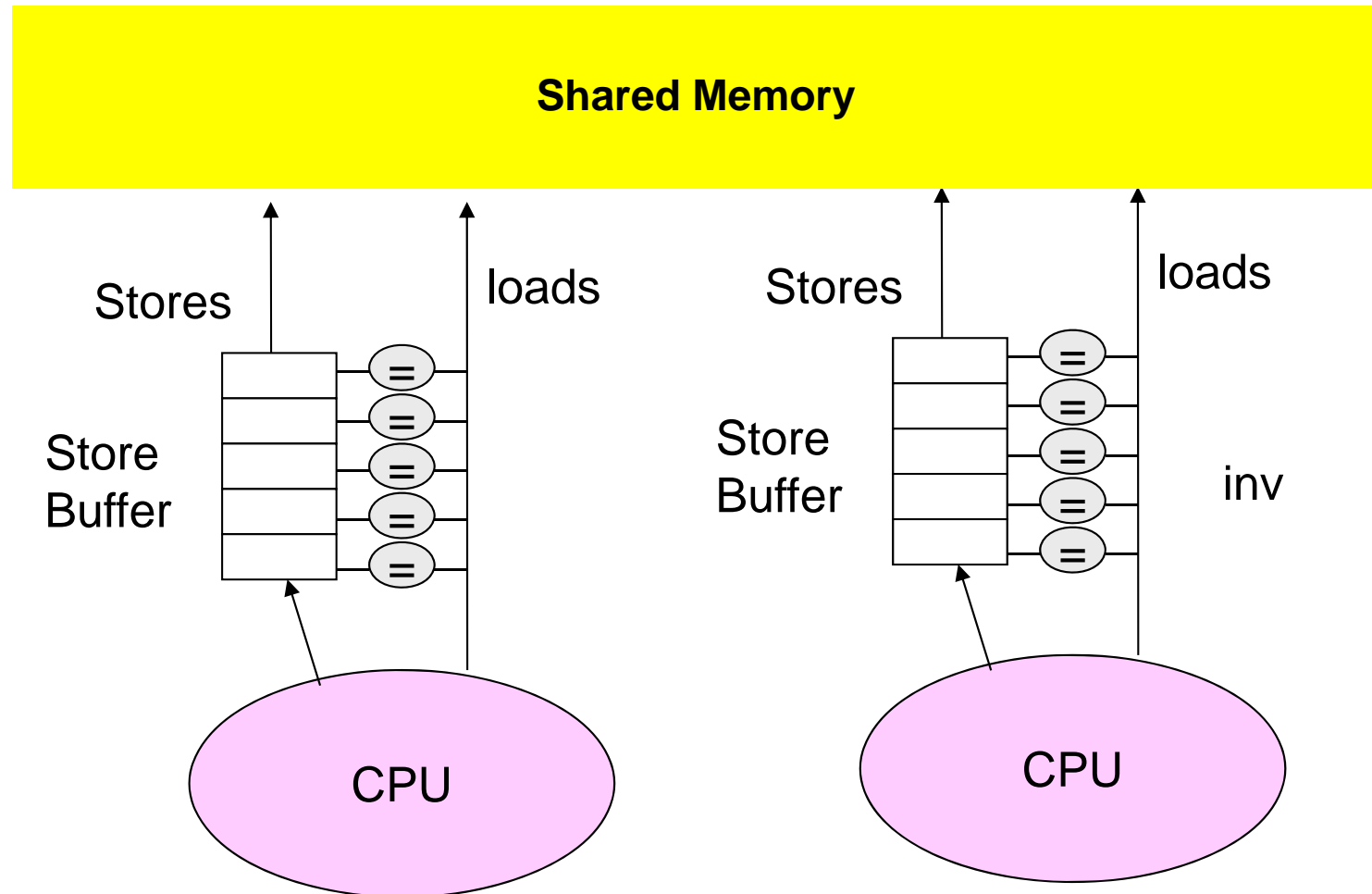# "Almost intuitive memory model" Total Store Ordering [TSO] (P. Sindhu)



* Global *interleaving* [order] for <u>all</u> stores from different threads (own stores excepted)

* "Programmer's intuition is almost maintained"
  * Flag synchronization? Yes
  * Store causality? Yes
  * Does Dekker work? No

* Unnecessarily restrictive ==> performance penalty

# TSO HW Model



→Stores are moved off the critical path
Coherence implementation can be the same as for SC

# Where Memory Models Matters

- ## Flag synchronization

  (initially flag = 0 and A = 0 )

  ```
  …                              …
  A = 1;                         while (flag != 1)  {};
  flag = 1;                      X = A;
                                 print(X);
  ```

- ## Causality (Causal correctness)

  **(Initially A = 0 and  B = 0)**

```
…                    …                    Read A
A = 1;               …                    …
…                    while (A==0) {};     …
                     B = 1;               …
                                          while (B==0) {};
                                          X = A;
                                          print (X);
```

# Dekker's Algorithm (mutual exclusion)

**Initially A = B = 0**

**"fork"**

**Does the write become globally visible before the read is performed?**

**A := 1**

**if (B== 0) print("A won")**

**B := 1**

**if (A == 0)  print("B won")**

**Given Total Store Order:**
**Is it possible that both threads "win"?**
- ❏ Yes
- ❏ No
- ❏ Undefined

**AVDARK**
**2013**

# Dekker's Algorithm for TSO

**Initially A = B = 0**

**"fork"**

**A := 1**

**Memory barrier**

**if (B== 0) print("A won")**

**B := 1**

**Memory barrier**

**if (A == 0)  print("B won")**

Memory barrier: Tells the HW to not start the LD until all previous stores have been "globaly ordered"

➔ behaves like SC

➔ Dekker's algorithm works!

AVDARK
2013

# Weak/release Consistency
## (M. Dubois, K. Gharachorloo)



**Shared Memory**

**loads stores**

Thread  *Thread*  *Thread*  *Thread*

* **Most accesses are unordered**
* **"Programmer's intuition is not maintained"**
  - Flag synchronization? No
  - Causal correctness? No
  - Dekker? No
* **Global order <u>only</u> established when the programmer explicitly inserts memory barrier instructions**
* ++ Better performance!!
* -- Interesting bugs!!
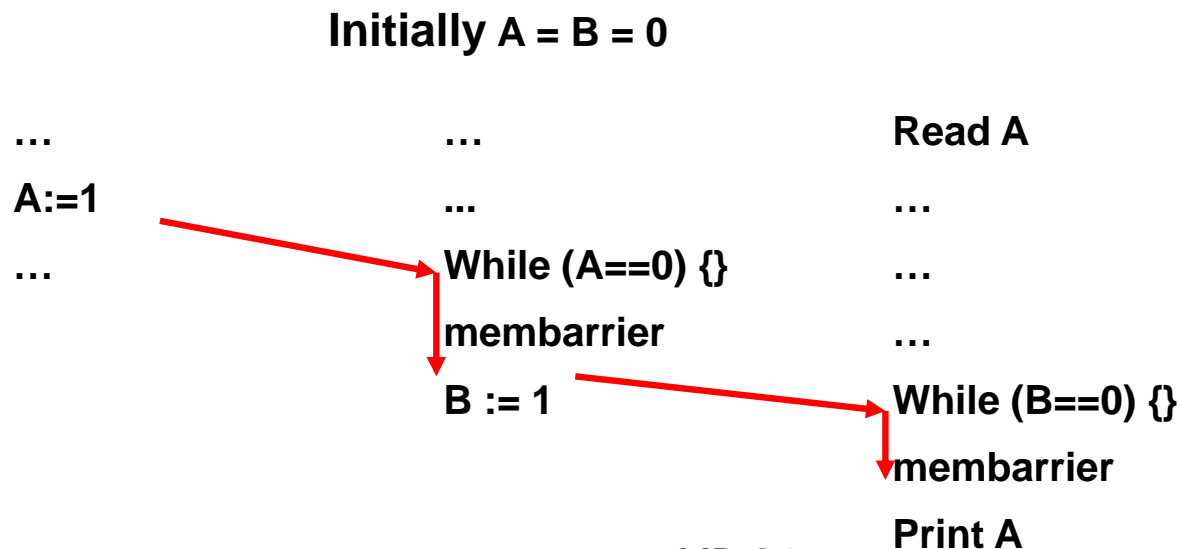
# Weak/Release consistency

- New flag synchronization needed
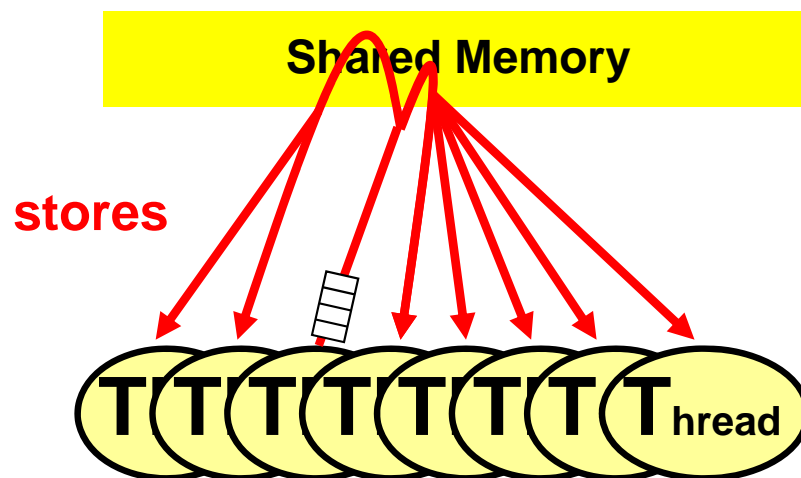
  A := data;                    while (flag != 1)  {};

  membarrier;                   membarrier;

  flag := 1;                    X := A;

- Dekker's: same as TSO

- Causal correctness provided for this code

**Initially A = B = 0**

| ... | ... | **Read A** |
|-----|-----|------------|
| **A:=1** | ... | ... |
| ... | **While (A==0) {}** | ... |
| | **membarrier** | ... |
| | **B := 1** | **While (B==0) {}** |
| | | **membarrier** |
| | | **Print A** |

AVDARK
2013

# Processor Consistency [PC] (J. Goodman)



**Shared Memory**

**stores**

T T T T T T T T**hread**

* PC: The stores from a processor appears to others in program order.

  - Flag synchronization? Yes

  - Causal correctness? Not clearly defined by Goodman. (yes, for PC "with causal correctness")

  - Dekker? No

AVDARK
2013

# Learning more about memory models

*Shared Memory Consistency Models: A Tutorial*
by Sarita Adve, Kouroush Gharachorloo
in IEEE Computer 1996 **(in the "Papers" directory)**

RTFM: Read the F*****n Manual of the system you are working on!
(Different microprocessors and systems supports different memory models.)

## Issue to think about:
What code reordering may compilers really do?
What does "volatile" declarations in C mean?

AVDARK
2013

# X86's new memory model

- Processor consistency with causual correctness for non-atomic memory ops

- TSO for atomic memory ops

- (Academia says the x86 mem model is TSO

- Link to the Intel manual (Section 8.2)

http://download.intel.com/products/processor/manual/325462.pdf

- Video presentation:
  http://www.youtube.com/watch?v=WUfvvFD5tAA&hl=sv

AVDARK
2013

# **Synchronization**

Erik Hagersten

Uppsala University

Sweden

**Execution on a sequentially consistent shared-memory machine:**

```
PSEUDO ASM CODE
            LD R1, #N
LOOP:       LD R2, (sum)
            SUB R1, R1, R2
            BGZ R3, CONT:
            LD R2, (sum)
            ADD R2, R2, #1
            ST R2, (sum)
            BR LOOP:


CONT:
```

"thread_create"

```
while (sum < N

    sum := s
```

```
while (sum < N

    sum := s
```

```
while (sum < N

    sum := s
```

```
while (sum < N)

    sum := sum + 1
```

"join"

`printf (sum)`

# Need to introduce synchronization

- Locking primitives are needed to ensure that only one process can be in the critical section:

```
LOCK(lock_variable)  /* wait for your turn */
if (sum > threshold) {
    sum := my_sum + sum
}
UNLOCK(lock_variable) /* release the lock*/
```

Critical Section

```
if (sum > threshold) {
    LOCK(lock_variable)  /* wait for your turn */
    sum := my_sum + sum
    UNLOCK(lock_variable) /* release the lock*/
}
```

Critical Section

AVDARK
2013

# Components of a Synchronization Event

- **Acquire method**
  - Acquire right to the synchronization (enter critical section, go past sync event)

- **Waiting algorithm**
  - Wait for synch to become available when it isn't

- **Release method**
  - Enable other processors to acquire right to the synch

# Atomic Instruction to Acquire

**Atomic example: test&set "TAS R1, (location)"**

The value at Mem(location) is loaded into R1, and

the constant "1" atomically stored into Mem(location)

(Other constant could be implemented, e.g., SPARC: "FF")

**Looks like a "store" to the coherence protocol**

**Implementation:**

1. Get a writable exclusive copy of the cache line (state M in MOSI)
2. Make the atomic modification to that cached copy

**Examples of other atomic primitives:**

**SWAP R1, (location):** atomically swap the values of R1 with Mem(location)

**CAS R1, R2, (location):** (Compare&Swap) SWAP if Mem(location)=REG2

AVDARK
2013

UPPSALA
UNIVERSITET

# Waiting Algorithms

## Blocking

- Waiting processes/threads are de-scheduled
- High overhead
- Allows processor to do "other things"

## Busy-waiting

- Waiting processes repeatedly test a lock_variable until it changes value
- Lower overhead, but consumes processor resources
- Can cause coherence network traffic

## Hybrid methods:  busy-wait a while, then block

# Release Algorithm

- Typically just a store "0"
- More complicated locks may require a conditional store or a "wake-up".

UPPSALA
UNIVERSITET

# A Bad Example: "POUNDING"

**How is TAS treated by the coherence protocol?**
- ❑ Like a CPU read operation
- ❑ Like a CPU write opreration
- ❑ By performing the "SWAP" atomically in DRAM

```
proc lock(lock_variable) {
    while (TAS[lock_variable]==1) {}    /* pound on the lock until free */
}

proc unlock(lock_variable) {
    lock_variable := 0
}
```

**If two threads are waiting for the lock**
- ❑ They will both spin locally in their cache
- ❑ They will create coherence traffic by invalidating each other
- ❑ They will both block and need to be woken up by the OS

*Assume: The function TAS[addr] returns the current memory value at addr and **atomically** writes the busy pattern "1" to the memory*

**Spinning threads produce traffic!**

UPPSALA
UNIVERSITET

# Optimistic Test&Set Lock "spinlock"

```
proc lock(lock_variable) {
    while true {
      if (TAS[lock_variable] ==0)  break;      /* pound on the lock once, done if TAS==0 */
      while(lock_variable != 0) {}       /* spin locally in your cache until "0" observed*/
  }
}


proc unlock(lock_variable) {
      lock_variable := 0
}
```

**If two threads are waiting for the lock**
- ❑ They will mostly spin locally in their cache
- ❑ They will create coherence traffic all the time by invalidating each other
- ❑ They will both block and need to be woken up by the OS

**Much less coherence traffic!!**
**-- still lots of traffic at lock handover!**
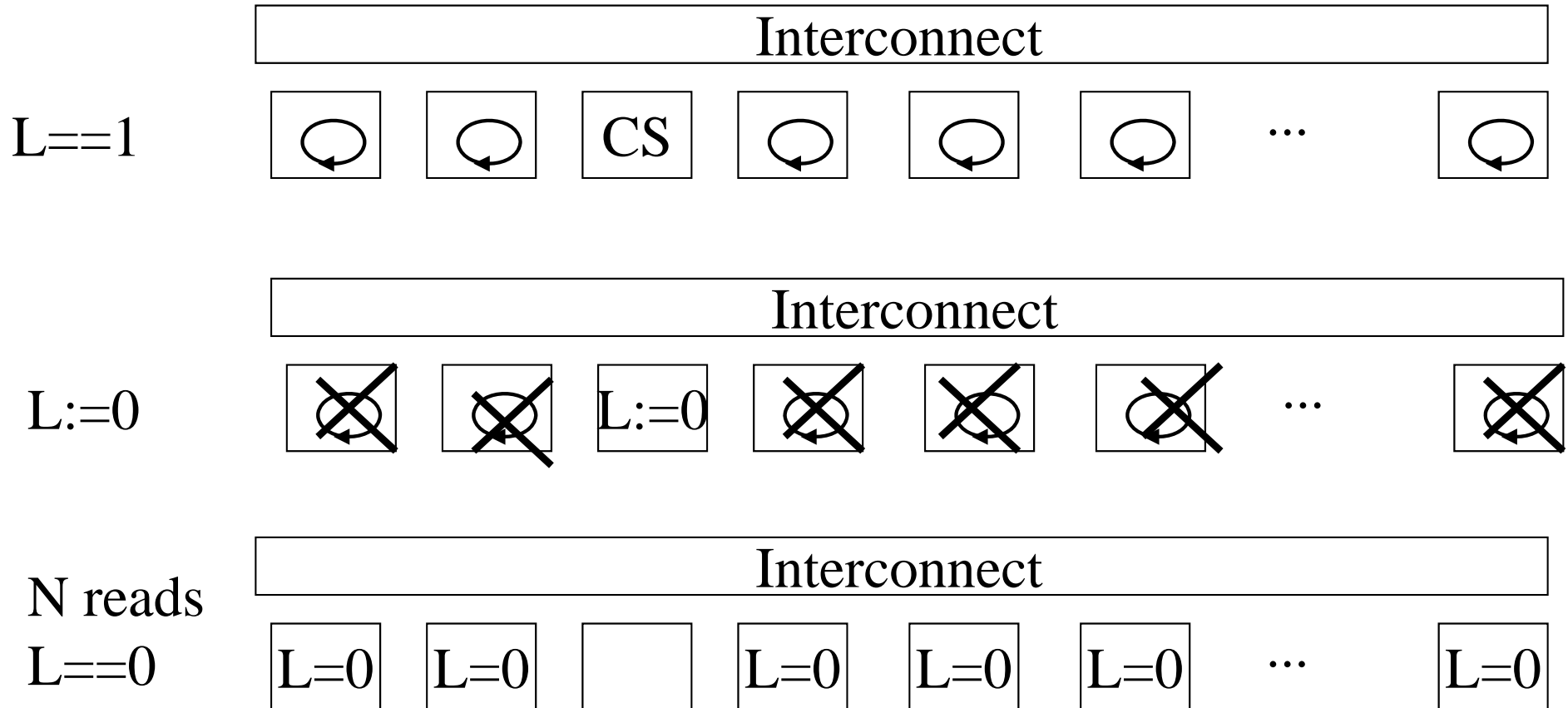**More on this during Scalable Synchronization**

AVDARK
2013

# Pesimistic Test&Set Lock "spinlock"

```
proc lock(lock_variable) {
    while true {
        while(lock_variable != 0) {}   /* spin locally in your cache until "0" observed*/
        if (TAS[lock_variable] ==0)  break; /* pound on the lock once, done if TAS==0
    }
}


proc unlock(lock_variable) {
    lock_variable := 0
}
```

**Slightly less traffic than Optimistic for contended locks
-- still lots of traffic at lock handover!
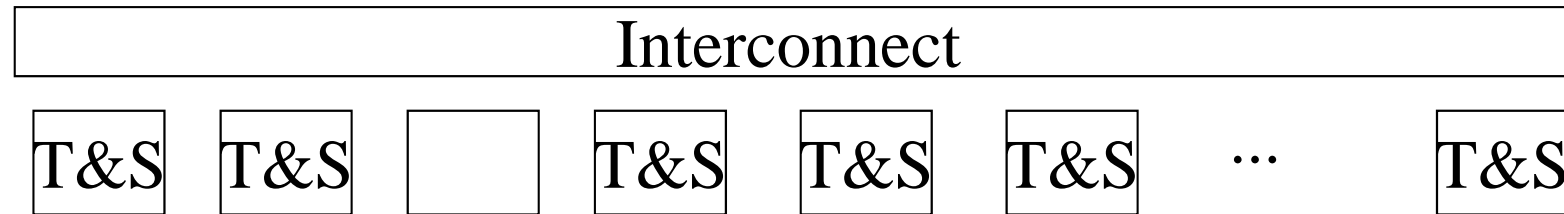More solutions during Scalable Shared Memory**

AVDARK
2013

UPPSALA
UNIVERSITET
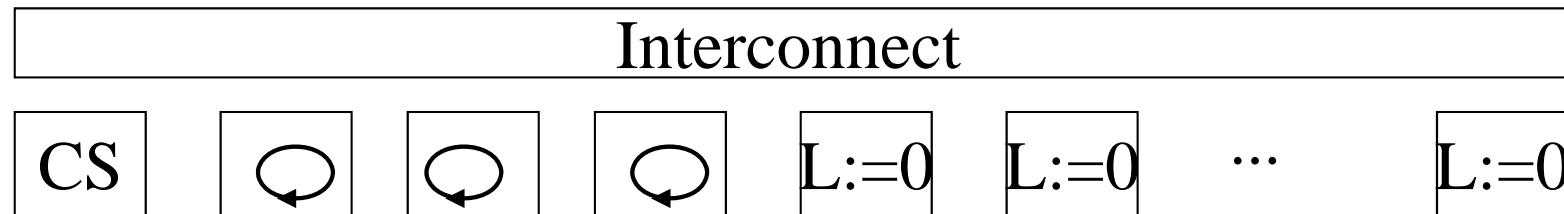
# It could still get messy!



L==1

Interconnect

| $\circlearrowleft$ | $\circlearrowleft$ | CS | $\circlearrowleft$ | $\circlearrowleft$ | $\circlearrowleft$ | $\cdots$ | $\circlearrowleft$ |

L:=0

Interconnect

| ⊗ | ⊗ | L:=0 | ⊗ | ⊗ | ⊗ | $\cdots$ | ⊗ |

N reads
L==0

Interconnect

| L=0 | L=0 | | L=0 | L=0 | L=0 | $\cdots$ | L=0 |

UPPSALA
UNIVERSITET

AVDARK
2013

# ...messy (part 2)

**N-1 Test&Set (i.e., N writes)**

Interconnect

| T&S | T&S | | T&S | T&S | T&S | ... | T&S |
|-----|-----|--|-----|-----|-----|-----|-----|

**L== 1**

Interconnect

| CS | ↺ | ↺ | ↺ | L:=0 | L:=0 | ... | L:=0 |
|----|---|---|---|------|------|-----|------|

potentially: ~N*N/2 reads :-(

Problem1: Contention on the interconnect slows down the CS execution
Problem2: The lock hand-over time is N*read_throughput
Fix1: Some back-off strategy, bad news for hand-over latency
Fix2: Queue-based locks

AVDARK 2013
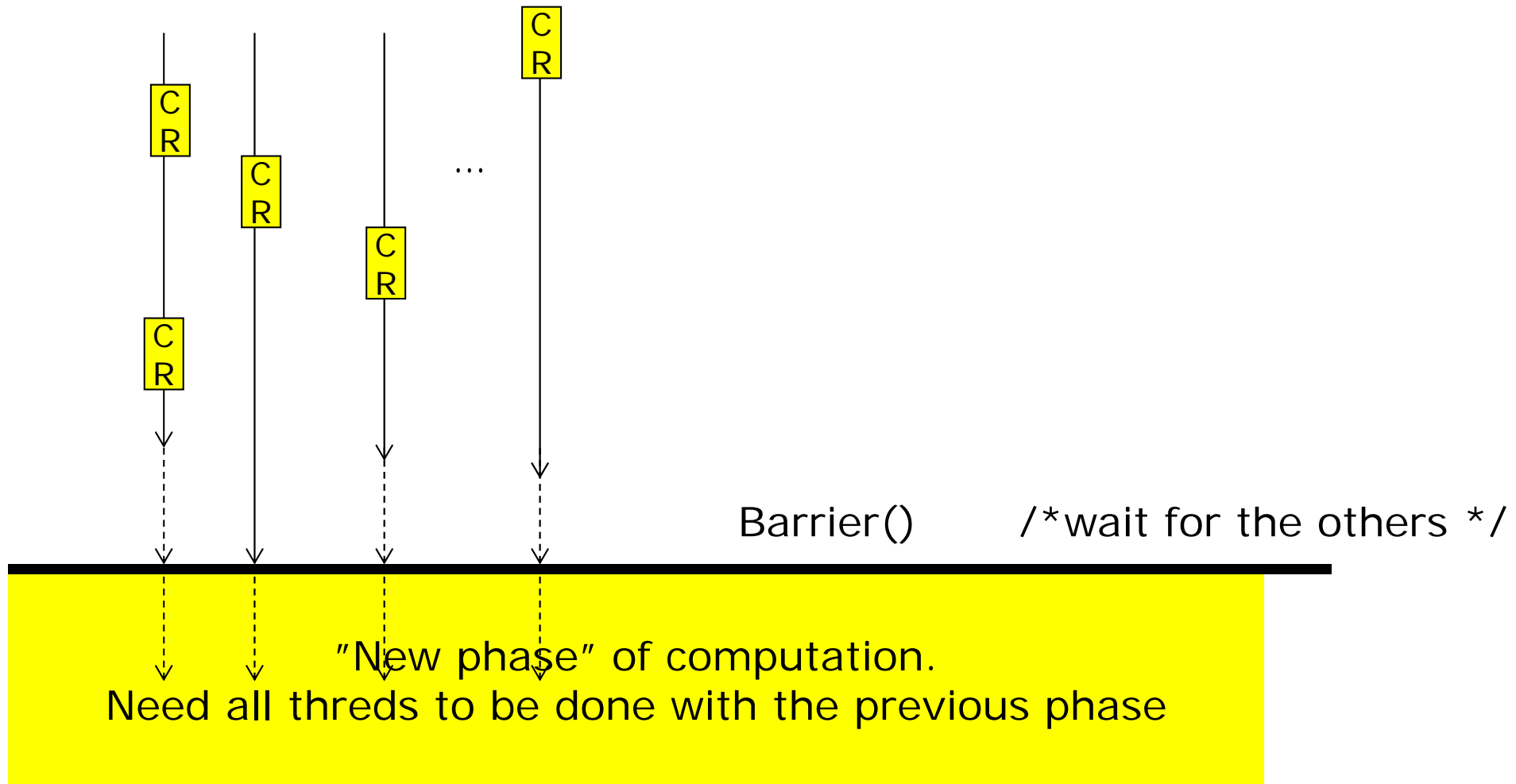
# Barrier Synchronization

Erik Hagersten
Uppsala University

# Barrier Synchronization



Barrier()        /*wait for the others */

"New phase" of computation.
Need all threds to be done with the previous phase

AVDARK
2013

# Barriers: Make the first threads wait for the last thread to reach a point in the program

1. Software algorithms implemented using locks, flags, counters

2. Hardware barriers
   * "Wired-AND" line separate from address/data bus
   * Set input high when arrive, wait for output to be high to leave
   * (In practice, multiple wires to allow reuse)
   * Difficult to support arbitrary subset of processors

# A Naiive Centralized Barrier

```
BARRIER (bar_name, p) {

  LOCK(bar_name.lock) {
      if (bar_name.counter == p) bar_name.counter = 0;  /* init count*/
      bar_name.counter++;                               /* globally increment the barrier count */
  }
  UNLOCK(bar_name.lock)

  while (bar_name.counter < p) {};      /* wait for the last thread */

}
```

# A More Complicated Centralized Barrier

```
BARRIER (bar_name, p) {
int loops;
loops = 0;

local_sense = !(local_sense)  ;              /* toggle private sense variable
                                               each time the barrier is used */

 LOCK(bar_name.lock);
    bar_name.counter++;                       /* globally increment the barrier count */
    if (bar_name.counter == p)  {             /* everybody here yet ? */
         bar_name.flag = local_sense;         /* release waiters*/
         UNLOCK(bar_name.lock)
     }
    else
      {  UNLOCK(bar_name.lock);
         while (bar_name.flag != local_sense) {      /* wait for the last guy */
           if (loops++ > UNREASONABLE) report_warning(pid)}
      }
```

AVDARK
2013