# Prog-order + serial...uniprocessor

- Memory operations = read/write within an instruction are assumed to execute atomically with respect to each other in a specified order

- Memory operation issues when it leaves the processor
  - But goes through the cache, write-buffer, memory-modules
  - i.e. it takes long time to complete

- **Only way processor observes state of memory is by issuing memory operations, e.g. Reads**

- Ex: Processor knows that write operation is performed if subsequent read returns the written value or value of a later write

- Ex: Processor completes a read operation means that subsequent writes to that location cannot affect the read value

- **"Subsequent"** is well defined in sequential process => i.e. they are in program order

# Concepts of program order and serialisation in sequential and parallel case

- We will discuss the concepts in
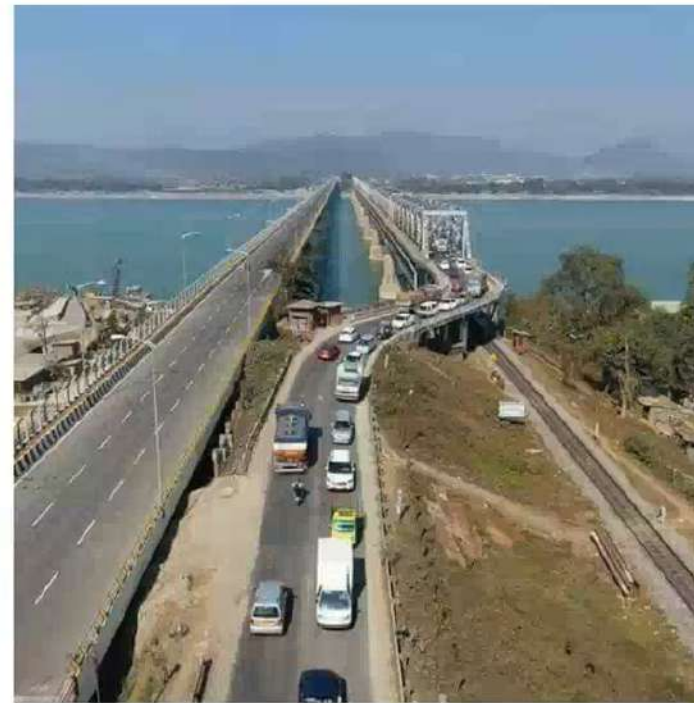  - A uniprocessor - sequential
  - In multiprocessor setup - parallel

# Parallel case ...

- Thus memory system is a point in which hardware determines the order of access

- There is not fixed sequence of access among the processes

- There is arbitrary interleaving

- Fairness -> as each process will eventually access memory

- Our understanding of "last" or "subsequent" access will be defined in this hypothetical serial order

- As the serial order must be consistent, the processes see the writes to a location in the same order

Hemangee K. Kapoor                    22

# Merging of cars
# =
# Merging of mem accesses

# What happens in practice?

- In practice we do not want to construct this serial order. In the presence of caches, ordering is varied

- We just need to make sure that the program behaves as if some serial order was enforced

- *FORMALLY*: Multiprocessor system is coherent if the results of any execution of a program are such that, for each location, it is possible to construct a hypothetical serial order of all operations to that location, that is consistent with the results of execution and in which

  - (1) operations issued by a process occur in its program order, AND

  - (2) value returned by read operation is the value written by the last write to that location in the serial order

# Defining Correctness Metrics

# Definition: Coherence

- Informally, we could say that a memory system is coherent if any read of a data item returns the most recently written value of that data item

- Easy in uniprocessor. But too vague and simplistic => involves 2 aspects of memory system behaviour

- Coherence and Consistency

| COHERENCE | CONSISTENCY |
|---|---|
| (1) Defines what values can be returned by a read | (1) Determines when a written value will be returned by a read |
| (2) Defines behaviour of same location | (2) Defines behaviour to other locations |

# Coherent view of Memory

A read by a processor to a location X that follows a write by another processor to X returns the written value if the read and write are sufficiently separated in time and no other writes to X occurs between the two accesses

P1

P2

P3

P4

WR $X = 2$

WR $X = 5$

Sees
$X = 2$
then
$X = 5$

Sees
$X = 5$
then
$X = 2$

Write
serialisation

Mem
$X = 2$
$X = 5$

time
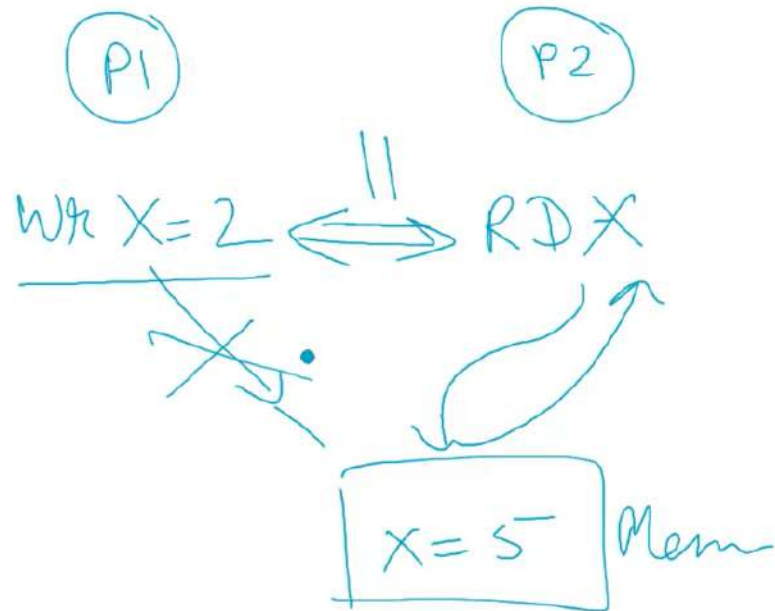
Mem
$X = 5$
$X = 2$

time

Hemangee K. Kapoor          33

# When a written value will be read? Write Consistency

- If P1 writes to X and then P2 reads X. But the read/write are so close in time that P1 has not updated memory and so P2 gets stale data value

- In other words, it may be impossible to ensure that the read returns the value of the data written, since the written data may not even have left the processor at that point.

- We deal with this in the section on memory consistency models

Hemangee K. Kapoor          34

# Two assumptions before consistency

- We make following two assumptions until we deal with memory consistency:

  (1) A write does not complete, and allow the next write to occur, until all processors have seen the effect of that write

  (2) The processor does not change the order of any write with respect to any other memory access

- e.g. If processor writes to location A followed by location B, any processor that sees the new value of B must also see the new value of A

  => These restrictions allows the processor to reorder reads, but writes finish in program order

# 2 properties

- Two properties implicit in the definition of coherence

  - Write propagation
    - Writes become visible to other processes

  - Write serialisation
    - All writes to a location (from same or different processes) are seen in the same order by all processes
    - e.g. if P1 reads the values to same location in the order: w1 then w2 ; then all other processes see the writes in the same order. i.e. no one can see them as w2 then w1

# Defining Coherent Memory System

3 conditions:

(1) Preserve Program Order

(2) Coherent View of memory

(3) Write Serialisation

# Cache Coherence Protocols

- Key to implementation is tracking the state of any sharing of data block

- Two classes of protocol:
  - Snooping-based and directory-based

# Basic Implementation Techniques (Architectural Building Blocks)

- All processors continuously snoop on the bus to check if address on the bus is in their cache **[CACHE]**

  <u>State</u> can be changed to <u>valid, invalid, exclusive</u>

- When a write occurs to a shared block --> acquire broadcast medium (i.e. bus) to **send inv** request

  If 2 processors try to write, they have to <u>arbitrate for the bus</u>

- The processor which gets the bus first, writes the value and the second processor has to invalidate its own copy

- The 2nd processor can later get the bus and do the write

- => **[BUS]** enforces <u>write-serialisation</u>

- When new cache miss occurs for an invalidated shared block, where to locate up-to-date copy of data? **[LATEST COPY]**

Hemangee K. Kapoor

8

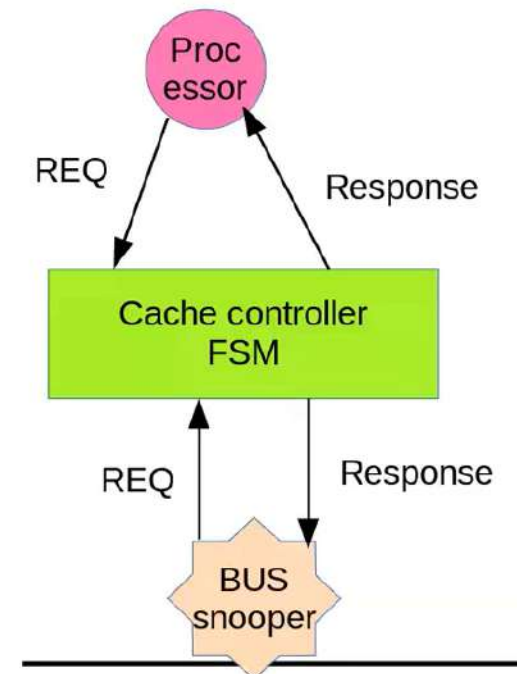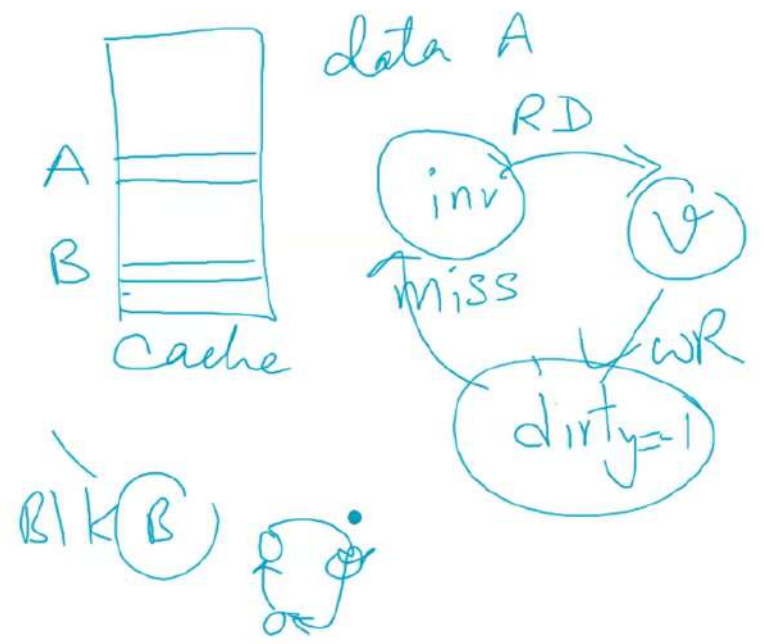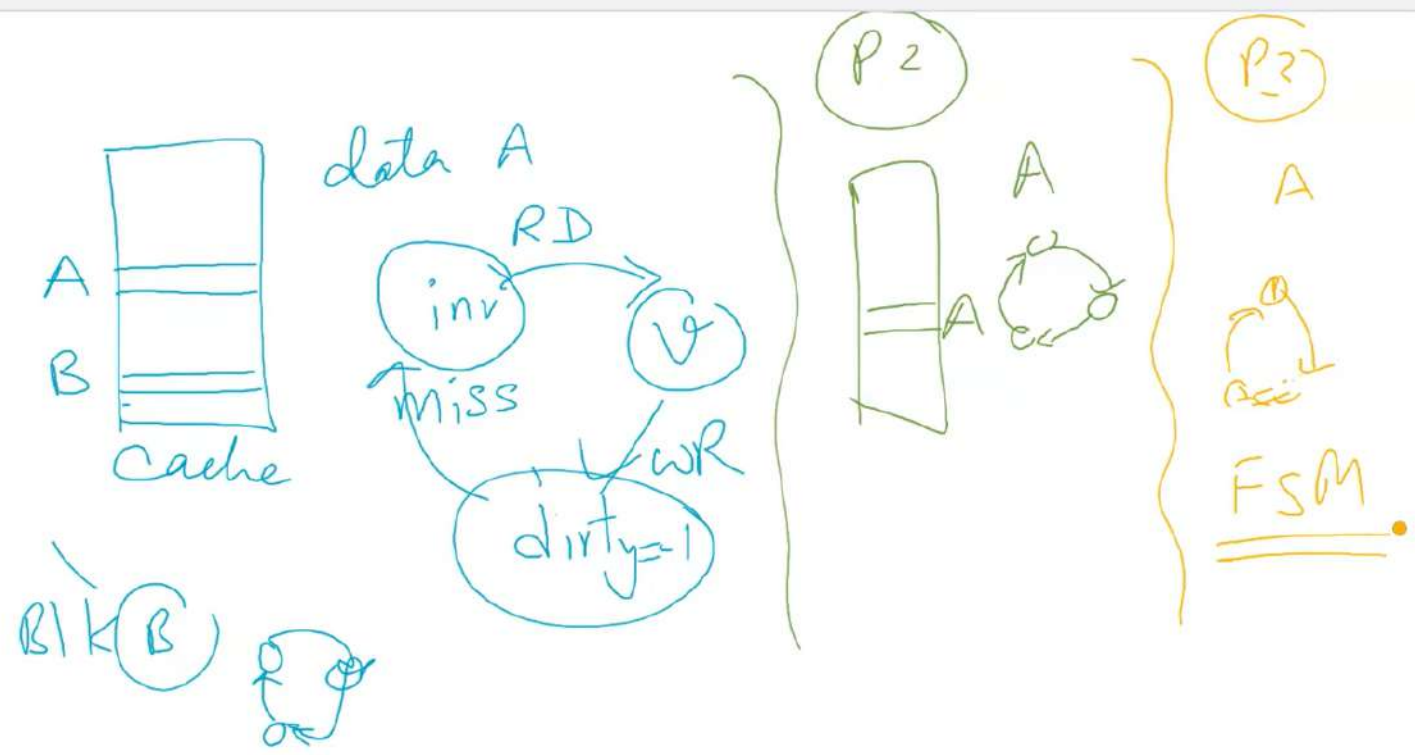# Protocol Implementation/Modelling

- In multiprocessor, each block has a state in each cache and the states change according to the state transition diagram

- If there are 'p' caches, then the block state is a vector of size=p

- Cache state is manipulated by a set of 'p' distributed finite state machines, implemented by the cache controllers

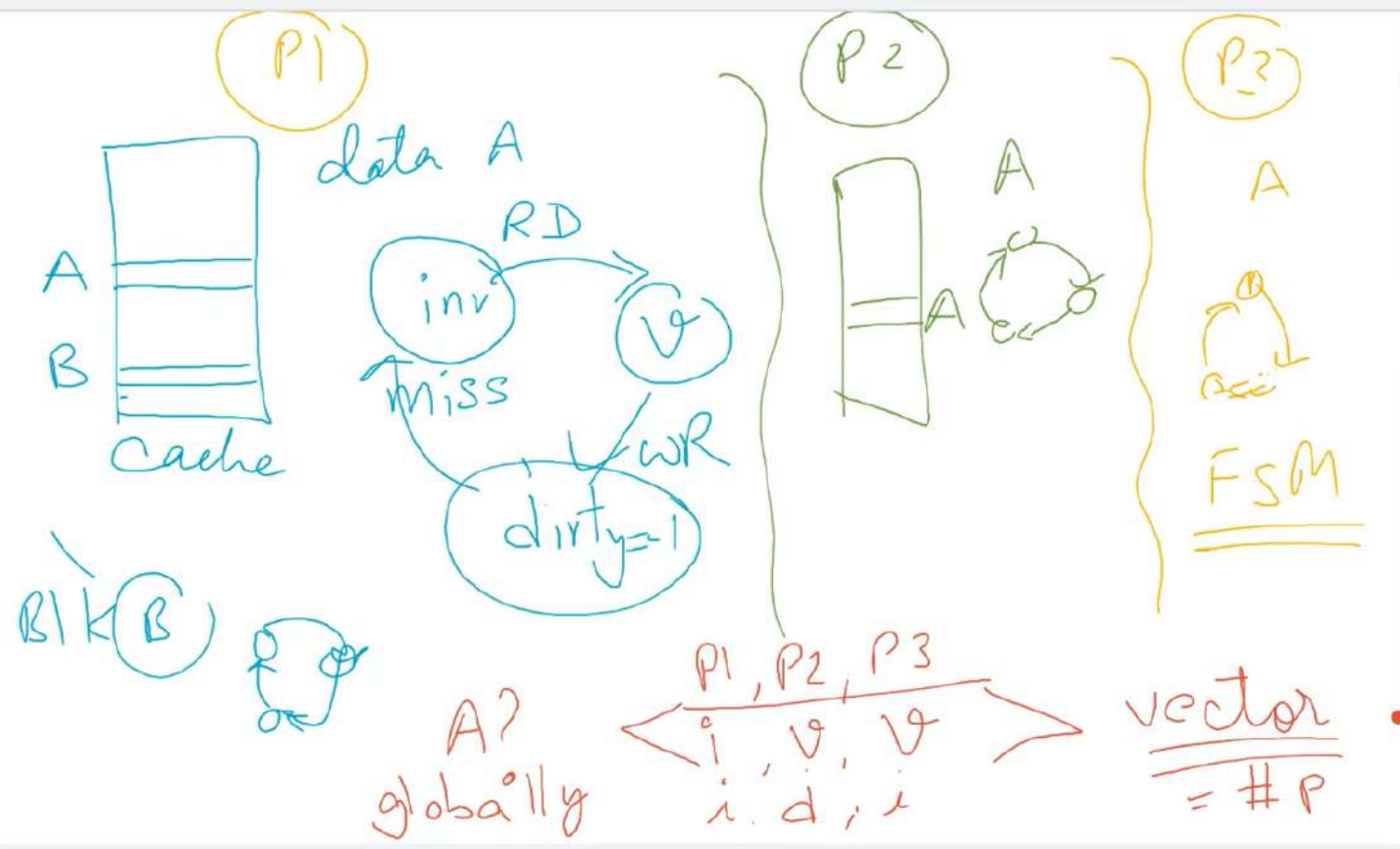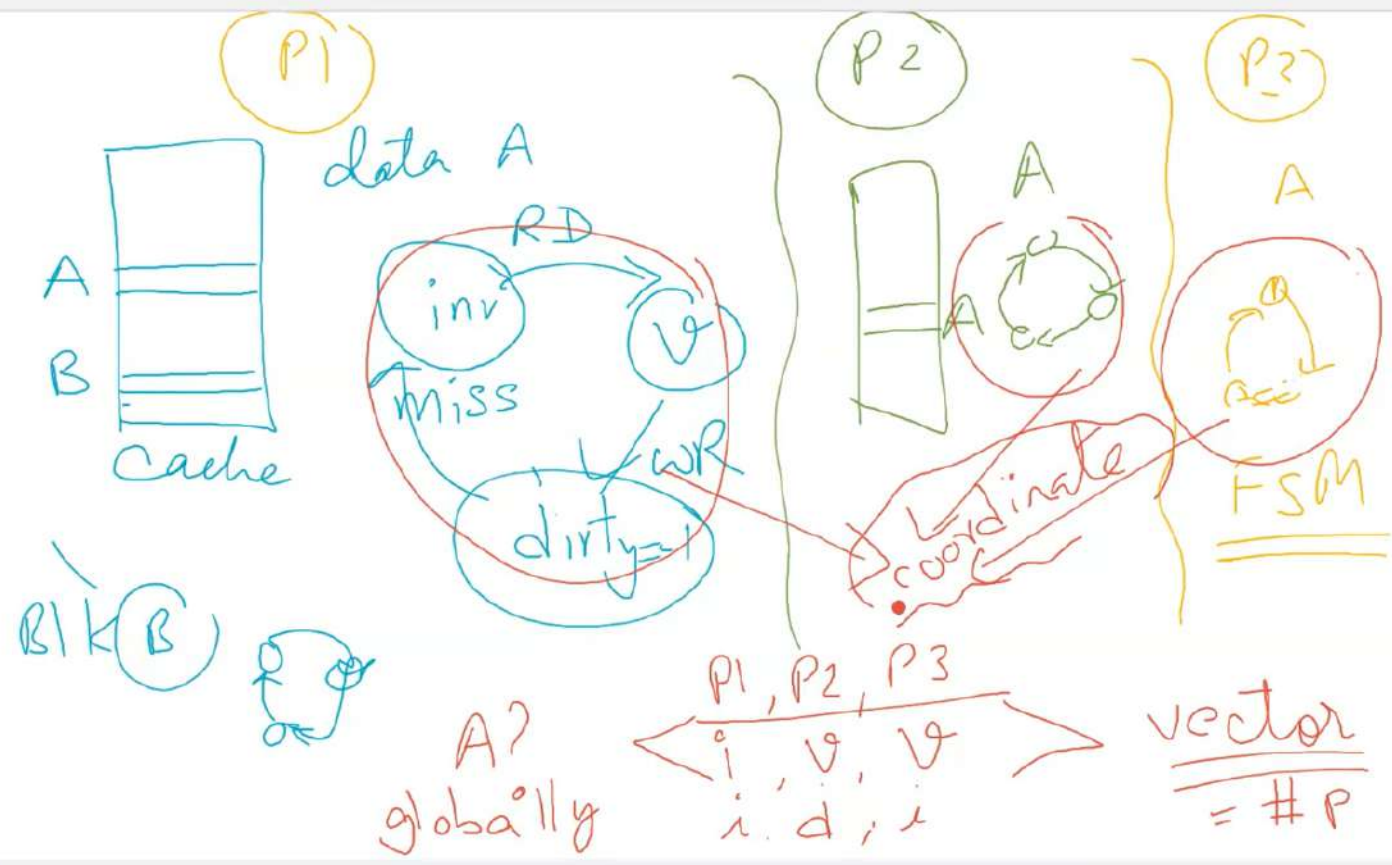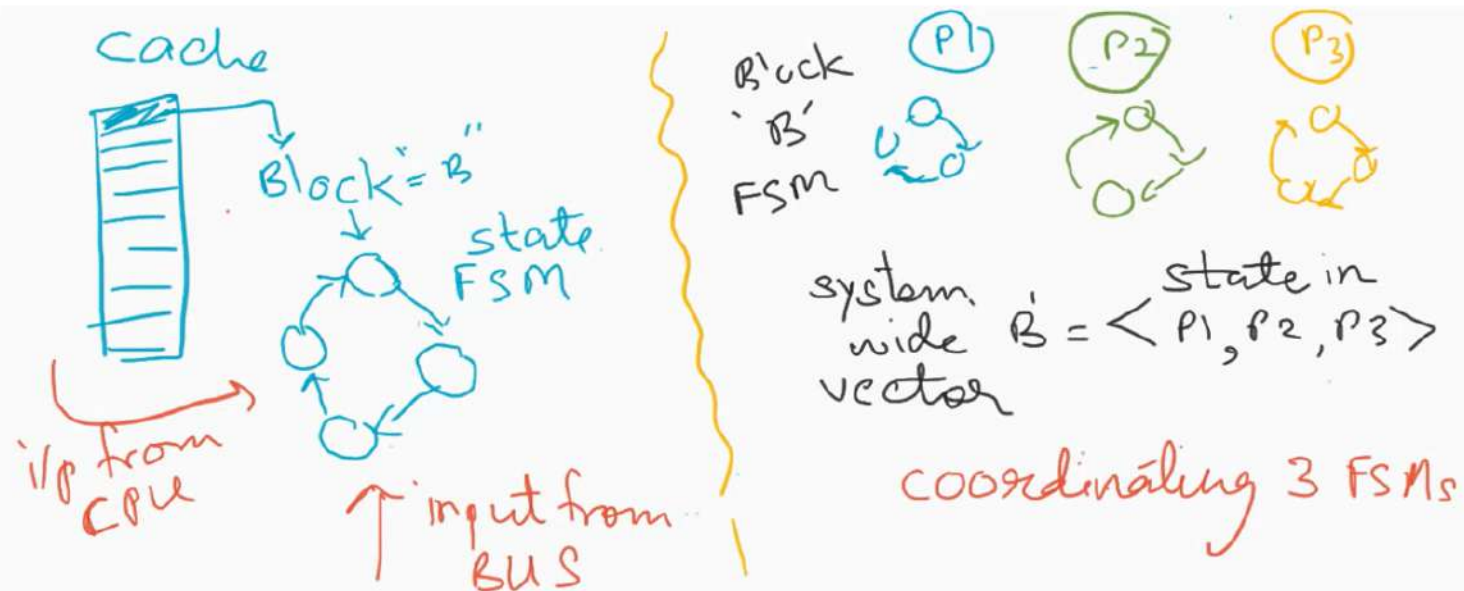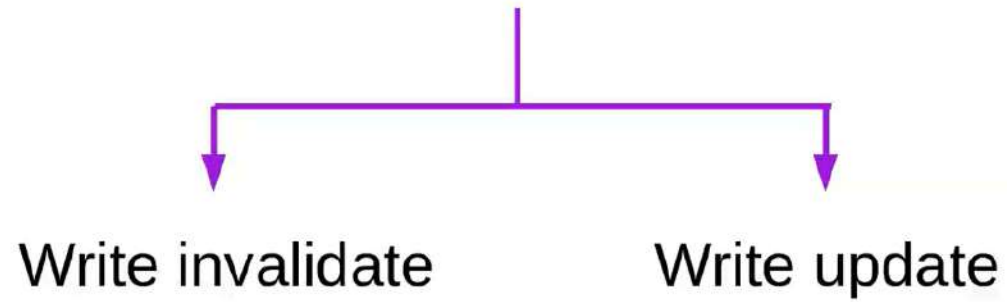- Controller FSMs will receive inputs from processor or bus snooper, and they act accordingly

# Snooping protocol=distributed FSM

- Thus the snooping protocol is a distributed algorithm represented by a collection of cooperating FSMs

- All these FSMs are cooridnated by bus transactions

- It is specified by the following components
  - Set of **states** of memory blocks in local caches
  - State **transition diagram**, which takes input from processor or bus and produces the next state
  - The **actions** associated with each state transition that tell what actions to do on the bus, the cache and the processor

# Example: Wr-inv protocol

State | Address | Data

Cache

x = 1

x : 0

Cache

x : 0

Invalidate all X at (3)

x : 0

Memory

(1) A reads X=0

(2) B read X=0

(3) A writes X=1;
   Inv all copies,
   memory updated

(4) B reads X
   => results in
      cache miss

Cache containts
<state, addr, data>

# Ex: write-thru cache

- 200 Mhz dual issue, CPI=1, 15% stores of 8 bytes

  How many processors will 1 GB/s bus support before saturation?

- ANS=

  Single processor will generate how many writes? Stores?

  =0.15 stores/instr * 1 instr/cycle * 200 * 10^6 cycle/sec

  = 30 * 10^6 stores/sec

  Each store = 8 bytes

  Per processor  8 * 30 * 10^6 bytes/sec = 240 MB/sec

  1 GB/s / 240 MB = 4 processors will results in saturation of bus

  Max support up to 4 processors

Hemangee K. Kapoor                18