# Parallel Programming

## Part 2:
## Shared Memory Parallel Architectures

**foils by C. Kessler, T. Fahringer**

1

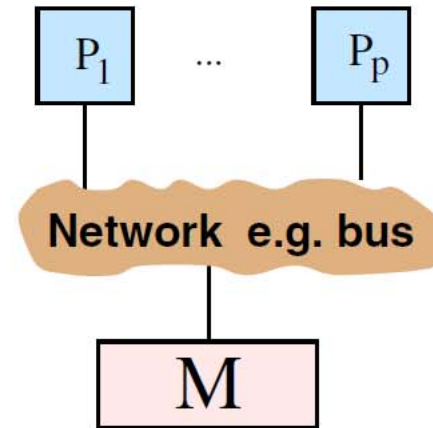# Outline

- Memory hierarchy

- Consistency issues and coherence protocols

- Performance issues, e.g. false sharing

- Optimizations for data locality

# Distributed Memory vs. Shared Memory
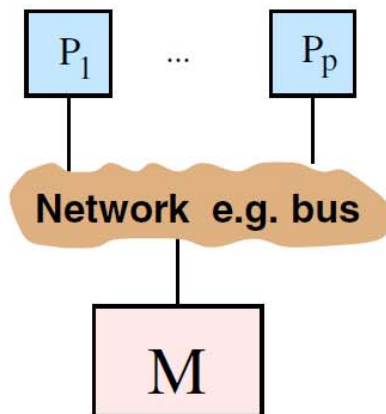


Distributed memory system
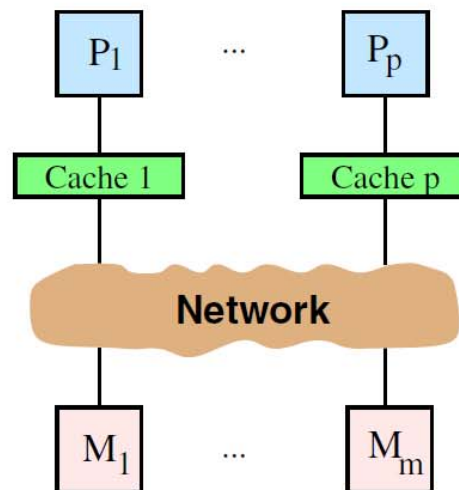
Shared memory system

# Shared Memory – Variants

- Single shared memory module (UMA) with identical costs for accessing a memory location quickly becomes a performance bottleneck.

- Often implemented with caches to leverage access locality

- Can even be realized on top of physically distributed memory system (NUMA – non-uniform memory access)

**SMP (Symmetric Multiprocessing)**
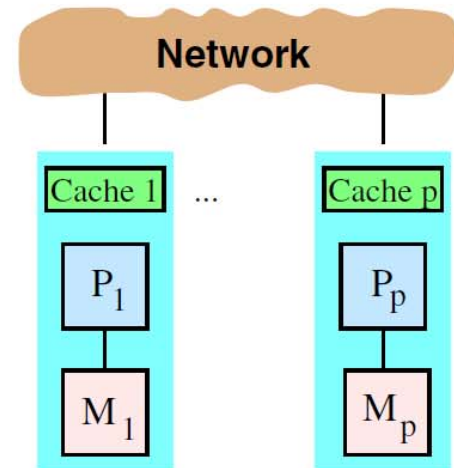
$P_1$ ... $P_p$

Network e.g. bus

$M$

**SMP with caches**

$P_1$ ... $P_p$

Cache 1    Cache p

Network

$M_1$ ... $M_m$

**Distributed shared memory (DSM / NUMA)**

Network

Cache 1 ... Cache p

$P_1$    $P_p$

$M_1$    $M_p$

# Cache

**Cache** = <u>small</u>, <u>fast</u> memory (SRAM) between processor
and main memory, today typically on-chip

- contains copies of main memory words

  - cache hit = accessed word already in cache, get it fast.

  - cache miss = not in cache, load from main memory (slower)

- **Cache line** holds a copy of a **block** of adjacent memory words

  - size: from 16 bytes upwards, can differ for cache levels

  - currently, typically 64 bytes (8 words 64 bit each)

In the following, we mainly refer to **data cache**

  - Other cache types in a processor: instruction cache, TLB

# Cache (cont.)

Mapping memory blocks → cache lines:

- direct mapped: $\forall j\, \exists! i:\ B_j \mapsto C_i,$

- fully-associative: any memory address may be placed in any cache line

- n-way set-associative

  – caches with n cache lines per set.

  – every memory address maps to a specific set but can be placed in any of the n cache lines of this set.

– modern processors use direct-mapped, 2-way set-associative or 4-way set associative caches.

# Cache (cont.)

- Cache-based systems profit from

  - spatial access locality

    - access also other data in same cache line

  - temporal access locality

    - access same words in cache line multiple times

- HW-controlled cache line replacement (which cache line) (for fully and set-associative caches)

  - E.g., LRU – Least Recently Used (usually, approximations) → dynamic adaptivity of cache contents

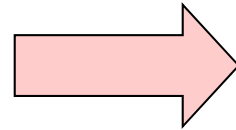- Suitable for applications with high (also dynamic) data locality

# Classification of Cache Misses

- **Cold miss** – data was never in cache for this process
  (first access in program execution)

  – Usually unavoidable

- **Capacity miss** – data was evicted earlier
  when cache capacity was exceeded
  (for fully associative caches)

  – Might be (partly) avoided by redesigning the programs'
  algorithm to improve access locality.

- **Associativity miss** – data was evicted earlier because of
  cache set conflict (for set-associative
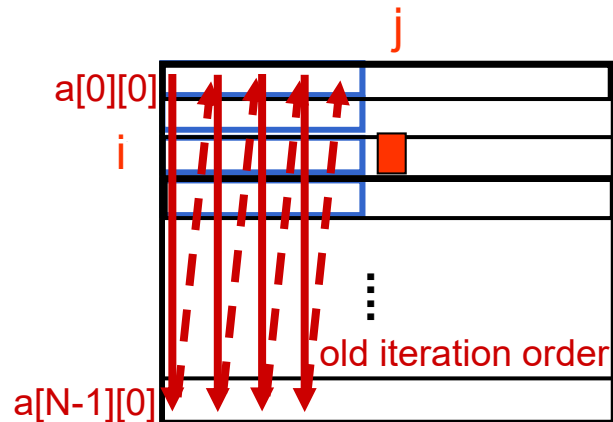  and direct-mapped caches)

# Optimizing Programs for Improved Access Locality

- **Example:**
  Loop Interchange



```
for (j=0;  j<M;  j++)
  for (i=0;  i<N;  i++)
    a[ i ][ j ] = 0.0 ;
```

row-wise storage of 2D-arrays in C, Java

```
for (i=0;  i<N;  i++)
  for (j=0;  j<M;  j++)
    a[ i ][ j ] = 0.0 ;
```

a[0][0]

i

a[N-1][0]

old iteration order

a[0][0]

i

a[0][M-1]

new iteration order
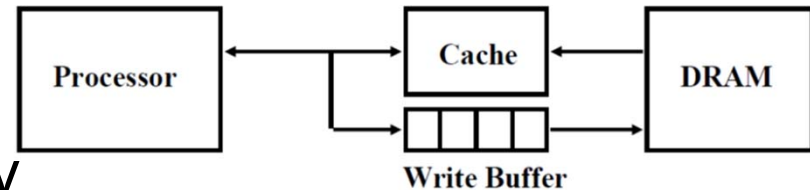
a[N-1][0]

- Can improve spatial locality of memory accesses (fewer cache misses).

12

# Caches: Memory Update Strategies

Write-through



   + write both cache and memory

   + consistency (data identical in cache and memory )

   &minus;  slow, write stall       ($\rightarrow$ write buffer)

Write-back

   + update only cache

   + write back to memory only when replacing cache line

   + write only if modified, marked by „dirty" bit for each $C_i$

   &minus;  not consistent,

      DMA access (I/O, other procs) may access stale values

      $\rightarrow$ write back on request (flush) or if cache line is evicted from cache.

# Memory Hierarchy Example

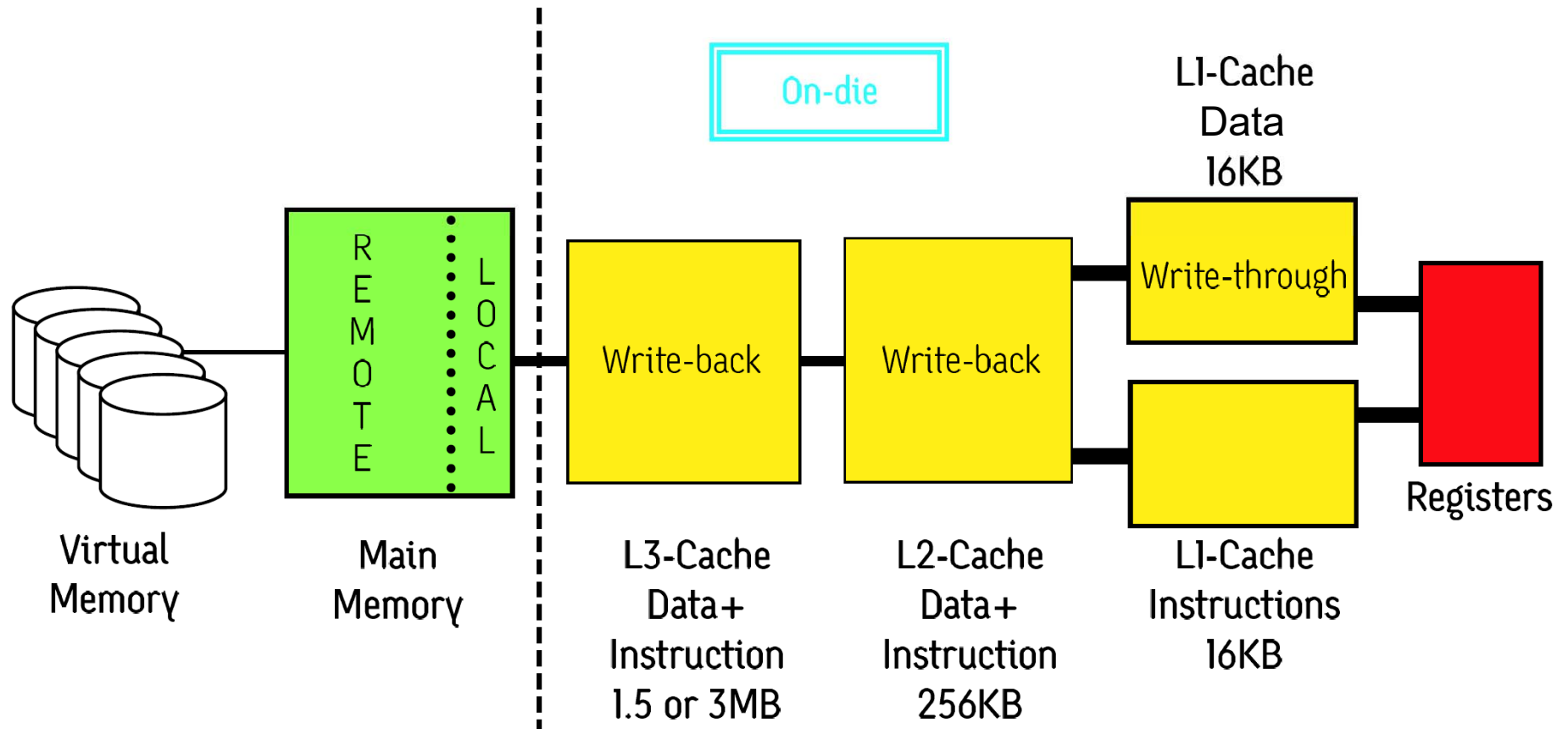Memory hierarchy of the Itanium2 processors in SGI Altix 3700 Bx2



Fig. 2. Itanium 2 memory and cache hierarchy diagram

# Cache Coherence and Memory Consistency

Caching of (shared) variables leads to consistency problems.

A cache management system is called coherent

> if a read access to a (shared) memory location x reproduces always the value corresponding to the *most recent write access* to x. Cache coherence deals with ordering of writes to a single memory location.

> → no access to stale values

A memory system is consistent (at a certain time)

> if all copies of shared variables in the main memory and in the caches are identical. Deals with ordering of reads/writes to all memory locations.

Permanent memory-consistency implies cache-coherence.

# Cache Coherence – Formal Definition

What does „most recent write access" to $x$ mean?

Formally, 3 conditions must be fulfilled for coherence:

(a) Each processor sees its *own* writes and reads in program order.

*P*1 writes $v$ to $x$ at time $t1$, P1 reads from $x$ at $t2 > t1$, no other processor writes to $x$ between $t1$ and $t2$ → read yields $v$

(b) the written value is eventually visible to *all* processors. *P*1 writes $v$ to $x$ at $t1$, *P*2 reads from $x$ at $t2 > t1$, no other processor writes to $x$ between $t1$ and $t2$, and $t2 > t1$ sufficiently large, then *P*2 reads $v$.

(c) All processors see the same total order of all write accesses. *(total store ordering)*

# Cache Coherence Protocols

Inconsistencies occur when modifying only the copy of a shared variable in a cache, not in the main memory and all other caches where it is held.

Write-update protocol (word basis)

> At a write access, all other copies in the system must be updated as well. Updating must be finished before the next access.

Write-invalidate protocol (cache line basis)

> Before modifying a copy in a cache, all other copies in the system must be declared as „invalid".

Most cache-based SMPs use a write-invalidate protocol.

Updating / invalidating straightforward in bus-based systems (bus-snooping), otherwise, a directory mechanism is necessary

# Details:  Write-Invalidate Protocol

Implementation: multiple-reader-single-writer sharing

At any time, a data item (usually, entire cache line) may either be:

accessed in read-only mode by one or more processors

read and written (exclusive mode) by a single processor

Items in read-only mode can be copied indefinitely to other processors.

Write attempt to read-only mode data $x$:

1. Broadcast invalidation message to all other copies of $x$

2. Await acknowledgements before the write can take place

3. Any processor attempting to access $x$ is blocked if a writer exists.

4. Eventually, control is transferred from the writer and other accesses may take place once the update has been sent

$\rightarrow$ all accesses to x processed on first-come-first-serve basis.

Achieves sequential consistency.

# Write-Invalidate Protocol (cont.)

+ parallelism (multiple readers)

+ updates propagated only when data are read

    + several updates by the same processor can take place before communication is necessary

– Cost of invalidating read-only copies before a write can occur

    + ok if read/write ratio is sufficiently high

    + for small read/write ratio: single reader-single-writer scheme (at most one process gets read-only access at a time)

# Write-Update Protocol

Protocol is implemented on the basis of words.

Write $x$:

- done locally + broadcast new value to all who have a copy of $x$

  these update their copies immediately.

- several processors may write the same data item at the same time (multiple-reader-multiple-writer sharing); blocking writes.

- Sequential consistency if broadcasts are totally ordered and blocking

  $\rightarrow$ all processors agree on the order of updates.

  $\rightarrow$ the reads between writes are well defined

Read $x$:

  read local copy of $x$, no need for communication.
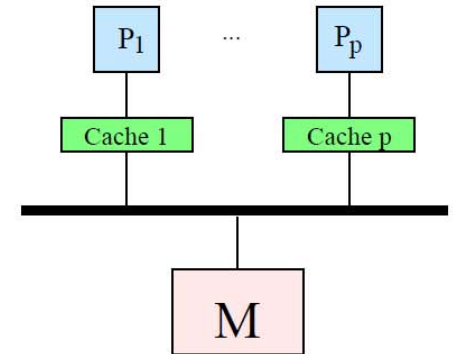
  $\rightarrow$ multiple readers

  + Reads are cheap

– totally ordered broadcast protocols quite expensive to implement

# Bus-Snooping

For bus-based SMP with caches and write-through strategy.

All relevant memory accesses go via the central bus.

Cache-controller of each processor listens to addresses on the bus:

    write access to main memory is recognized
    and committed to the own cache.

– bus is performance bottleneck $\rightarrow$ poor scalability

# Write-back invalidation protocol
# (MSI-protocol)

MSI protocol applied for write-back and write-invalidate caches and uses bus snooping technique.

A block held in cache has one of 3 states:

M (modified)

  only this cache entry is valid, all other copies + MM location are not.

S (shared)

  cached on one or more processors, all copies are valid.

I (invalid)

  this cache entry contains invalid values.

# MSI-Protocol: State Transitions

State transitions:

triggered by observed bus operations and local processor reads/writes

Bus read (BusRd)

Place read miss on bus
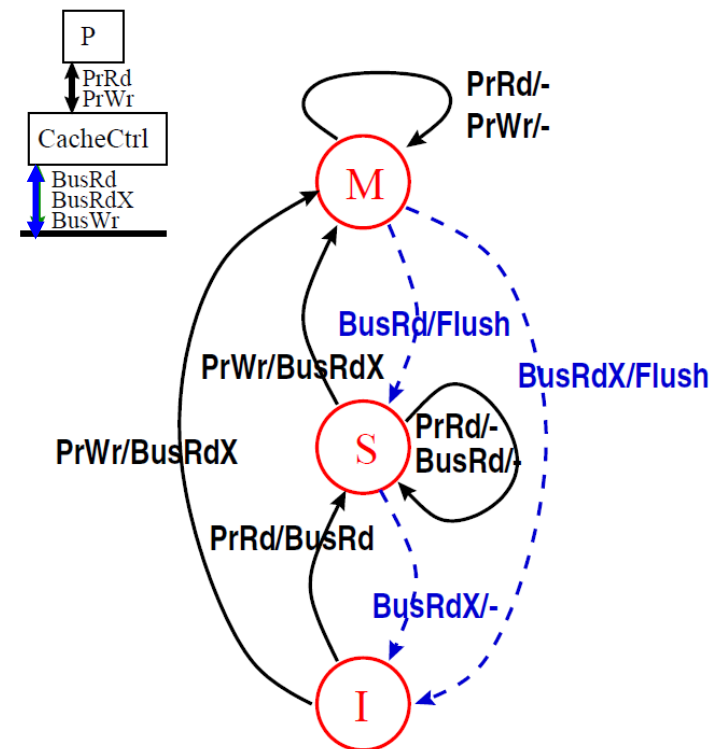
Bus read exclusive (BusRdX)

Place write miss on bus

→ must invalidate other copies

Write back (BusWr), due to replacement

→ Flush

Processor reads (PrRd)

Processor writes (PrWr)



P

PrRd
PrWr

CacheCtrl

BusRd
BusRdX
BusWr

M

PrRd/-
PrWr/-

BusRd/Flush

PrWr/BusRdX

BusRdX/Flush

PrRd/-
BusRd/-

S

PrWr/BusRdX

PrRd/BusRd

BusRdX/-

I

Processor operation / Cache controller operation

Observed operation / Cache controller operation

Flush = desired value put on the bus

Missing edges - no change of state

30

# Example  Initial situation…

## Core *i*

Cache lines / state

I

## Core *j*

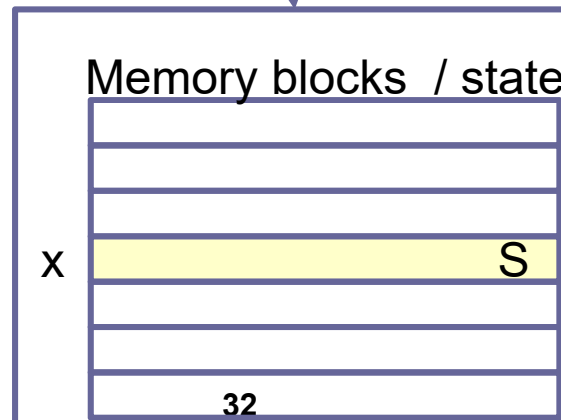Cache lines / state



Processor operation / Cache controller operation

Observed operation / Cache controller operation

## Memory blocks / state

x   S

32

# Core *i*

**Load x:**

Cache lines  /  state

| | |
|---|---|
| | |
| | |
| | |
| S | |
| | |
| | |
| | |

# Core *j*

Cache lines  / state

| |
|---|
| |
| |
| |
| |
| |
| |
| |

BusRd x

## Memory blocks  / state

| | |
|---|---|
| | |
| | |
| | |
| x | S |
| | |
| | |
| **33** | |

P

PrRd
PrWr

CacheCtrl

BusRd
BusRdX
BusWr

PrRd/-
PrWr/-

M

BusRd/Flush

PrWr/BusRdX

BusRdX/Flush

PrWr/BusRdX

S

PrRd/-
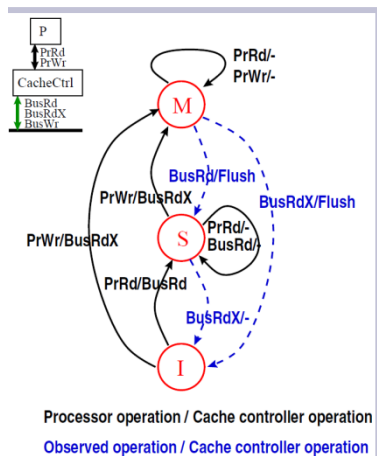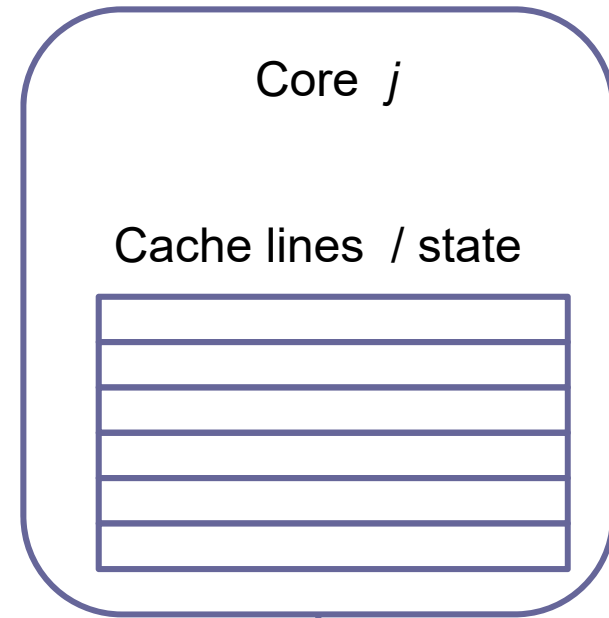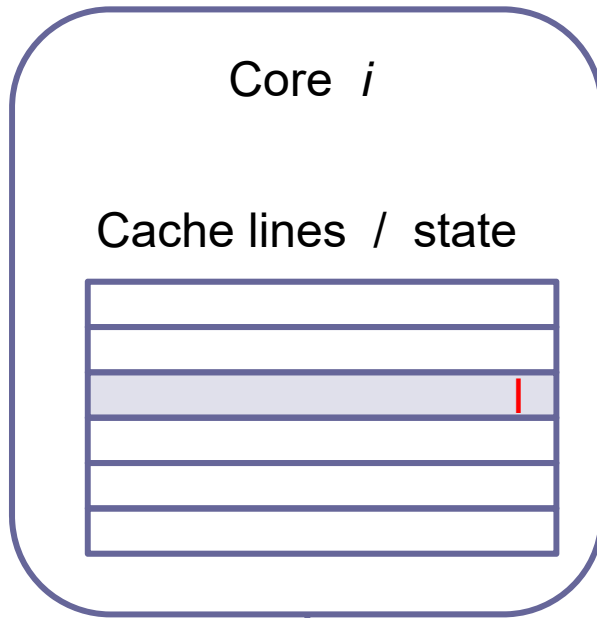BusRd/-

PrRd/BusRd

BusRdX/-

I

Processor operation / Cache controller operation

Observed operation / Cache controller operation

# Core *i*

## Cache lines / state

| | |
|---|---|
| | |
| | |
| | |
| | S |
| | |
| | |
| | |

# Core *j*

**Load x:**

## Cache lines / state

| | |
|---|---|
| | |
| | |
| | S |
| | |
| | |
| | |

BusRd x

## Memory blocks / state

| | |
|---|---|
| | |
| | |
| x | |
| | |
| | |
| | **34** |

Processor read /
bus read observation:
state for other copies
remains S

Processor operation / Cache controller operation

Observed operation / Cache controller operation

# Write requires invalidation

**Core  i**

Cache lines  /  state

**Core  j**

✦ Invalidate others before
**Store x**
Cache lines  / state

M

BusRdX x

Memory blocks  / state

x

35

Must be exclusive owner before writing to local copy: first invalidate the others and update my copy, by BusRdX x

P

PrRd
PrWr

CacheCtrl

BusRd
BusRdX
BusWr

M

PrRd/-
PrWr/-

BusRd/Flush

PrWr/BusRdX

BusRdX/Flush

PrWr/BusRdX

S

PrRd/-
BusRd/-

PrRd/BusRd

BusRdX/-

I

Processor operation / Cache controller operation

Observed operation / Cache controller operation

# Write requires invalidation

# Read triggers updating

Core *i*

**Load x:**

Cache lines / state

Core *j*

Cache lines / state

S

S

BusRd x

Memory blocks

x

S

Main memory copy is updated from core j's cache when the value of x is on the bus

37

P

PrRd
PrWr

CacheCtrl

BusRd
BusRdX
BusWr

M

PrRd/-
PrWr/-

BusRd/Flush

BusRdX/Flush

PrWr/BusRdX

PrRd/-
BusRd/-

S

PrWr/BusRdX

PrRd/BusRd

BusRdX/-

I

Processor operation / Cache controller operation

Observed operation / Cache controller operation

# Another MSI Example

# Another scenario: Core *j* has written…



**Core  i**

Cache lines  /  state

I

**Core  j**

Cache lines  / state

M

Memory blocks  / state

x                                    I

39

P
PrRd
PrWr
CacheCtrl
BusRd
BusRdX
BusWr

PrRd/-
PrWr/-

M

BusRd/Flush

PrWr/BusRdX          BusRdX/Flush

S   PrRd/-
    BusRd/-

PrWr/BusRdX

PrRd/BusRd

BusRdX/-

I

Processor operation / Cache controller operation

Observed operation / Cache controller operation

# Now Core *i* wants to write…

Core *i*

Invalidate others before

**Store x**

Cache lines / state

M

Upgrade I → M requires update of the cache line in memory and requesting cache before overwriting

Core *j*

Cache lines / state

I

BusRdX x

Observation of BusRdX x: Cache controller flushes cache line to memory and invalidates own copy

Memory blocks / state

x    I

**40**

P

PrRd
PrWr

CacheCtrl

BusRd
BusRdX
BusWr

PrRd/-
PrWr/-

M

BusRd/Flush

PrWr/BusRdX        BusRdX/Flush

PrWr/BusRdX    S    PrRd/-
BusRd/-

PrRd/BusRd

BusRdX/-

I

Processor operation / Cache controller operation

Observed operation / Cache controller operation

# Now Core *i* can write…

**Core *i***

Cache lines / state

M

**Core *j***

Cache lines / state

I

**Memory blocks / state**

x                              I

41

P

PrRd
PrWr

CacheCtrl

BusRd
BusRdX
BusWr

M

PrRd/-
PrWr/-

BusRd/Flush

PrWr/BusRdX

BusRdX/Flush

PrWr/BusRdX

S

PrRd/-
BusRd/-

PrRd/BusRd

BusRdX/-

I

Processor operation / Cache controller operation

Observed operation / Cache controller operation

# MESI-Protocol

MSI protocol:

2 bus operations (BusRd, BusRdX) required

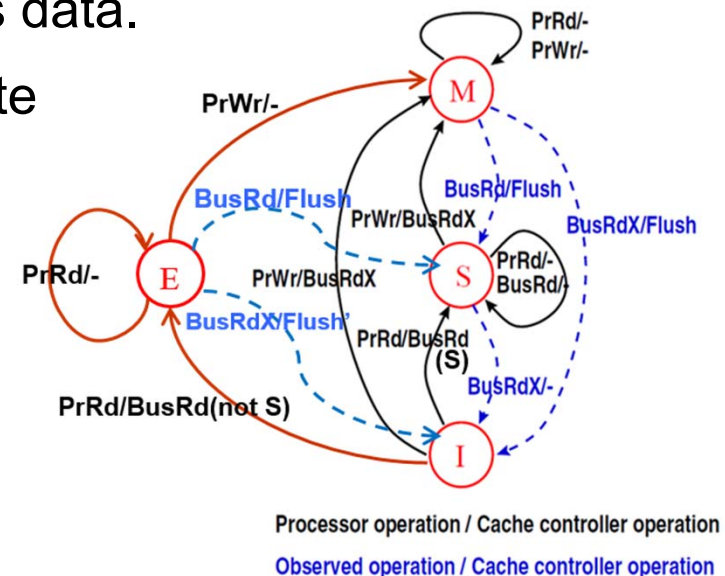if a processor first reads ($\rightarrow S$), then writes ($\rightarrow M$) a memory location,

even if no other processor works on this data.

$\rightarrow$ generalization to MESI-protocol with new state

E (exclusive)

no other cache has a copy of this block,

this copy is not modified and same data

in MM.



Processor operation / Cache controller operation
Observed operation / Cache controller operation
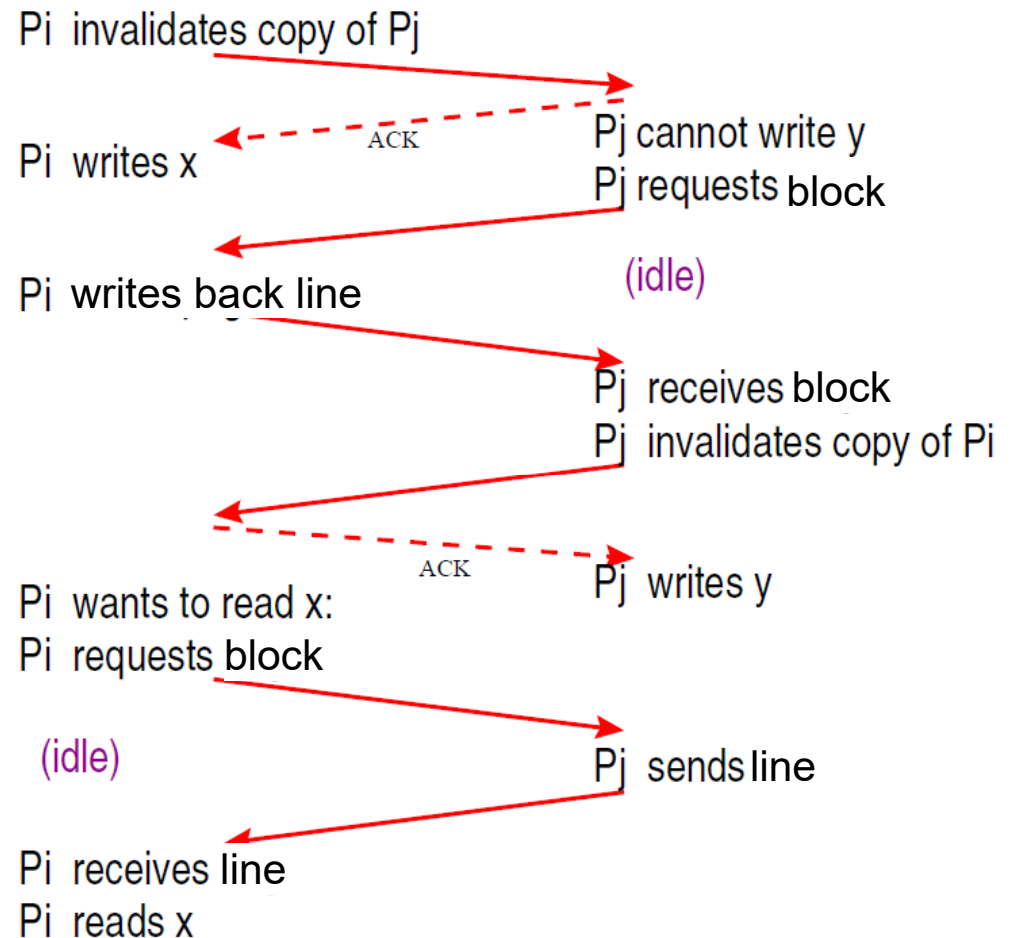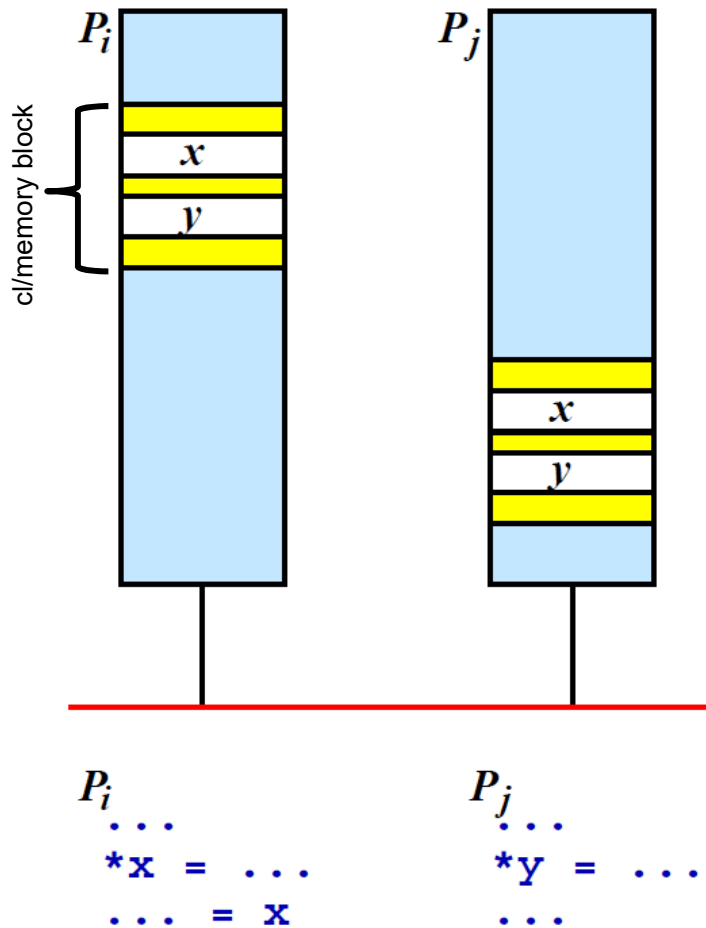
Modifications in MSI-protocol:

+ PrRd to a <u>non-cached</u> address (BusRd): $\rightarrow E$ (not $S$)

+ PrWr to $E$-address: local change to $M$, write (**no bus operation**)

+ read access from another processor to $E$-address (BusRd/Flush): $\rightarrow S$

MESI supported by Intel Pentium, MIPS R4400, IBM, PowerPC, …

45

# Performance Issue: False Sharing

cl/memory block

$P_i$

$P_j$

| x |
| y |

| x |
| y |

$P_i$
...
*x = ...
... = x

$P_j$
...
*y = ...
...

Pi invalidates copy of Pj

Pi writes x          ACK

Pj cannot write y
Pj requests block

Pi writes back line          (idle)

Pj receives block
Pj invalidates copy of Pi

ACK          Pj writes y

Pi wants to read x:
Pi requests block

(idle)          Pj sends line

Pi receives line
Pi reads x

Invalidation causes an extra cache miss
Cache lines or memory blocks are the units of sharing
→ sequentialization, thrashing

# How to Avoid False Sharing?

- Smaller cache lines

    - false sharing less probable, but

    - more administrative effort

- Programmer or compiler gives hints for data placement

    - padding, data privatization

- Time slices for exclusive use of CL:

    each line stays for $\geq d$ time units at one processor

How to reduce performance penalty of false sharing?

- Use weaker consistency models

    - programming more complicated and error-prone

# Shared Memory Consistency Models

- Problem of memory inconsistencies when memory-access order differs from the program execution-order.

- Improve performance at the cost of error-prone programming

- different memory consistency models:
  Strict consistency
  Sequential consistency
  Causal consistency
  Superstep consistency
  "PRAM" consistency          [Culler et al.'98, Ch. 9.1], [Gharachorloo/Adve'96]
  Weak consistency
  Release consistency / Barrier consistency
  Lazy Release consistency
  Entry consistency
  Others (processor consistency, total/partial store ordering etc.)
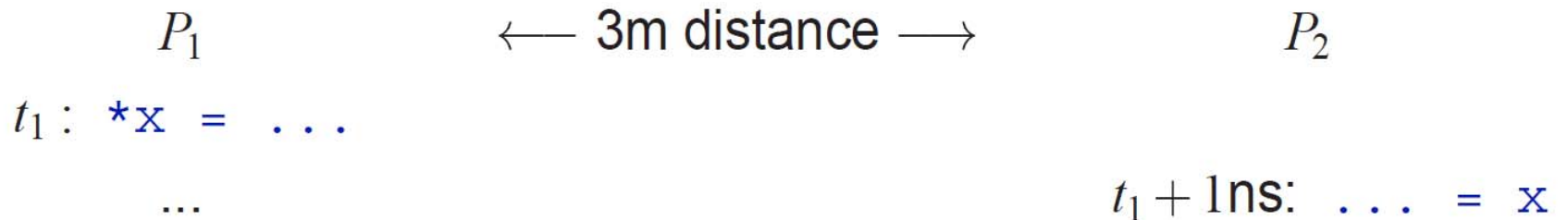
# Consistency Models: Strict Consistency

Strict consistency: no memory inconsistency

Read (x) returns the value that was most recently ($\rightarrow$global time) written to *x*.

Writes must be seen instantaneously by all processors.

realized in classical uniprocessors and SB-PRAM;

in DSM physically impossible without additional synchronization

$$P_1 \qquad\qquad \longleftarrow \text{ 3m distance } \longrightarrow \qquad\qquad P_2$$

$t_1: \quad *x \; = \; \ldots$

$\qquad\qquad \ldots \qquad\qquad\qquad\qquad\qquad\qquad\qquad t_1 + 1\text{ns:} \; \ldots \; = \; x$

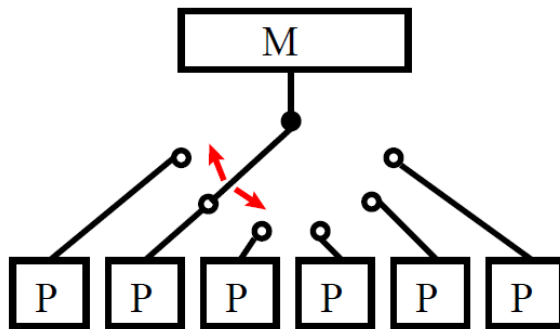Transport of *x* from P$_1$ to P$_2$ with speed $10c$

# Consistency Models: Sequential Consistency

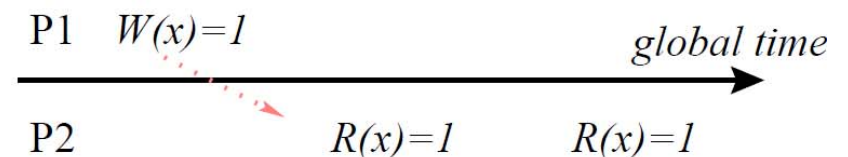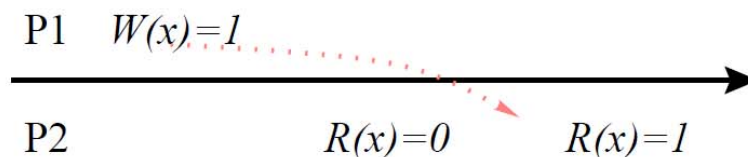Sequential consistency  [Lamport'79] no memory inconsistency

\+ all memory accesses are ordered in some sequential order and
   seen the same way by all processors

\+ all read and write accesses of a processor appear in program order

\+ otherwise, arbitrary delays possible

Read/writes are atomic.

Acces
mode.



Not deterministic:

| P1 | $W(x)=1$ | | |
|----|----------|--|--|
| P2 | | $R(x)=0$ | $R(x)=1$ |

| P1 | $W(x)=1$ | global time |
|----|----------|-------------|
| P2 | $R(x)=1$ | $R(x)=1$ |

# Consistency Models: Weak Consistency

- memory consistency must be guaranteed by programmer not HW
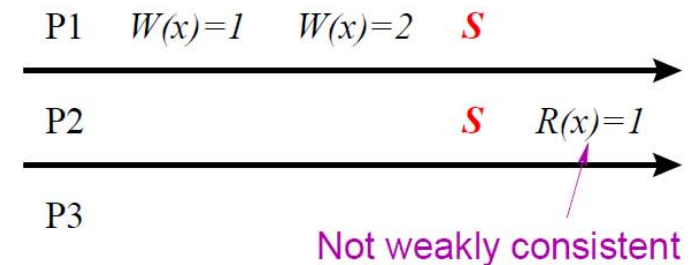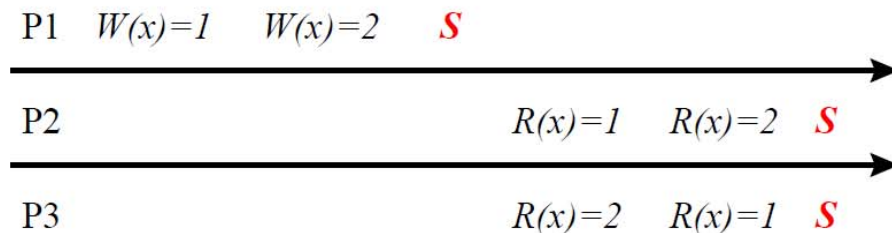
+ Classification of shared variables (and their accesses):

   synchronization variables (locks, semaphores)

   $\rightarrow$ always consistent, atomic access

   other shared variables

   $\rightarrow$ kept consistent by the user, using synchronizations

+ Accesses to synchronisation variables are sequentially consistent

- Accesses to other shared variables may be seen in different order on different processors.

| P1 | W(x)=1 | W(x)=2 | S | |
|----|--------|--------|---|---|
| P2 | | | R(x)=1 | R(x)=2 | S |
| P3 | | | R(x)=2 | R(x)=1 | S |

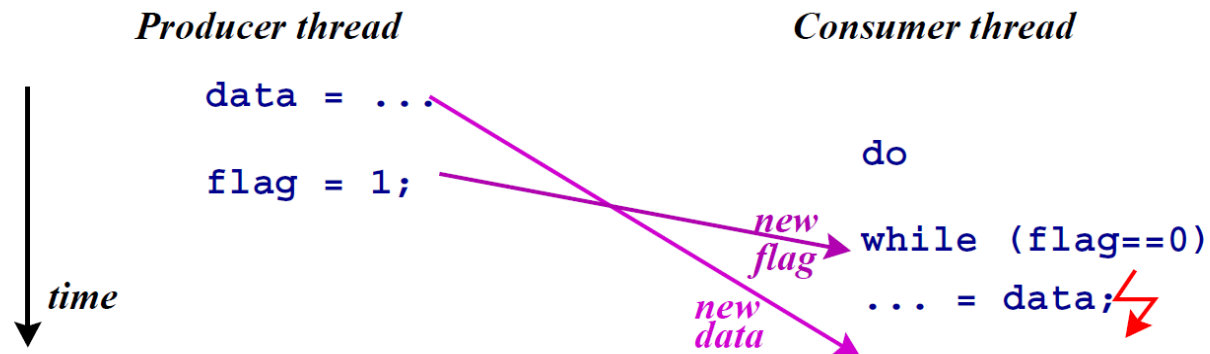| P1 | W(x)=1 | W(x)=2 | S | |
|----|--------|--------|---|---|
| P2 | | | S | R(x)=1 |
| P3 | | | | |

Not weakly consistent

S is a flush operation.

# Consistency Models:
# Weak Consistency in OpenMP

OpenMP implements weak consistency. Inconsistencies may occur due to

+ register allocation

+ compiler optimizations

+ caches with write buffers

**Producer thread**

```
data = ...
flag = 1;
```

*time*

**Consumer thread**

```
            do

new
flag   while (flag==0)

new    ... = data;
data
```

Need explicit „memory fence" to control consistency: flush directive

• write back register contents to memory

• forbid code moving compiler optimizations

• flush cache write buffers to memory
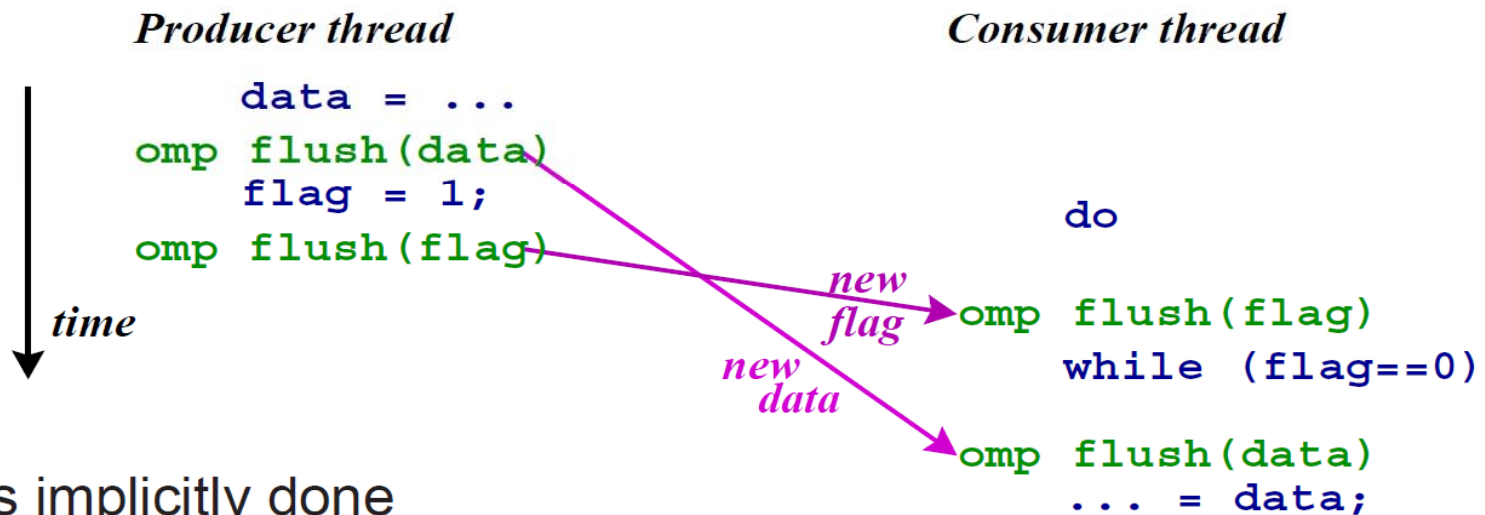
• re-read flushed values from memory

# Consistency Models: Weak Consistency in OpenMP  (cont.)

`!$omp flush ( `*shvarlist*` )`

creates for the executing processor a consistent memory view
for the shared variables in *shvarlist*.

If no parameter: create consistency of all accessible shared variables.

**Producer thread**

```
        data = ...
omp flush(data)
        flag = 1;
omp flush(flag)
```

*time*

**Consumer thread**

```
                do
```

*new flag*

```
omp flush(flag)
        while (flag==0)
```

*new data*

```
omp flush(data)
        ... = data;
```

A flush is implicitly done
at `barrier`, `critical`, `end critical`, `end parallel`,
and at `end do`, `end section`, `end single`
if no `nowait` parameter is given

59

# Summary

- Overview of architectures for shared memory parallel computers

- Caches to reduce latency for data acesses

- Cache coherence to keep data consistent that is stored in caches and main memory

- False sharing and what to do about it

- Memory consistency models.

# Further Reading

- D. Culler et al. *Parallel Computer Architecture, a Hardware/Software Approach.* Morgan Kaufmann, 1998.

- J. Hennessy, D. Patterson: *Computer Architecture, a Quantitative Approach,* Second edition (1996) or later. Morgan Kaufmann.

- S. Adve, K. Gharachorloo: Shared memory consistency models: a tutorial. IEEE Computer, 1996.

# Sources for Slides (Acknowledgements)

- Christoph Kessler, PELAB / IDA,Linköping University, Sweden

- Thomas Fahringer, UIBK