# Background to Multicores

Erik Hagersten

Uppsala University
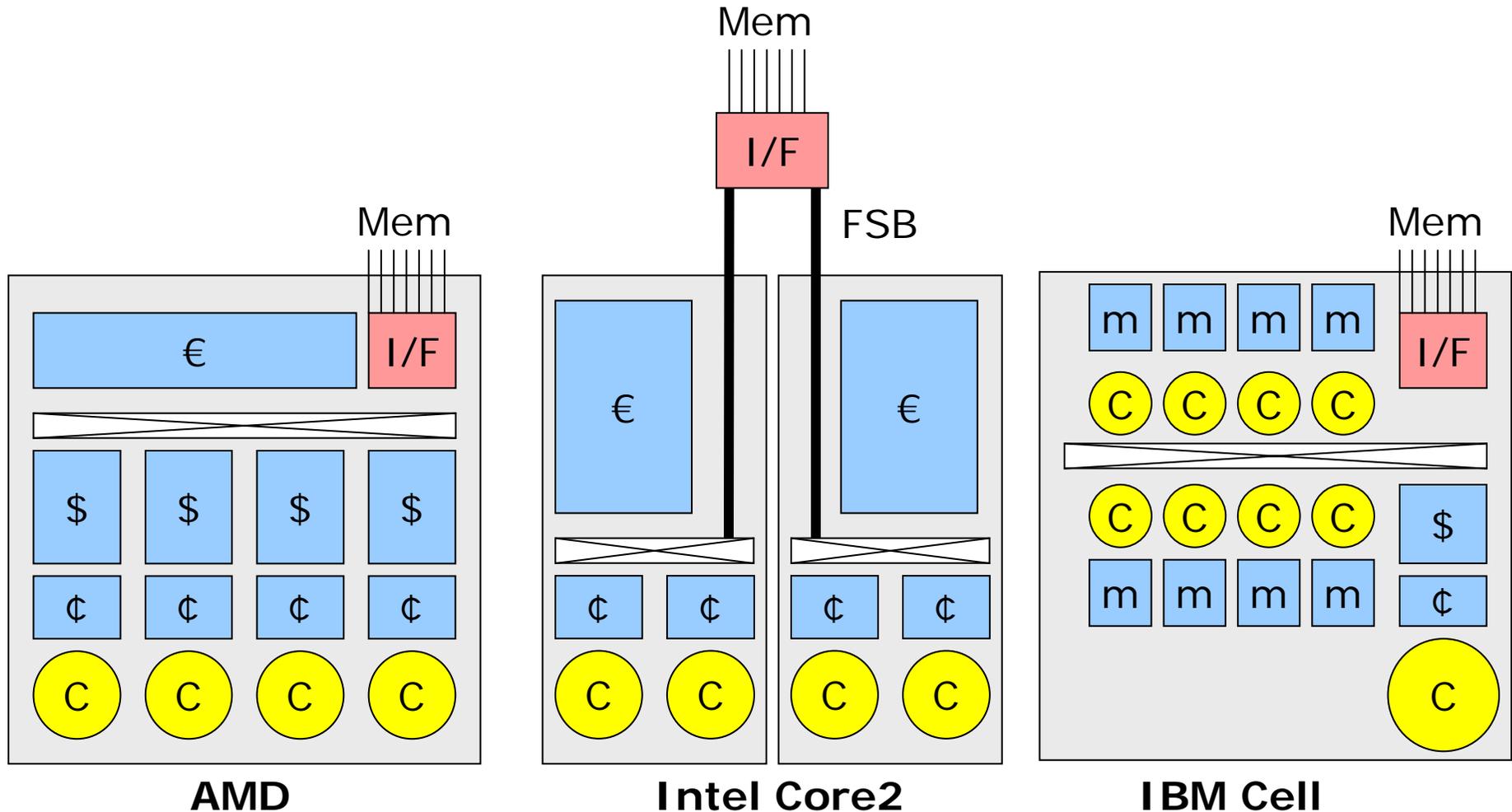
# The era of the "supercomputer" multiprocessors

- The one with the most blinking lights wins
- The one with the strangest languages wins
- The niftier the better!

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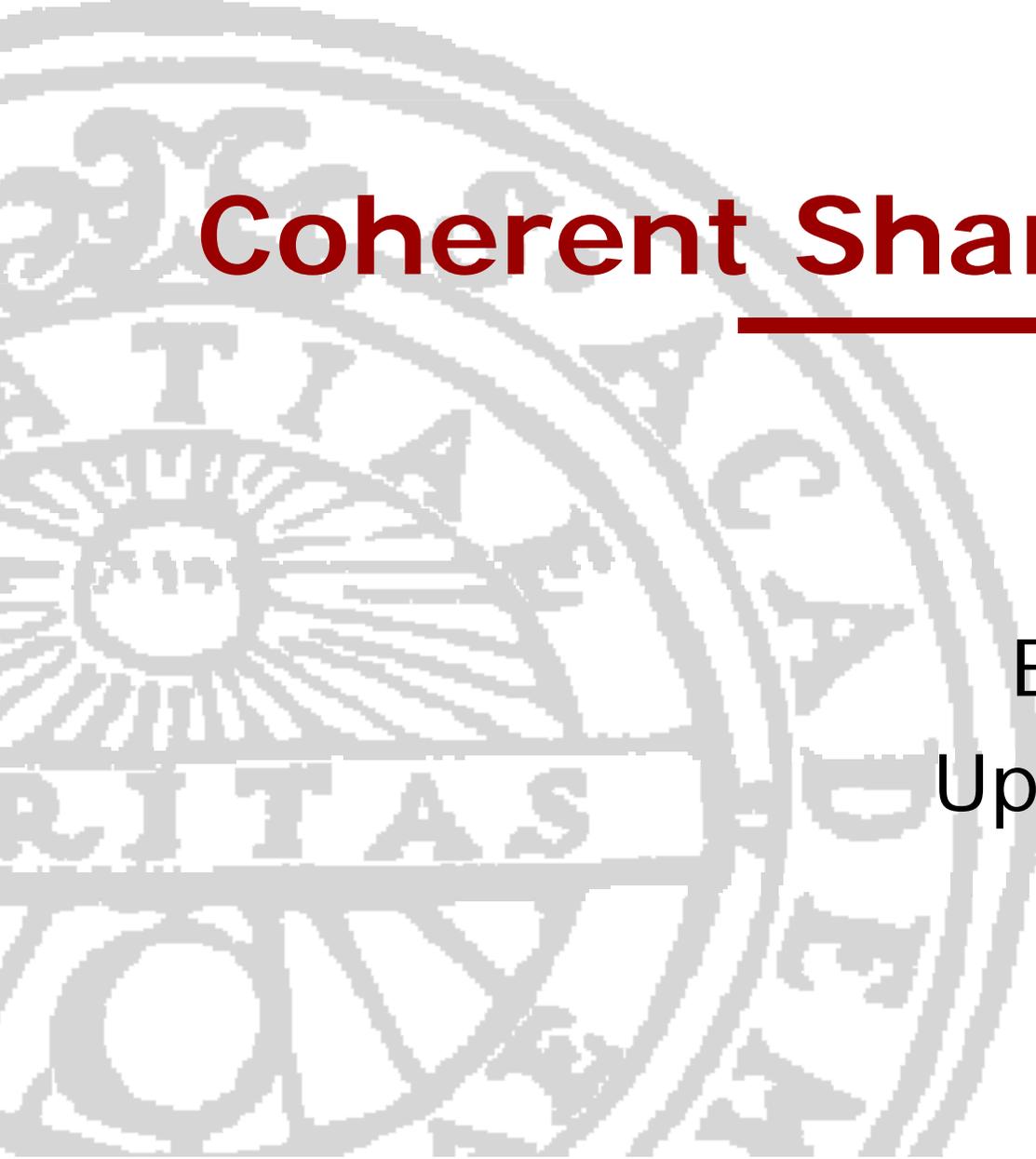


**AMD**

**Intel Core2**

**IBM Cell**

MP 3

# Taxomy for Architectures [Flynn]

*SIMD*        *MIMD*

*Message-passing*        *Shared Memory*

*Focus of this multicore lecture*

**Fine-grained**    **Coarse-grained**        **UMA**    **NUMA**    **COMA**

PDC
Summer
School
2010

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

MP 4

© Erik Hagersten| user.it.uu.se/~eh

# Outline

- Recap
  - Caches and Coherence

- More insight
  - Coherence
  - Memory models

- Example
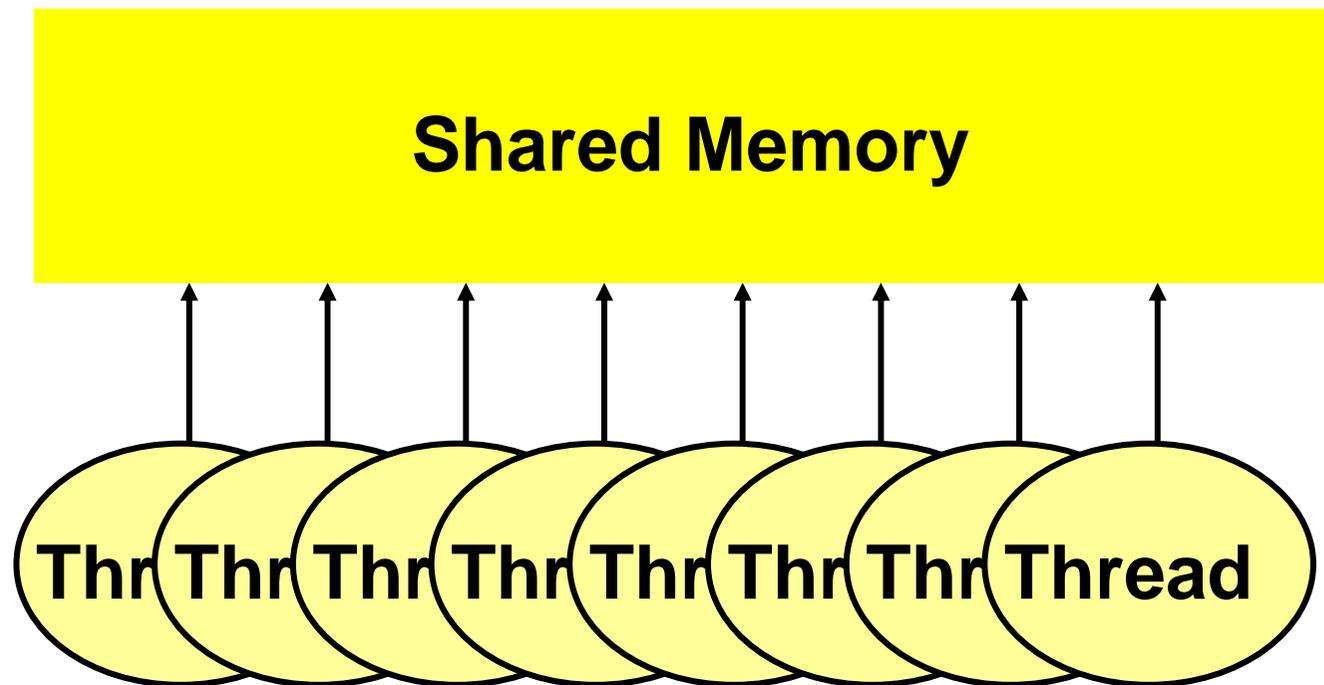  - Leveraging coherence properties for efficient synchronization

PDC
Summer
School
2010

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

MP 5

© Erik Hagersten| user.it.uu.se/~eh

# Coherent Shared Memory
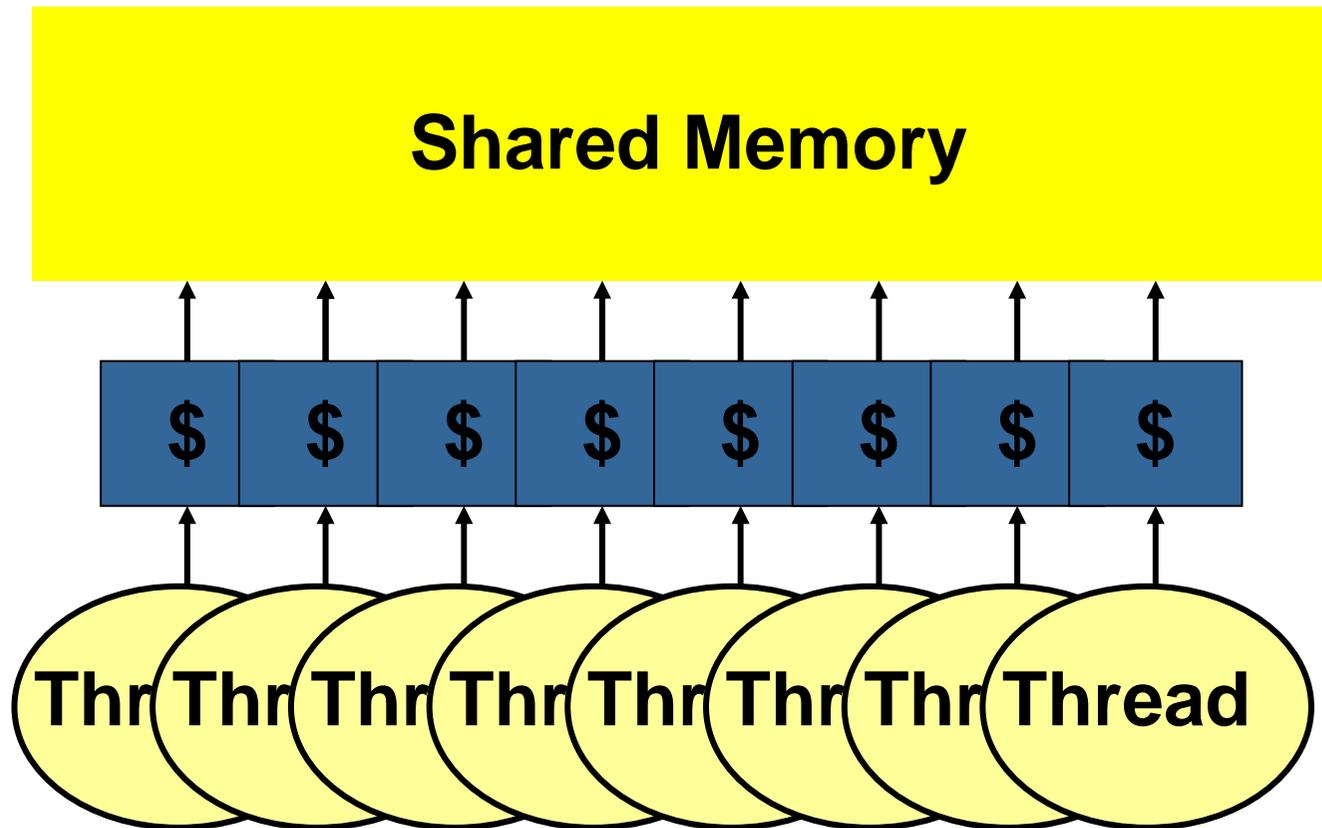
Erik Hagersten

Uppsala University

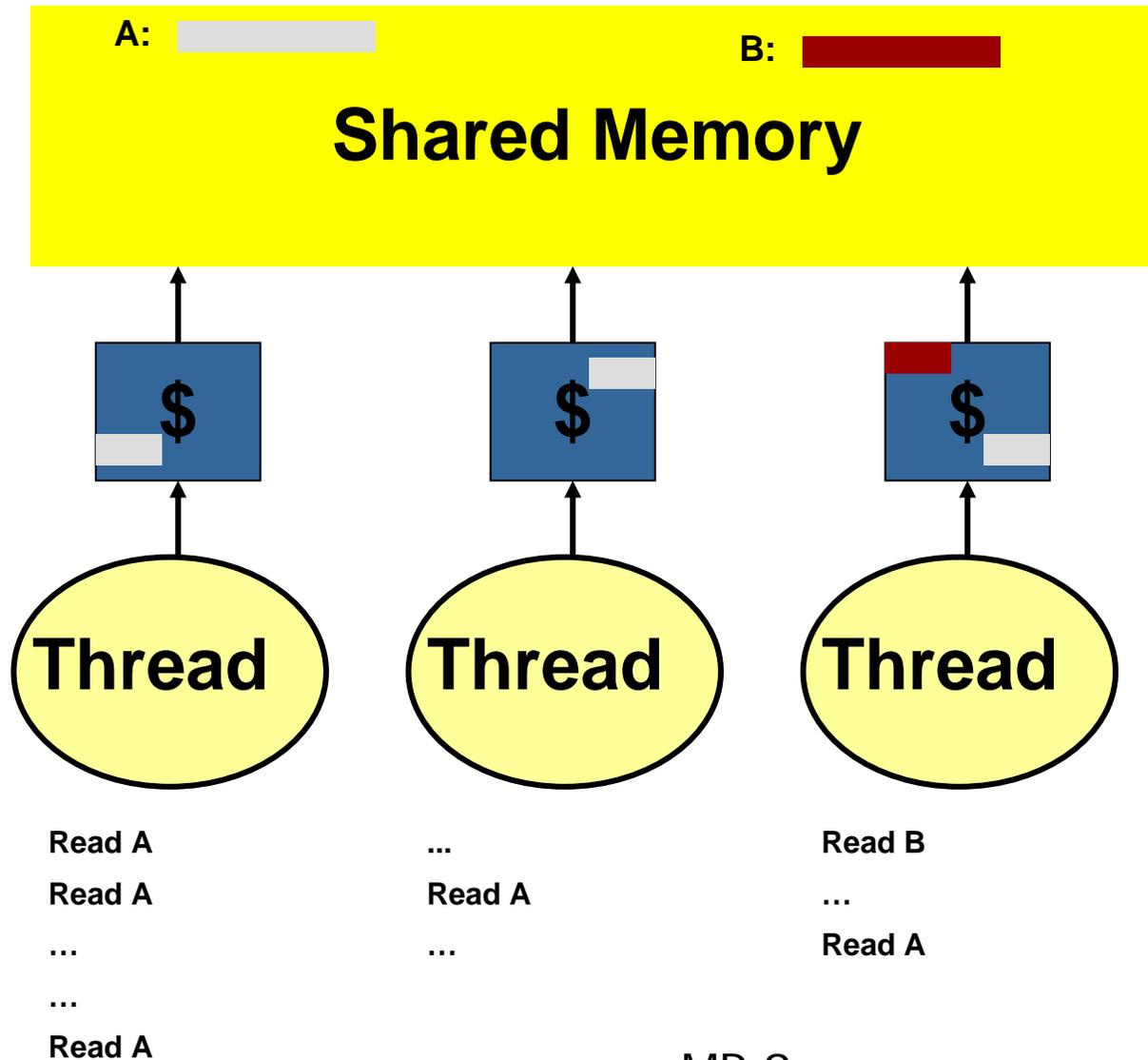# The Shared Memory Programming Model (Pthreads/OpenMP, ...)
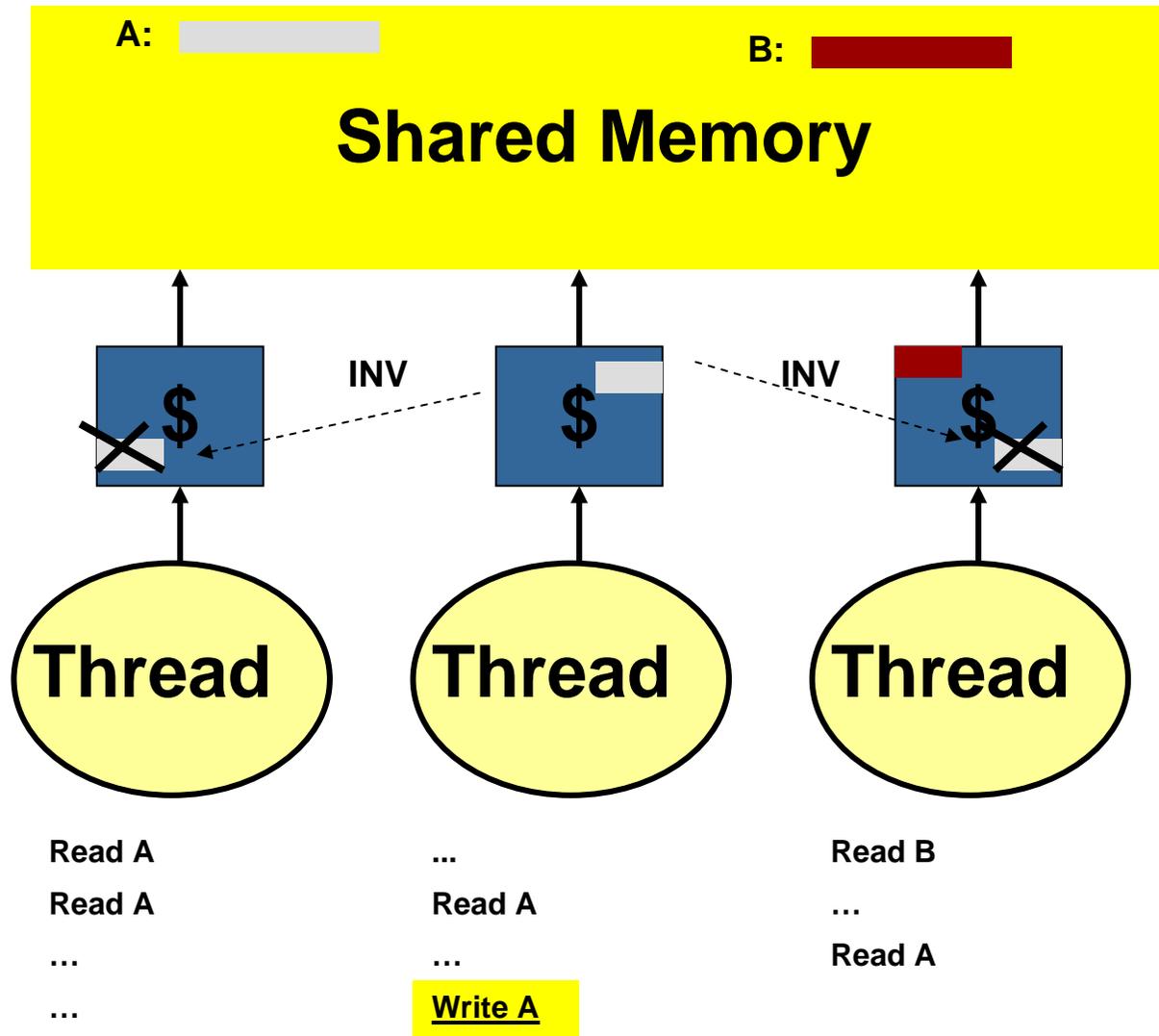
MP 7

# Adding Caches: More Concurrency

# Automatic Replication of Data



A:

B:

**Shared Memory**

$

$

$

**Thread**

**Thread**

**Thread**

Read A

...

Read B

Read A

Read A

…

...

...

Read A

...

...

Read A

# The Cache Coherent Memory System

A: █████████          B: ███████████

## Shared Memory

$    INV    $    INV    $

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

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

# The Cache Coherent Memory System

**A:** **B:**

## Shared Memory

$

$

$

**Thread**

**Thread**

**Thread**

Read A

...

Read B

Read A

Read A

…

…

…

Read A

…

Write A

PDC
Summer
School
2010

UPPSALA
UNIVERSITET

# The Cache Coherent Memory System

A: [     ]          B: [████]

**Shared Memory**

$

$

$

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

Read A            ...                 Read B

Read A            Read A              …

…                 …                   Read A

…                 Write A

**Read A**

# Snoop-based Protocol Implementation

**BUS**

**Shared Memory**

**BUS snoop**

**Cache**

**Bus transaction**

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

**CPU access**

**CPU**

# "Upgrade" in snoop-based

A: ▭   B: ▬

**BusINV**

Have to INV

My INV

Have to INV

$

$

$

**Thread**

**Thread**

**Thread**

Read A

Read A

...

...

...

Read A

...

**Write A**

Read B

…

Read A

# Cache-to-cache in snoop-based

A: ▭    B: ▬

**BusRTS**

**My RTS → wait for data**

**Gotta answer**

$ | $ | $

**Thread** | **Thread** | **Thread**

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

# Directory-based coherence: per-cachline info in the memory

A:

**Who has a copy**

B:

**Who has a copy**

**Directory Protocol**

$ | State

$ | State

$ | State

Cache access

Cache access

Cache access

**Thread**

**Thread**

**Thread**

MP 16

# "Upgrade" in dir-based



A:

Who has a copy

B:

Who has a copy

INV    ACK INV    INV

network

ACK    ACK

$    $    $

Thread    Thread    Thread

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

# Cache-to-cache in dir-based

A: [    ]    **Who has a copy**

B: [■■■■■]    **Who has a copy**

**ReadRequest**

**ReadDemand**

network

**Ack**

**Forward**

$  $  $

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

Read A      ...          Read B

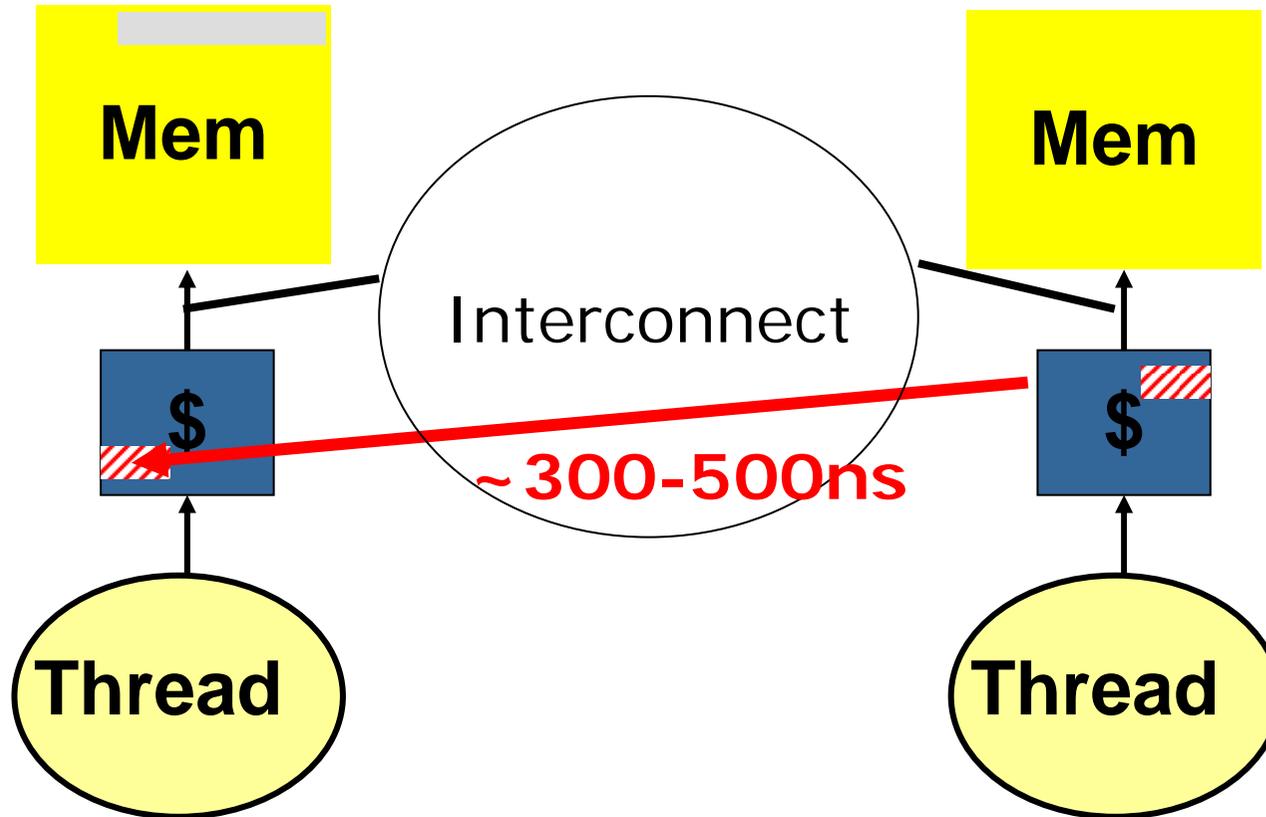Read A      Read A       …

…           …            Read A

…           **Write A**

**Read A**

# Non-uniform Architectures: NUMA & Multisocket→ Communication cost is much worse!
# A case for directory-based coherence

**Mem**

**Interconnect**

**Mem**

**$**

**~300-500ns**

**$**

**Thread**

**Thread**

Read A

Read A

...

...

**Read A**

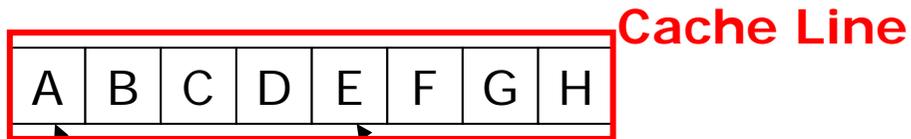...

Read A
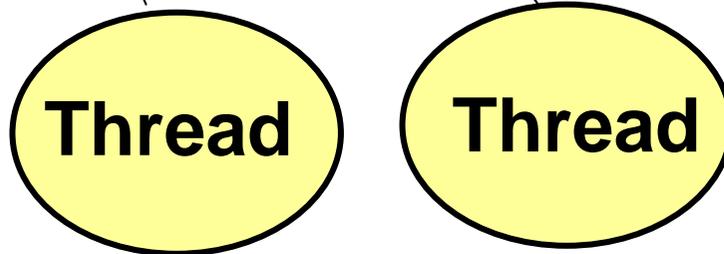
...

Write A

MP 19

# Why do you miss in a cache?

- Capacity misses – *the cache is too small*

- Conflict misses – *the cache organization is not perfect*

- Compulsory misses – *touching the data for the first time*

- Coherence misses – *caused by communication*

PDC
Summer
School
2010

# False sharing

**Cache Line**

| A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|

**Communication misses even though the threads do not share data "the cache line is too large"**

**Thread**

**Thread**

Read A
Write A
…
…
Read A

Read E
…
Write E

# So Far...

## Coherent shared memory

- "Snoopy-based coherence protocol"
  - All global accesses are broadcasted to all caches
  - Cache lines are automatically invalidated/fetched
- Directory-based coherence
- Ensures ordering and serialization for a single cache line

Sloppy definition: *There can be several copies of a datum, but only one value at "a given point in time"*

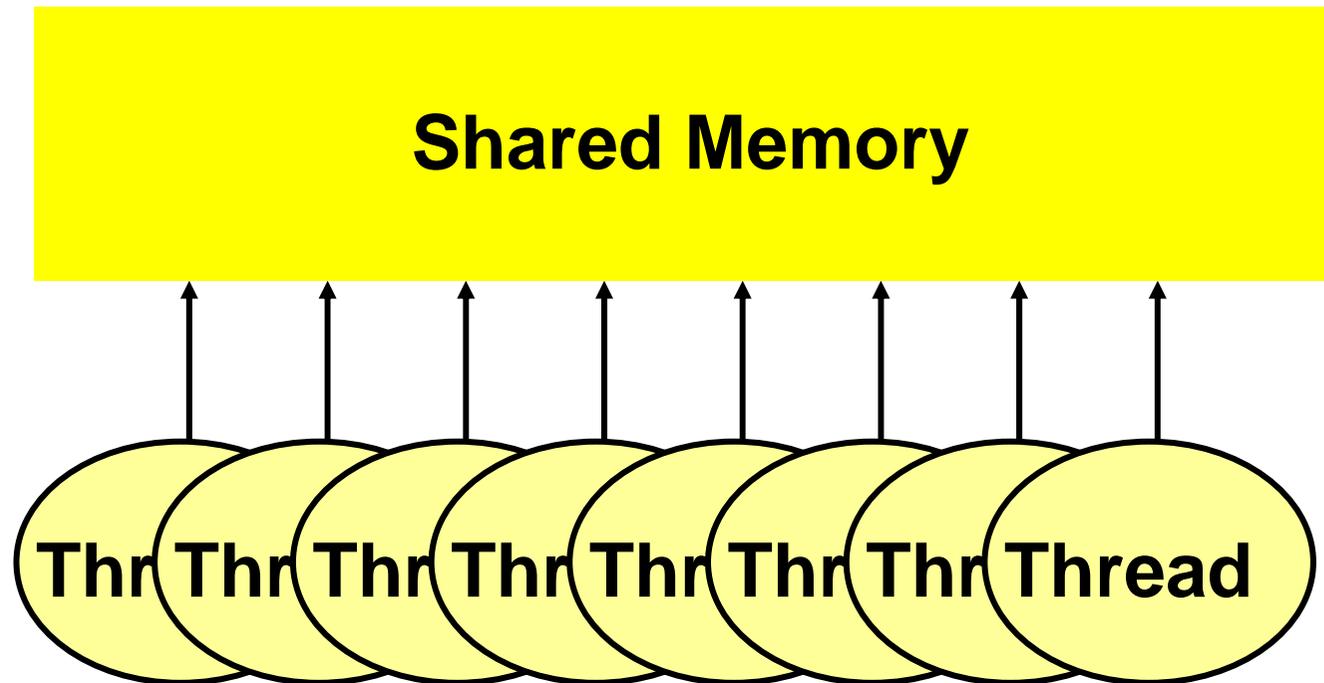Better definition: *Only a single value-change order can be observed for each datum.*

PDC Summer School 2010

UPPSALA UNIVERSITET

# Memory Ordering (aka Memory Consistency) -- tricky but important stuff

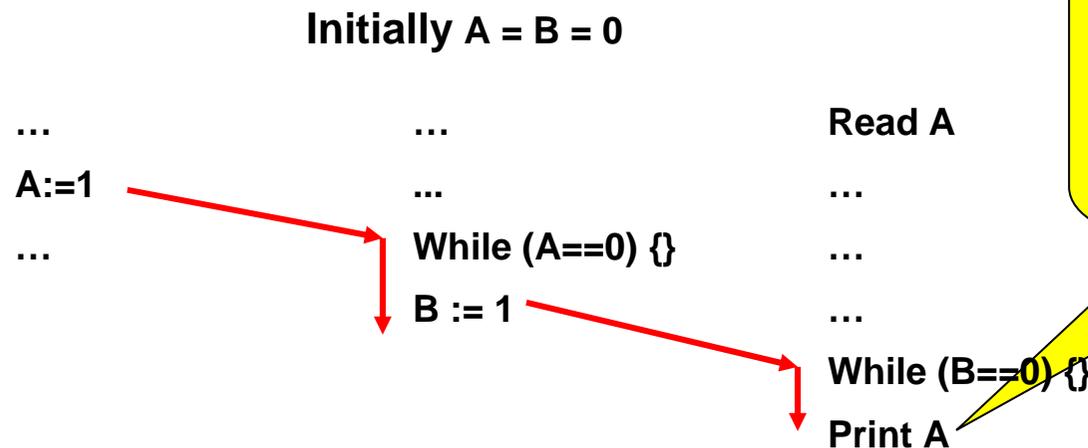Erik Hagersten

Uppsala University

Sweden

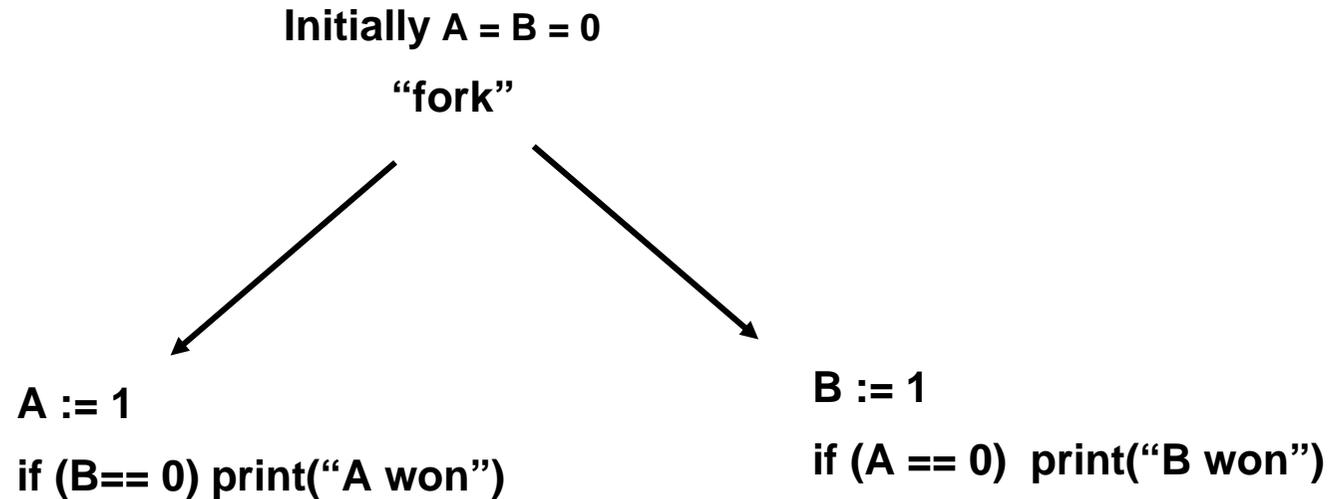# The Shared Memory Programming Model (Pthreads/OpenMP, ...)

# Memory Ordering

- Coherence defines a per-datum valuechange order

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

**Initially A = B = 0**

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

**Q:** What value will get printed?

# Dekker's Algorithm

Initially A = B = 0

"fork"

A := 1
if (B== 0) print("A won")

B := 1
if (A == 0)  print("B won")

## Q: Is it possible that both A and B win?

PDC
Summer
School
2010

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

MP 26
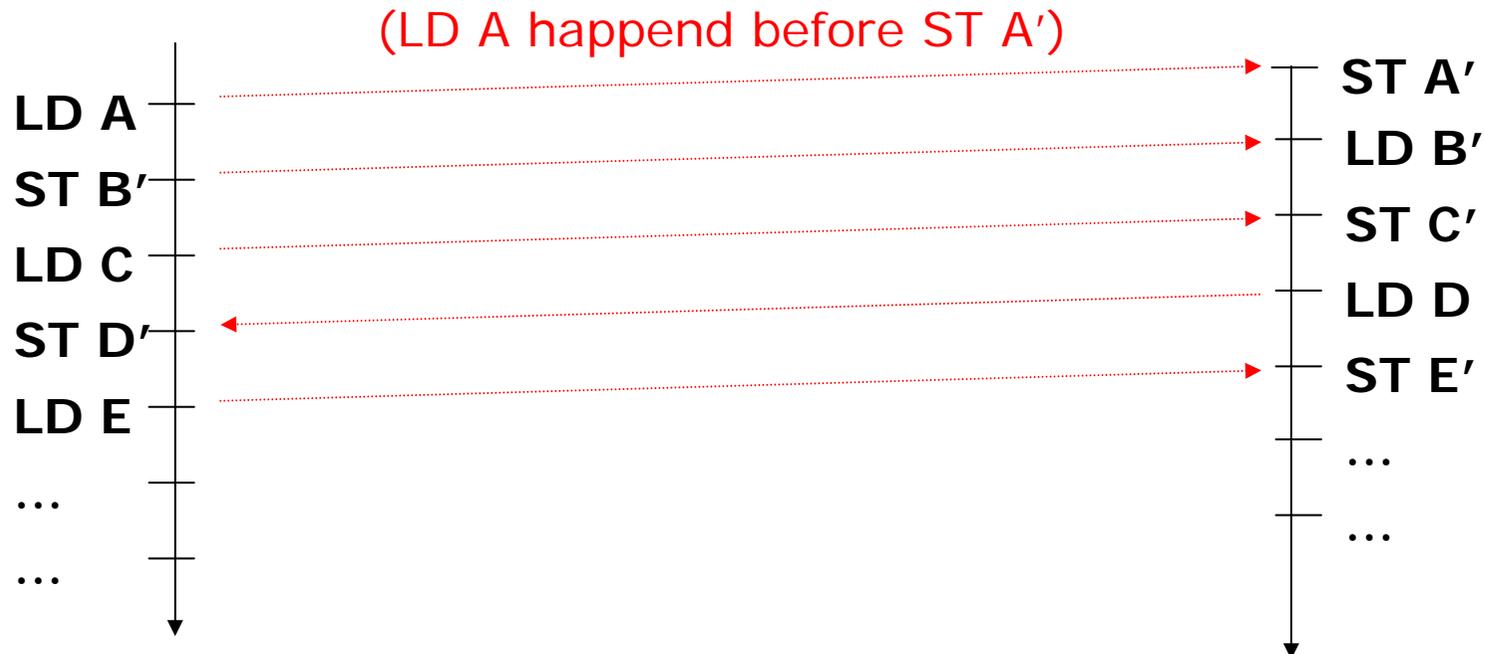
© Erik Hagersten| user.it.uu.se/~eh

# Memory Ordering

- Defines the guaranteed memory ordering: *If a thread has seen that A happens before B, what order can the other threads can observe?*

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

- Without it, you can not say much about the result of a parallel execution

MP 27

# Human intuition: There is one global order!

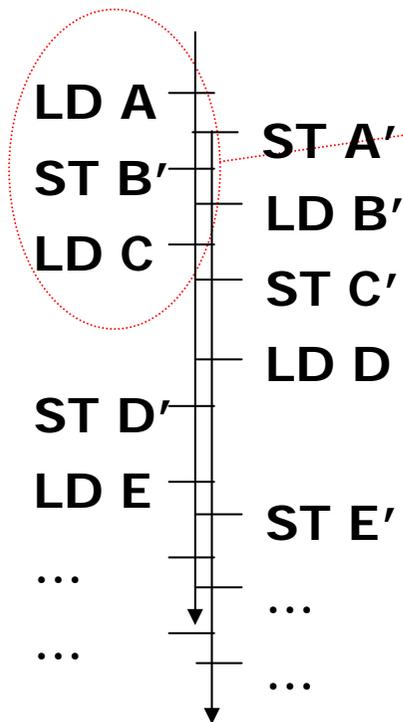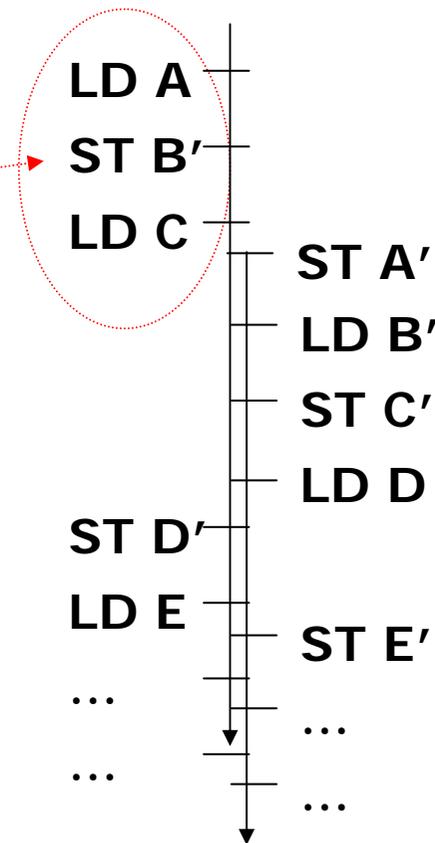**( A′ denotes a modified value to the data at addr A)**

**Thread 1**

**Thread 2**

(LD A happend before ST A′)

LD A

ST B′

LD C

ST D′

LD E

…

…

ST A′

LD B′

ST C′

LD D

ST E′

…

…

# One possible observed order

# Another possible observed order

**Thread 1 Thread 2**

Thread 1:
LD A
ST B'
LD C
ST D'
LD E
...
...

Thread 2:
ST A'
LD B'
ST C'
LD D
ST E'
...
...

**Thread 1 Thread 2**

Thread 1:
LD A
ST B'
LD C
ST D'
LD E
...
...

Thread 2:
ST A'
LD B'
ST C'
LD D
ST E'
...
...

PDC
Summer
School
2010

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

**Shared Memory**

loads, stores

T T T T T T T T T<sub>hread</sub>

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

* "Programmer's intuition is maintained"
  * Store causality? Yes
  * Does Dekker work? Yes

* Unnecessarily restrictive ==> performance penalty

PDC Summer School 2010

# Dekker's Algorithm

Initially A = B = 0

"fork"

A := 1
if (B== 0) print("A won")
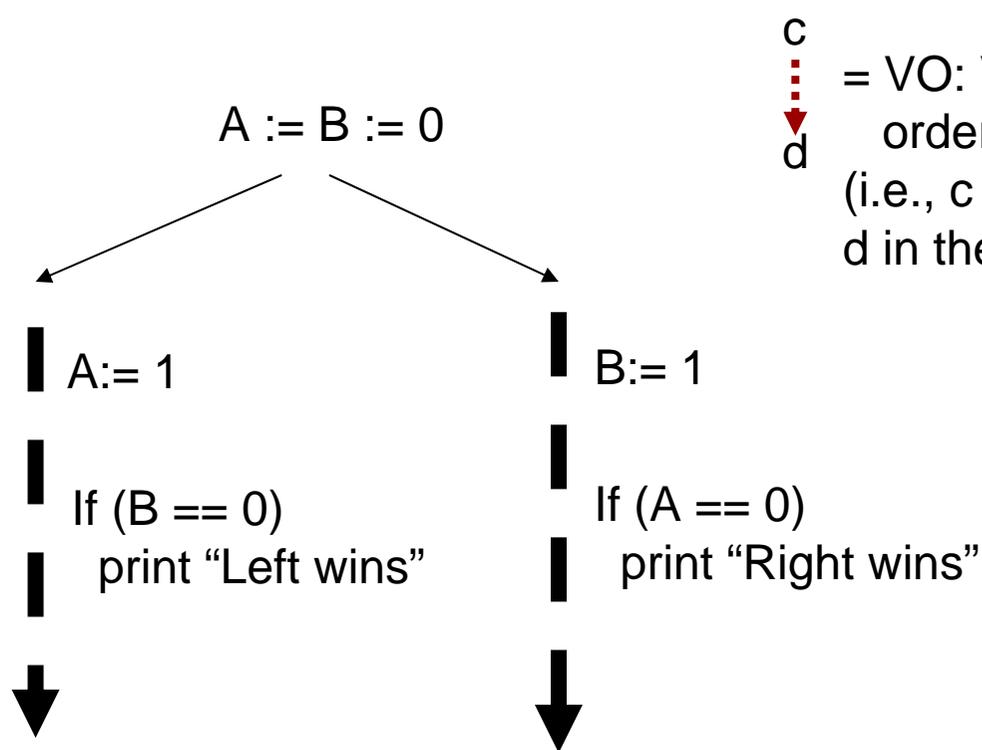
B := 1
if (A == 0)  print("B won")
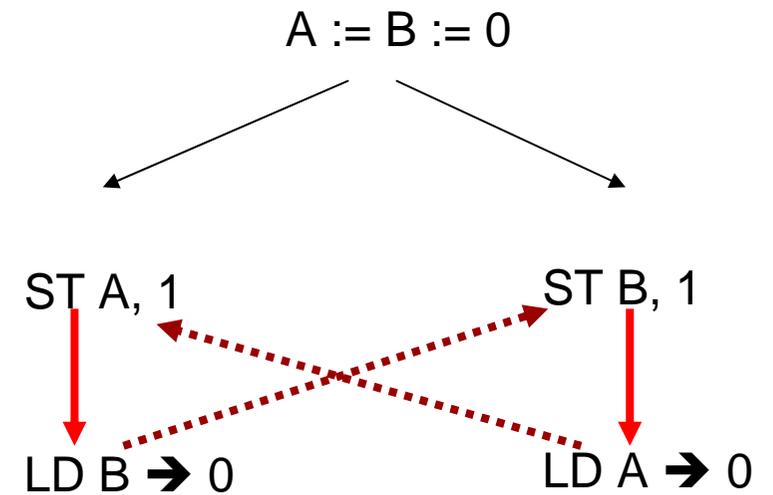
## Q: Is it possible that both A and B win?

MP 31

# Sequential Consistency (SC) Violation → Dekker: both wins

## Acess graph

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

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

A := B := 0

A:= 1

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

B:= 1

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

Both Left and Right wins →
SC violation

A := B := 0

ST A, 1                    ST B, 1

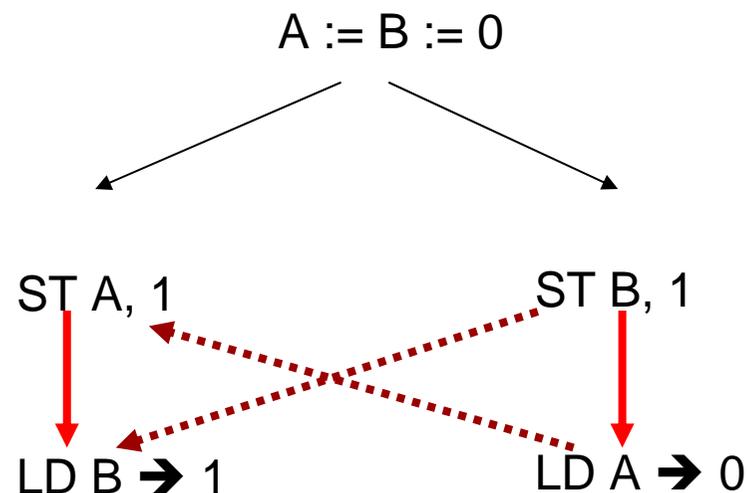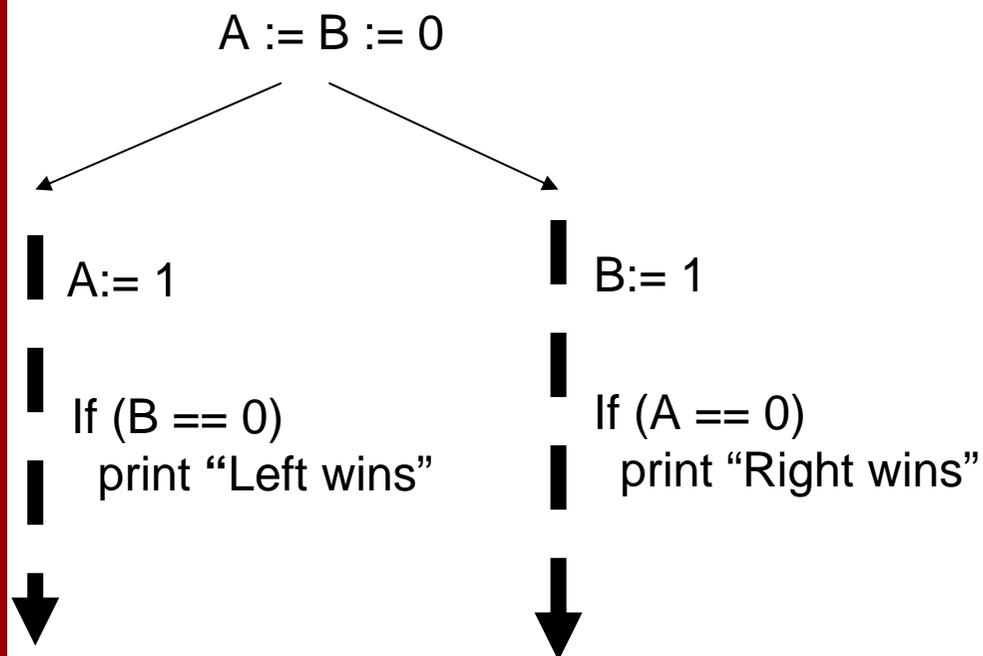LD B → 0                   LD A → 0

**Cyclic access graph → Not SC
(there is no global order)**

# SC is OK if one thread wins

Only Right wins ➔ SC is OK

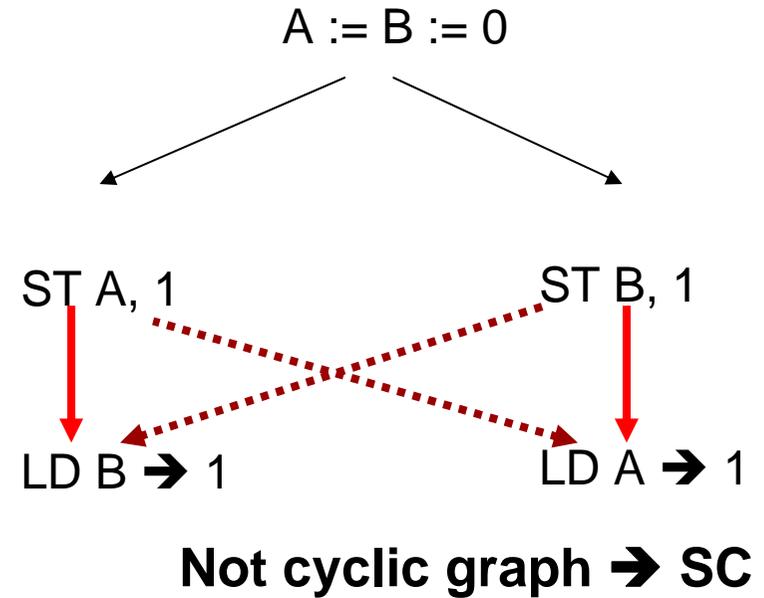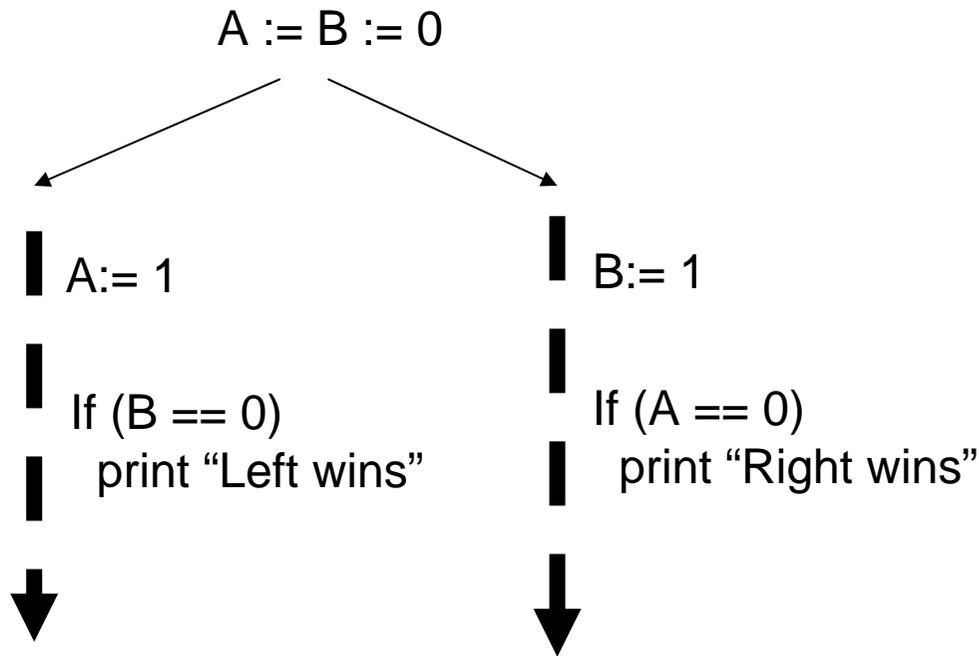A := B := 0

A:= 1

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

A := B := 0

ST A, 1

LD B ➔ 1

B:= 1

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

ST B, 1

LD A ➔ 0

**Not cyclic graph ➔ SC**

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

PDC
Summer
School
2010

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

MP 33

© Erik Hagersten| user.it.uu.se/~eh

# SC is OK if no thread wins

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**

UPPSALA
UNIVERSITET

# One implementation of SC in dir-based
## (....without speculation)

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

© Erik Hagersten| user.it.uu.se/~eh

PDC
Summer
School
2010

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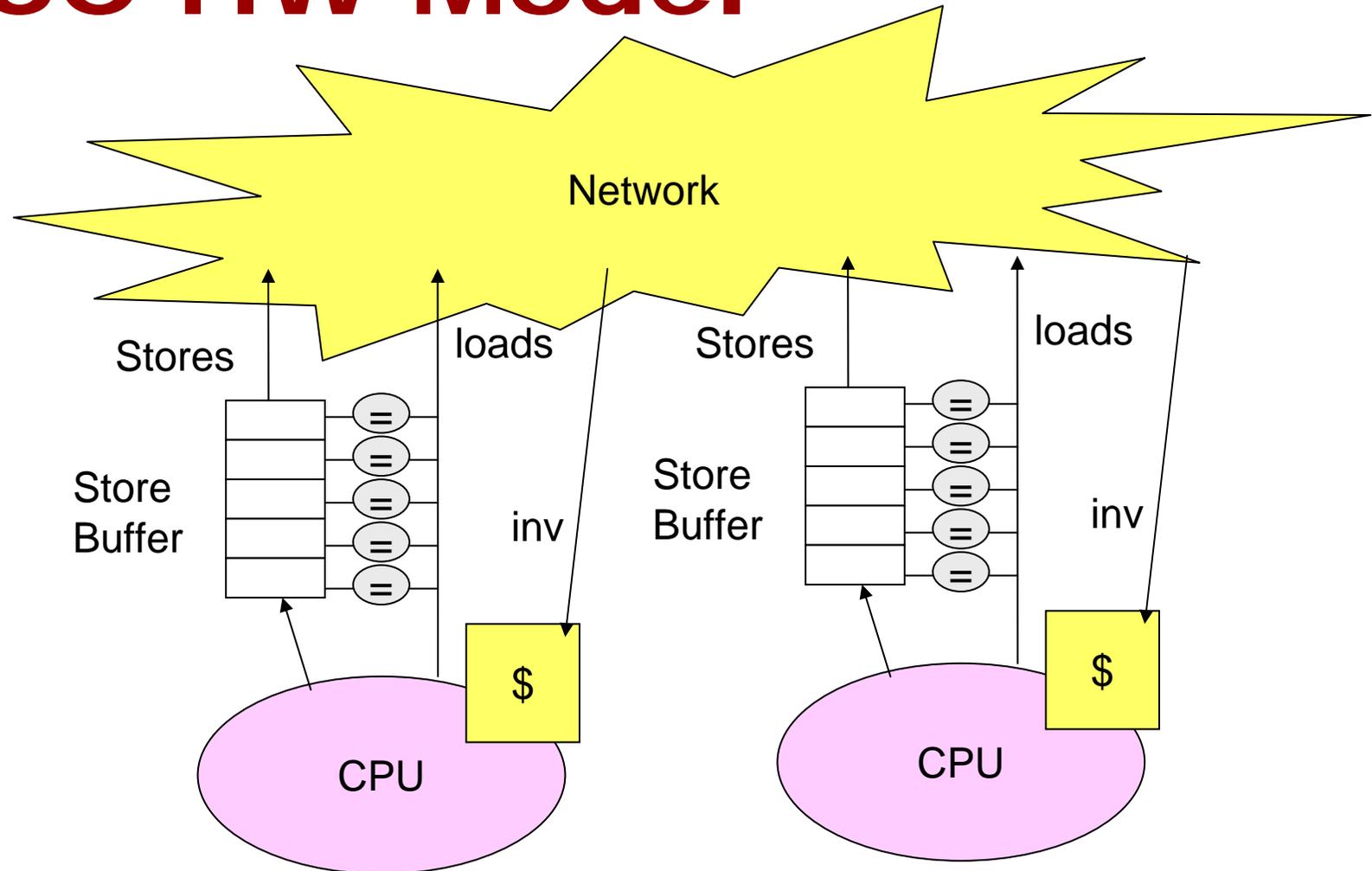


* Global *interleaving* [order] for <u>all</u> stores from different threads (own stores excepted)
* "Programmer's intuition is maintained"
  * Store causality? Yes
  * Does Dekker work? No
* Unnecessarily restrictive ==> performance penalty

# TSO HW Model



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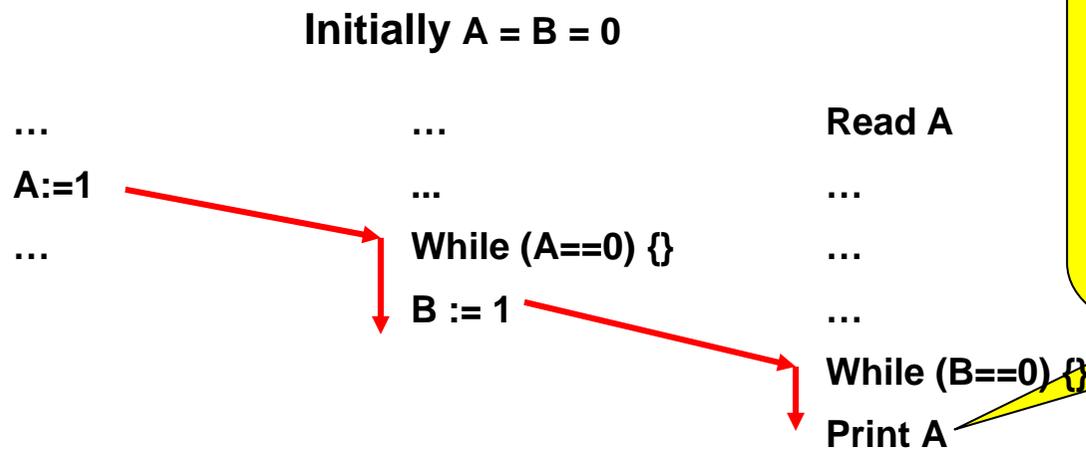Coherence implementation can be the same as for SC

MP 37

# TSO

- ## Flag synchronization works

  A := data                    while (flag != 1)  {};

  flag := 1                    X := A

- ## Provides causal correctness

  **Initially A = B = 0**

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

  **Q:** What value will get printed? Answer: 1

# Dekker's Algorithm, TSO

**Initially A = B = 0**

**"fork"**

**A := 1**

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

**B := 1**

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

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

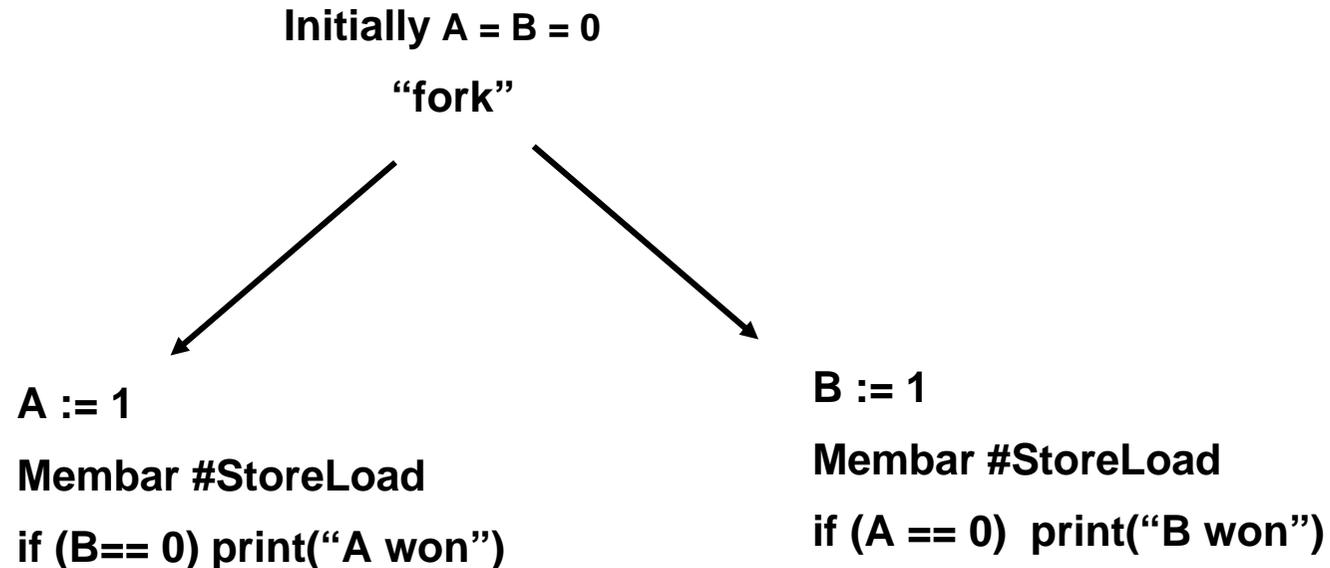## Q: Is it possible that both A and B wins?

**Left: The read (i.e., test if B==0) can bypass the store (A:=1)**
**Right: The read (i.e., test if A==0) can bypass the store (B:=1)**
**➔both loads can be performed before any of the stores**
**➔yes, it is possible that both wins**
**➔➔ Dekker's algorithm breaks**

# Dekker's Algorithm for TSO

Initially A = B = 0

"fork"

A := 1

Membar #StoreLoad

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

B := 1

Membar #StoreLoad
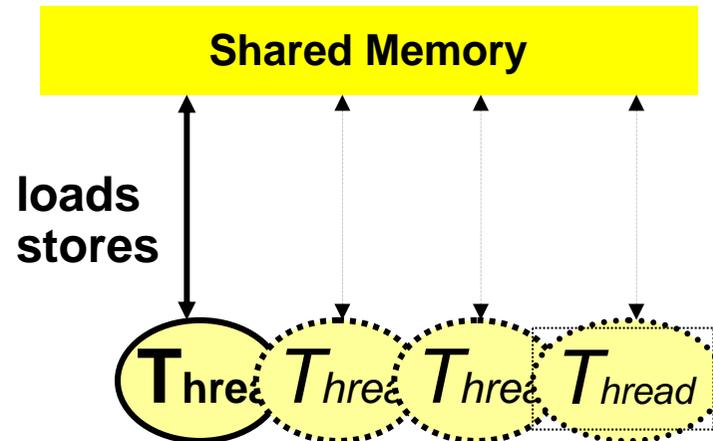
if (A == 0)  print("B won")

## Q: Is it possible that both A and B win?

Membar: The read is stared after all previous stores have been "globaly ordered"
➔ behaves like SC
➔ Dekker's algorithm works!

PDC
Summer
School
2010

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

MP 40

© Erik Hagersten| user.it.uu.se/~eh

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



* **Most accesses are unordered**
* **"Programmer's intuition is not maintained"**
  - Store causality? No
  - Does Dekker work? No
* **Global order <u>only</u> established when the programmer explicitly inserts memory barrier instructions**

++ Better performance!!
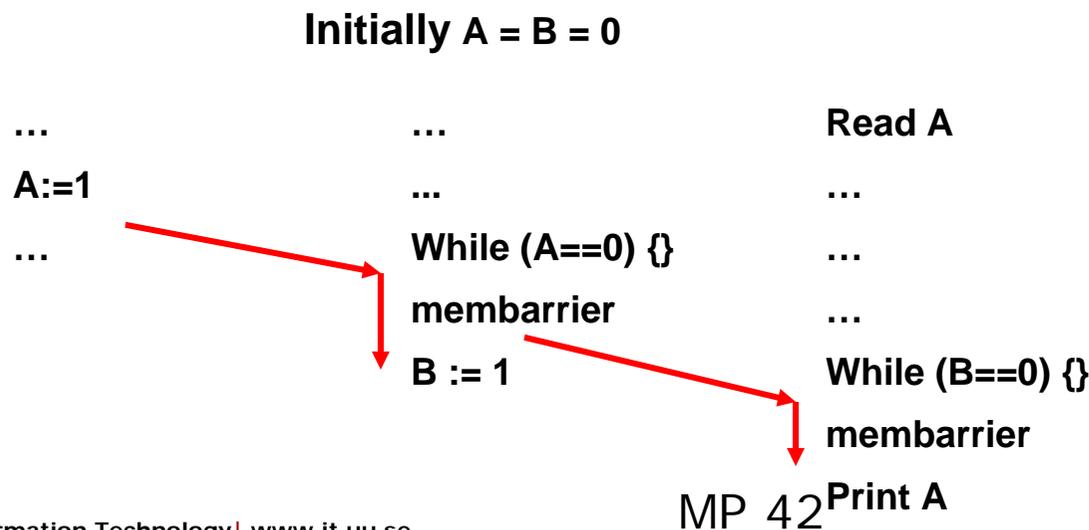
--- Interesting bugs!!

MP 41

# 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** |

MP 42

**Q:** What value will get printed? **Answer**: 1

# Learning more about memory models

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

RFM: 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?
Have to use "volatile" declarations in C.

PDC
Summer
School
2010

UPPSALA
UNIVERSITET

# X86's new memory model

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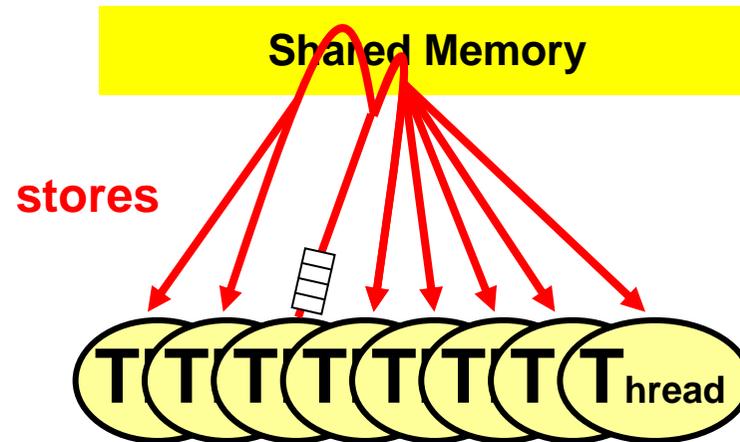- TSO for atomic memory ops

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

- See section 8.2 in this manual:
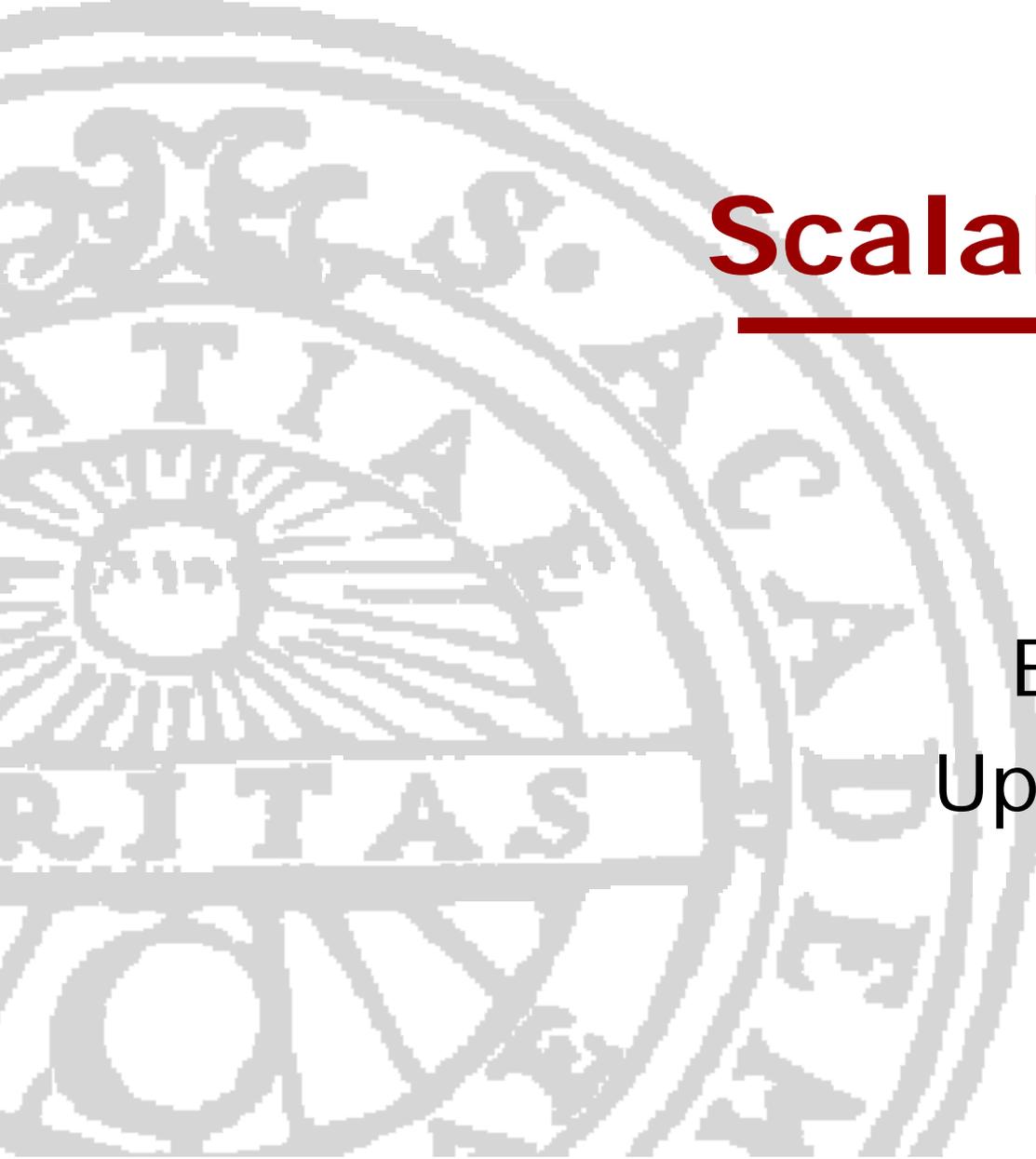  http://developer.intel.com/Assets/PDF/manual/253668.pdf

MP 44

# Processor Consistency [PC] (J. Goodman)



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

* Causal correctness (often added to PC): if a processor observes a store before performing a new store, the observed store must be observed before the new store by all processors

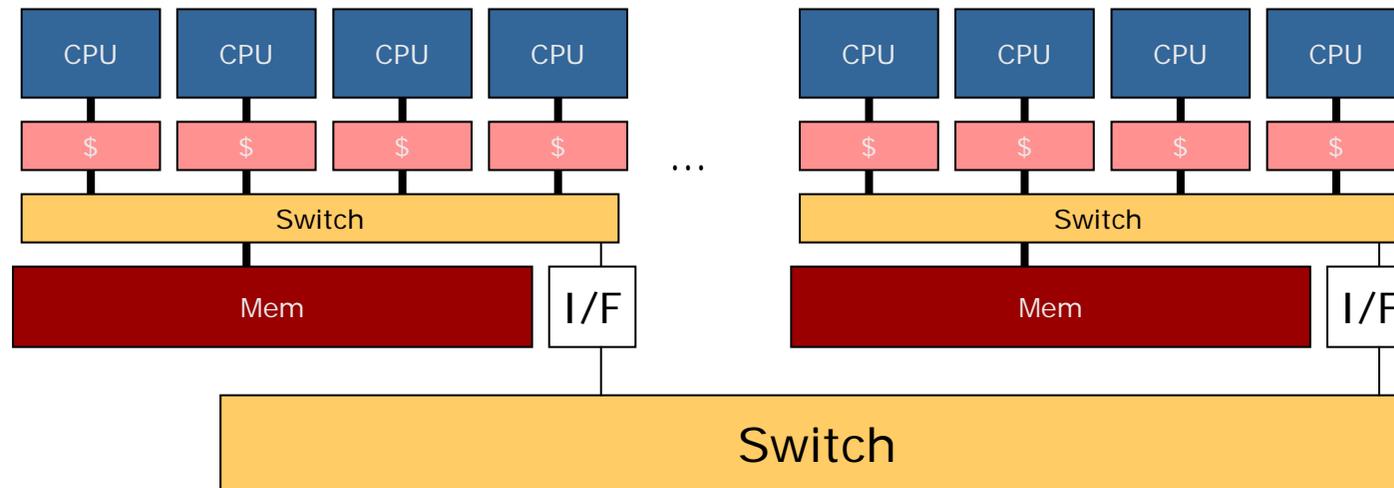➔ Flag synchronization works.
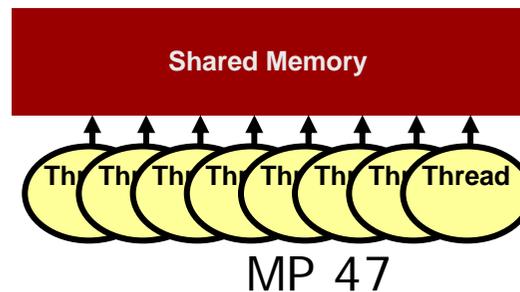
➔ No causal correctness issues

MP 45

# Scalable Systems

Erik Hagersten

Uppsala University

# NUMA:
# Non-uniform memory architecture

| CPU | CPU | CPU | CPU | | CPU | CPU | CPU | CPU |
|-----|-----|-----|-----|---|-----|-----|-----|-----|
| $ | $ | $ | $ | ... | $ | $ | $ | $ |

| Switch | | Switch |
|--------|---|--------|

| Mem | I/F | | Mem | I/F |
|-----|-----|---|-----|-----|

| Switch |
|--------|

Same SW view:

**Shared Memory**

Thr Thr Thr Thr Thr Thr Thr **Thread**

MP 47

# Non-uniform Architectures: NUMA
# ➔ Communication cost is much worse!

**Mem**

**Mem**

Interconnect

**$**

**$**

**~300-500ns**

**Thread**

**Thread**

**Read A**

**Read A**

**...**

**...**

**Read A**

**...**

**...**

**Write A**

**Read A**

PDC
Summer
School
2010

UPPSALA
UNIVERSITET

# Example
## Efficient us of coherent cache
## Synchronization

Erik Hagersten

Uppsala University

# The programmer's view



Shared Memory

Thr Thr Thr Thr Thr Thr Thr Thread

PDC
Summer
School
2010

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

MP 50

© Erik Hagersten| user.it.uu.se/~eh

# 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

**How to implement Lock/Unlock efficiently with coherent cache?**
1. **Read-replication is cheap;**
2. **Coherence traffic is expensive**

# A Bad Example: "POUNDING"

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


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

- *Assume: The <u>atomic</u> function TAS (test and set)*
  *returns the current memory value and atomically writes "1" to the memory location.*
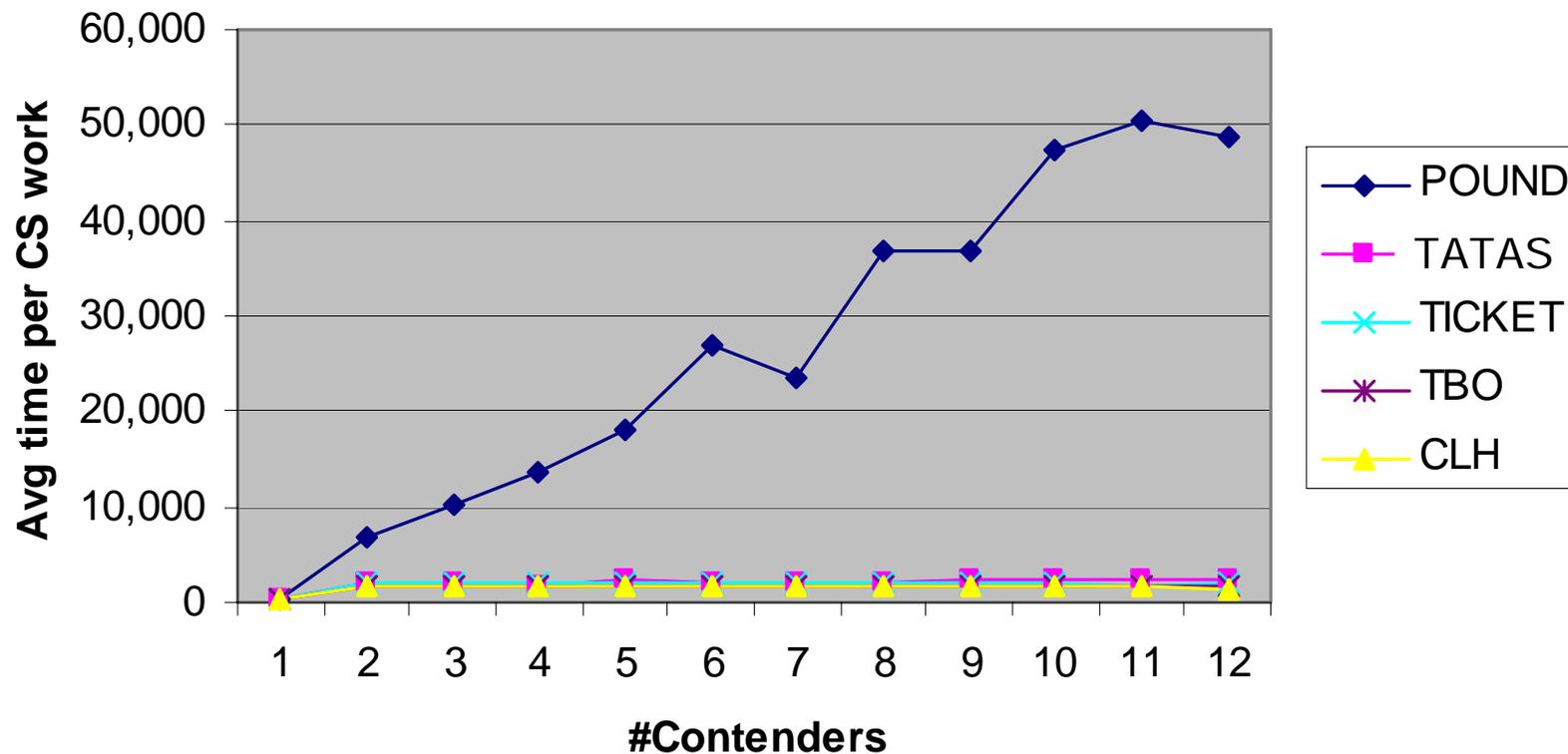
MP 52

# Optimistic Test&Set Lock "spinlock" (called "TATAS", test and test-and-set)

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

  }
}


proc unlock(lock_variable) {
    lock_variable := 0
}
```
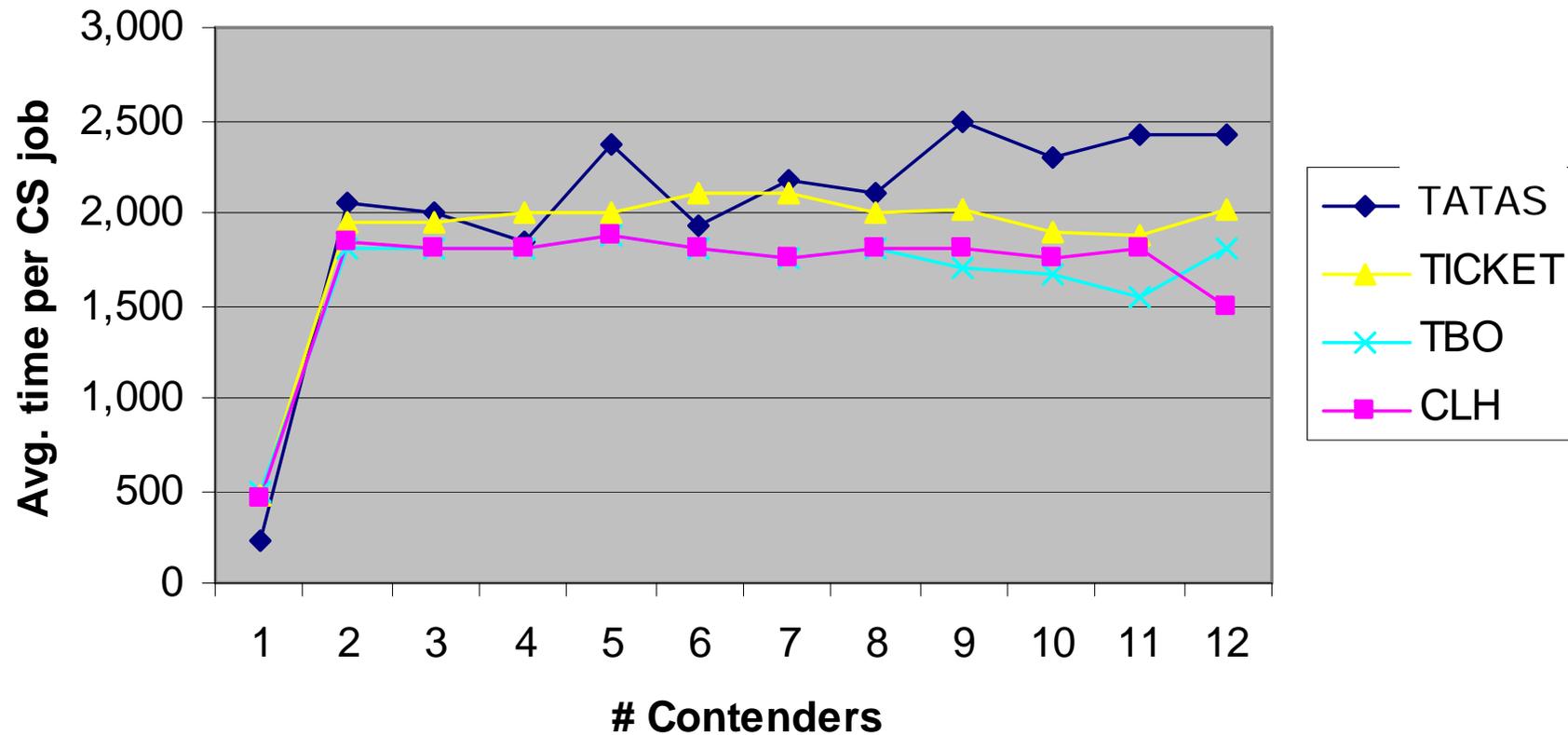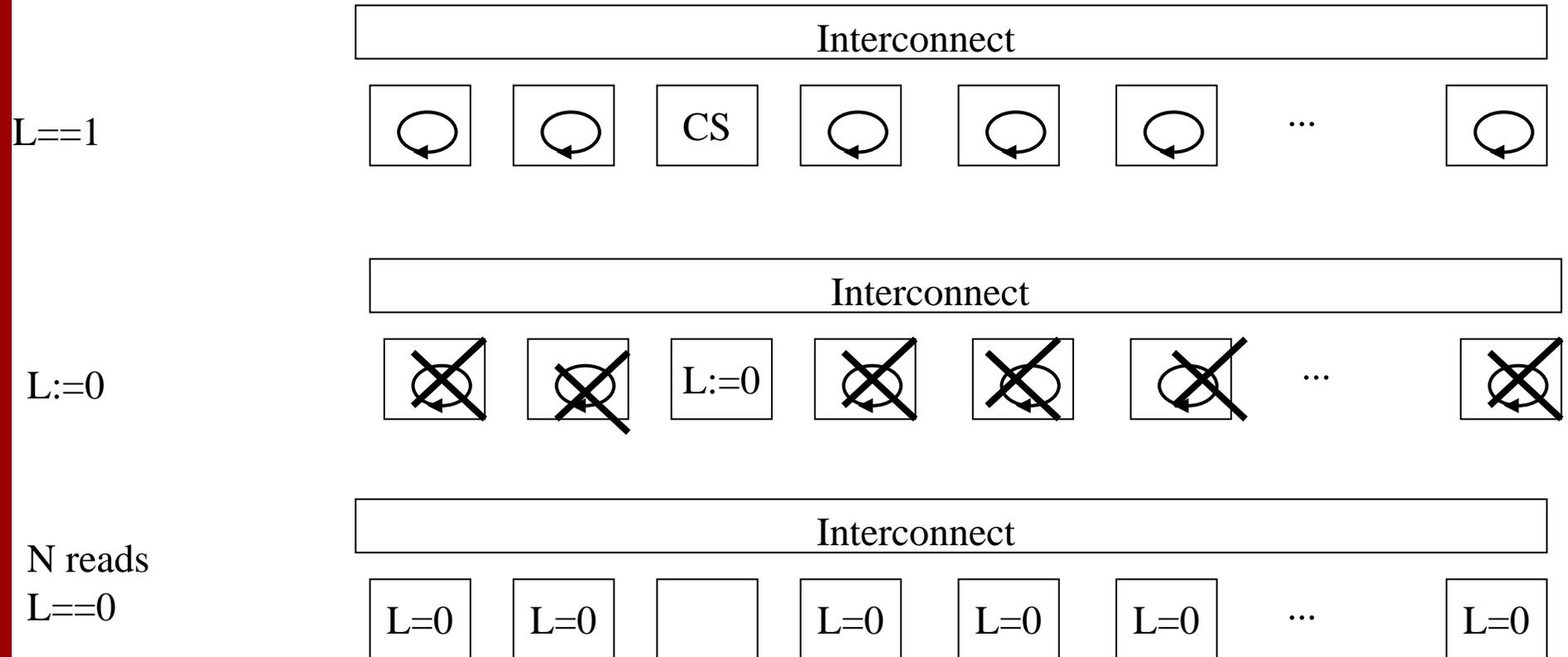
MP 53

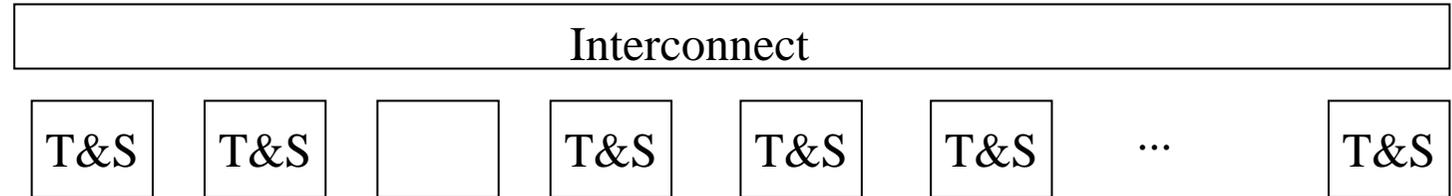# Uniform memory architecture: E6800 12 CPUs

# E6800 locks (exluding POUND)

MP 55

# TATAS could still get messy!

Interconnect

L==1

| $\circlearrowleft$ | $\circlearrowleft$ | CS | $\circlearrowleft$ | $\circlearrowleft$ | $\circlearrowleft$ | ... | $\circlearrowleft$ |

Interconnect

L:=0

| ⊠ | ⊠ | L:=0 | ⊠ | ⊠ | ⊠ | ... | ⊠ |

N reads
L==0

Interconnect

| L=0 | L=0 | | L=0 | L=0 | L=0 | ... | L=0 |

# ...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 proc
Problem2: The lock hand-over time is N*read_throughput
Fix1: some back-off strategy, bad news for hand-over latency
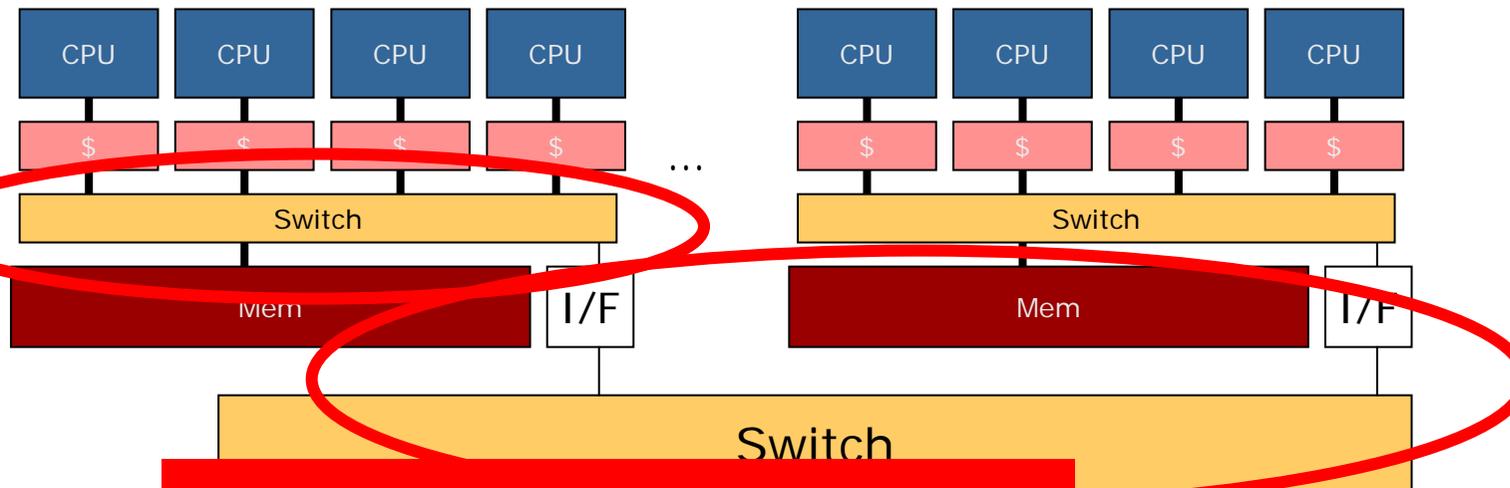Fix2: Queue-based locks
Fix3: Generate handover locality

# NUMA:

# NUCA: Non-uniform Comm Arch.

**WF**

**Snoop**

| CPU | CPU | CPU | CPU | ... | CPU | CPU | CPU | CPU |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| $ | $ | $ | $ | | $ | $ | $ | $ |

Switch                          Switch

Mem          I/F          Mem          I/F

Switch

**Directory-latency = 6x snoop i.e., roughly CMP NUCA-ness**

# Queue-based Locks



NUMA

- Inserts the contenders in a FIFO order
- Gauge the completion of your predecessor [only] while waiting

➔ Fairness
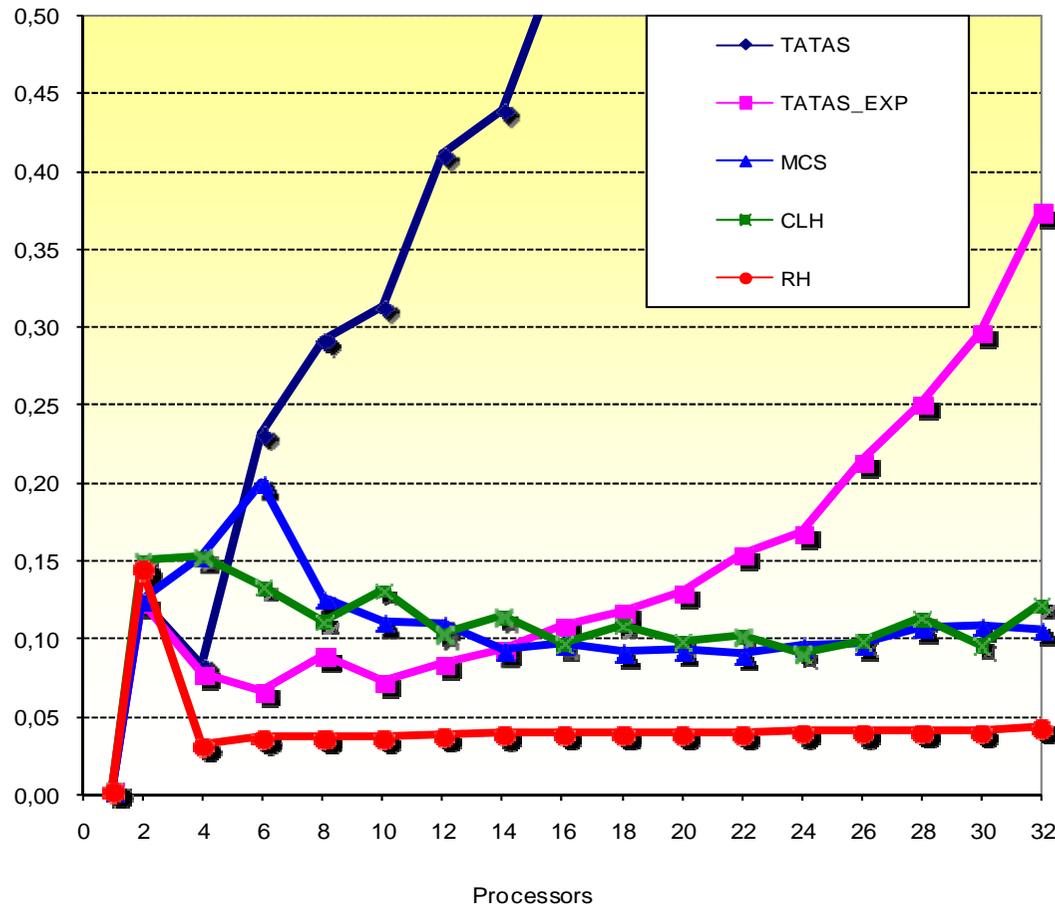
➔ Reduced Traffic

Examples: MCS locks and CLH locks

PDC
Summer
School
2010

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

MP 59

© Erik Hagersten| user.it.uu.se/~eh

# UART research: RH locks
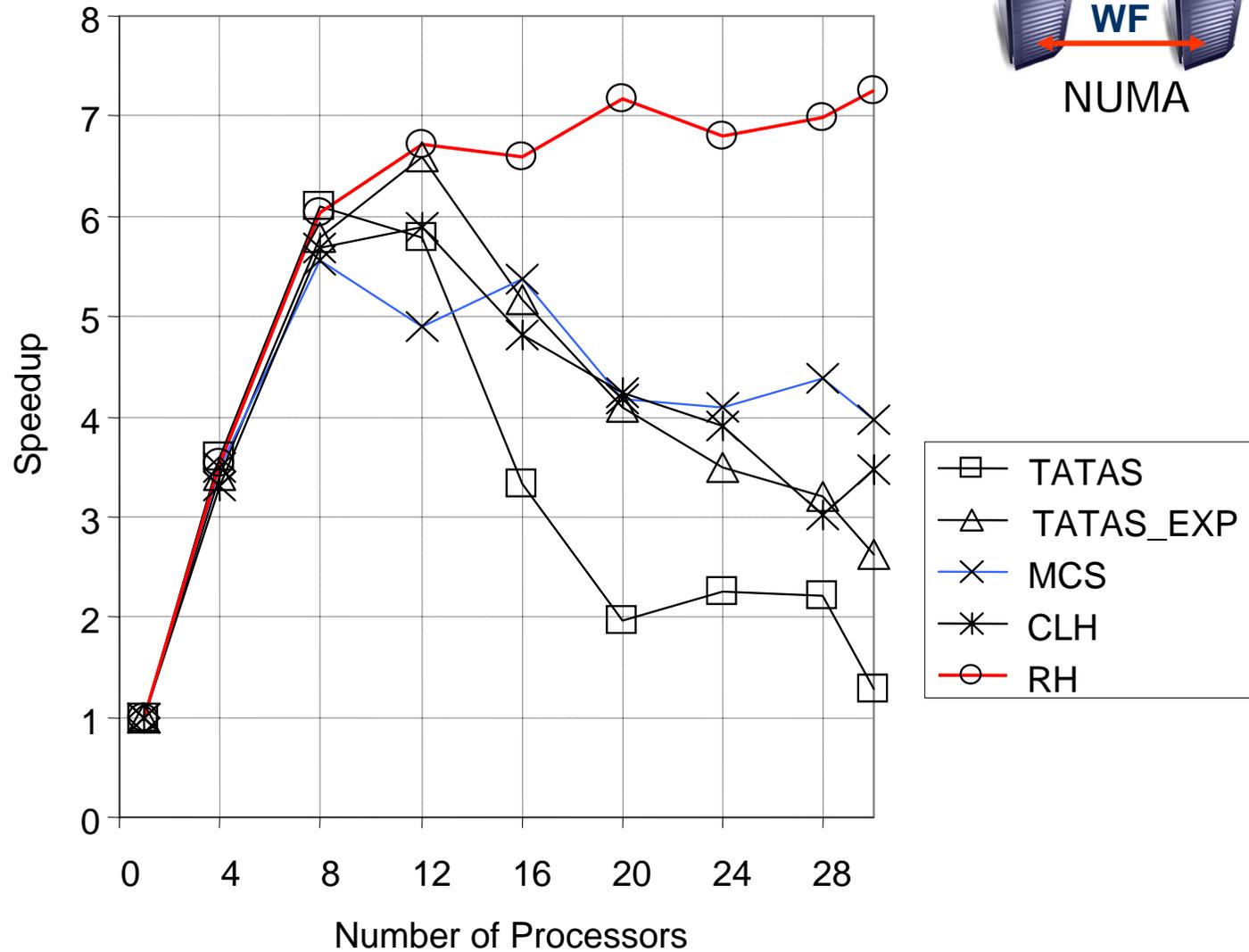
**Thanks: Zoran Radovic**

NUMA

Lock
Hand-over
Time

# Application performance

## Ex: Raytrace Speedup



NUMA

**WF**

Speedup vs. Number of Processors

Legend:
- □ TATAS
- △ TATAS_EXP
- ✕ MCS
- ✱ CLH
- ○ RH

PDC
Summer
School
2010