

NC STATE UNIVERSITY

# Other Bus-Based Coherence Protocols

### Lecture 15 (Chapter 7, cont.)

E. F. Gehringer, based on slides by Yan Solihin



# Lecture 15 Outline

- MSI protocol
  - State diagram
  - Animations
- MESI protocol
- Dragon protocol
- Firefly protocol

	Inval- idate	Update
3-state	MSI	Firefly
4-state	MESI	Dragon

2

### **Basic MSI Writeback Invalidation Protocol**

- States
  - Invalid (I)
  - Shared (S): one or more copies, and memory copy is up-to-date
  - Dirty or Modified (M): only one copy
- Processor Events:
  - PrRd (read), PrWr (write)
- Bus Transactions
  - BusRd: asks for copy with no intent to modify
  - BusRdX: asks for copy with intent to modify (instead of BusWr)
  - Flush: updates memory
- Actions
  - Update state, perform bus transaction, flush value onto bus





- On the following slides, we will display the state-transition diagrams
  - for processor-initiated transactions
  - for bus-initiated transactions
- We will see transitions of the following form:
  - Invalidation:  $\langle Any \rangle \rightarrow I$
  - Intervention: {Exclusive, Modified} → Shared

4

## **MSI: Processor-Initiated Transactions**



Ē

CSC/ECE 506: Architecture of Parallel Computers

5



# **MSI: Bus-Initiated Transactions**



NC STATE UNIVERSITY



# **MSI State Transition Diagram**



CSC/ECE 506: Architecture of Parallel Computers

#### NC STATE UNIVERSITY



# Lecture 15 Outline

- MSI protocol
  - State diagram
  - Animations
- MESI protocol
- Dragon protocol
- Firefly protocol

	Inval- idate	Update
3-state	MSI	Firefly
4-state	MESI	Dragon

8



### **MSI** Visualization – Start State



Trace

P <sub>1</sub>	P <sub>2</sub>	P <sub>3</sub>
Cache	Cache	Cache
Snooper	Snooper	Snooper
Bus	Controller $A = 1$ Main memory	

MSI	Firefly
MESI	Dragon

#### \_

9

### NC STATE UNIVERSITY





10

#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





### CSC/ECE 506: Architecture of Parallel Computers

#### NC STATE UNIVERSITY





### NC STATE UNIVERSITY CSC





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY
### MSI: Processor P<sub>2</sub> Reads A



### NC STATE UNIVERSITY



# MSI: Processor P<sub>2</sub> Reads A



### NC STATE UNIVERSITY

# MSI Example: Rd/Wr to a single line

Proc Action	State P1	State P2	State P3	Bus Action	Data From
R1	S	_	_	BusRd	Mem
W1	М	-	-	BusRdX*	Mem
R3	S	-	S	BusRd/Flush	P1 cache
W3	Ι	-	М	BusRdX*	Mem
R1	S	_	S	BusRd/Flush	P3 cache
R3	S	_	S	—	Own Cache
R2	S	S	S	BusRd	Mem

\*or, BusUpgr (data from own cache)



# Notes on MSI Protocol

• For  $M \rightarrow I$ , BusRdX/Flush: why flush?



# Notes on MSI Protocol

- For M → I, BusRdX/Flush: why flush? Because it is a read with intention to write, as opposed to write.
  - Thus, there is a possibility for a read before the write is performed.
  - In addition, the write could be to a different word in the line (so the whole line needs to be flushed).



# Notes on MSI Protocol

- For M → I, BusRdX/Flush: why flush? Because it is a read with intention to write, as opposed to write.
  - Thus, there is a possibility for a read before the write is performed.
  - In addition, the write could be to a different word in the line (so the whole line needs to be flushed).
- In case of a write to a shared block:
  - Cache already has latest data; can use upgrade (BusUpgr) instead of BusRdX
- Replacement changes state of two blocks: outgoing and incoming
- Flush has to modify both caches and main memory

*Note:* Coherence granularity is u (a single line). What happens when all the reads go to word 0 on line u, but write by P3 goes to word 1 on line u? False-sharing miss on the 2nd R1

CSC/ECE 506: Architecture of Parallel Computers

42

# **MSI: Coherence and SC**

- Coherence:
  - Write propagation:
    - through invalidation, and flush on subsequent BusRds
  - Write serialization?
    - Writes (BusRdX) that go to the bus appear in bus order (and handled by snoopers in bus order!)
    - Writes that do not go to the bus?
      - Only happen when the line state is M, i.e. when I am the only processor holding the line. Local writes are only visible to me, so they are serialized.
- To enforce SC:
  - Program order: enforced by following the bus transaction order
    - All writes appear on the bus
    - All local writes (within 1 processor) can follow program order
  - Write completion: Occurs when write appears on bus
  - Write atomicity: A read returns the latest value of a write. At that time, the value is visible to all others (on a bus transaction, or on a local write).



# Lecture 15 Outline

- MSI protocol
- MESI protocol
- Dragon protocol
- Firefly protocol

	Inval- idate	Update
3-state	MSI	Firefly
4-state	MESI	Dragon

44

# **Lower-Level Protocol Choice**

 What transition should occur when a BusRd is observed in state M?

- Should the state change to S or to I?

# MESI (4-state) Invalidation Protocol

- Here's a problem with the MSI protocol:
  - A {Rd, Wr} sequence causes two bus transactions
    - BusRd (I  $\rightarrow$  S) followed by BusRdX or BusUpgr (S  $\rightarrow$  M)
    - even when no one is sharing (e.g., serial program!)
    - In general, coherence traffic from serial programs is unacceptable
- To avoid this, add a fourth state, Exclusive:
  - Invalid
  - Modified (dirty)
  - Shared (two or more caches may have copies)
  - Exclusive (only this cache has clean copy, same value as in memory)



46

- How does the protocol decide whether  $I \rightarrow E$  or  $I \rightarrow S$ ?
  - Need to check whether someone else has a copy
  - "Shared" signal on bus: wired-or line asserted in response to BusRd

# **MESI: Processor-Initiated Transactions**



Ę



# **MESI: Bus-Initiated Transactions**

Flush' means flush only if cacheto-cache sharing is used; only the cache responsible for supplying the data will do a flush.



NC STATE UNIVERSITY

CSC/ECE 506: Architecture of Parallel Computers

48



# **MESI State Transition Diagram**



 BusRd(S) means shared line asserted on BusRd transaction

49



# **MESI** Visualization



#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





### NC STATE UNIVERSITY



# Processor $P_1$ Writes A = 2





# Processor $P_1$ Writes A = 2



### NC STATE UNIVERSITY





### CSC/ECE 506: Architecture of Parallel Computers

#### NC STATE UNIVERSITY





### CSC/ECE 506: Architecture of Parallel Computers

### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





### NC STATE UNIVERSITY



**NC STATE UNIVERSITY** 

# Processor P<sub>3</sub> Reads A





# Processor $P_3$ Writes A = 3



### CSC/ECE 506: Architecture of Parallel Computers

#### NC STATE UNIVERSITY



# Processor $P_3$ Writes A = 3



### NC STATE UNIVERSITY



# Processor $P_3$ Writes A = 3



### CSC/ECE 506: Architecture of Parallel Computers

#### NC STATE UNIVERSITY



**NC STATE UNIVERSITY** 

# Processor $P_3$ Writes A = 3







### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





### CSC/ECE 506: Architecture of Parallel Computers

#### NC STATE UNIVERSITY





### NC STATE UNIVERSITY





### NC STATE UNIVERSITY


**NC STATE UNIVERSITY** 

### Processor P<sub>3</sub> Reads A







#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





### CSC/ECE 506: Architecture of Parallel Computers

#### NC STATE UNIVERSITY



**NC STATE UNIVERSITY** 

### Processor P<sub>2</sub> Reads A



### MESI Example (Cache-to-Cache Transfer)

Proc Action	State P1	State P2	State P3	Bus Action	Data From
R1	E	-	—	BusRd	Mem
W1	М	-	_	_	Own cache
R3	S	—	S	BusRd/Flush	P1 cache
W3	I	-	М	BusRdX	Mem
R1	S	-	S	BusRd/Flush	P3 cache
R3	S	_	S	_	Own cache
R2	S	S	S	BusRd/Flush <sup>,</sup>	P1/P3 Cache*

\* Data from memory if no cache-to-cache transfer, BusRd/ -

### Change from MSI (Cache-to-Cache Transfer)

Proc Action	State P1	State P2	State P3	Bus Action	Data From
R1	E	-	-	BusRd	Mem
W1	М	-	-	-	Own cache
R3	S	-	S	BusRd/Flush	P1 cache
W3	I	_	М	BusRdX	Mem
R1	S	_	S	BusRd/Flush	P3 cache
R3	S	_	S	_	Own cache
R2	S	S	S	BusRd/Flush <sup>,</sup>	P1/P3 Cache*

\* Data from memory if no cache-to-cache transfer, BusRd/ -

### Change from MSI (Cache-to-Cache Transfer)

Proc Action	State P1	State P2	State P3	Bus Action	Data From
R1	E	-	-	BusRd	Mem
W1	М	-	-	_	Own cache
R3	S	-	S	BusRd/Flush	P1 cache
W3	I	_	М	BusRdX	Mem
R1	S	_	S	BusRd/Flush	P3 cache
R3	S	_	S	_	Own cache
R2	S	S	S	BusRd/Flush <sup>,</sup>	P1/P3 Cache*

\* Data from memory if no cache-to-cache transfer, BusRd/ -

### Change from MSI (Cache-to-Cache Transfer)

Proc Action	State P1	State P2	State P3	Bus Action	Data From
R1	E	-	-	BusRd	Mem
W1	М	-	-	_	Own cache
R3	S	-	S	BusRd/Flush	P1 cache
W3	I	_	М	BusRdX	Mem
R1	S	_	S	BusRd/Flush	P3 cache
R3	S	_	S	_	Own cache
R2	S	S	S	BusRd/Flush <sup>,</sup>	P1/P3 Cache*

\* Data from memory if no cache-to-cache transfer, BusRd/ -

### MESI Example (Cache-to-Cache Transfer+BusUpgr)

Proc Action	State P1	State P2	State P3	Bus Action	Data From
R1	E	-	-	BusRd	Mem
W1	М	-	-	-	Own cache
R3	S	-	S	BusRd/Flush	P1 cache
W3	I	-	М	BusUpgr	Own cache
R1	S	-	S	BusRd/Flush	P3 cache
R3	S	-	S	-	Own cache
R2	S	S	S	BusRd/Flush'	P1/P3 Cache*

\* Data from memory if no cache-to-cache transfer, BusRd/ -

### MESI Example (Cache-to-Cache Transfer+BusUpgr)

Proc Action	State P1	State P2	State P3	Bus Action	Data From
R1	E	-	-	BusRd	Mem
W1	М	-	-	-	Own cache
R3	S	-	S	BusRd/Flush	P1 cache
W3	I	-	М	BusUpgr	Own cache
R1	S	-	S	BusRd/Flush	P3 cache
R3	S	-	S	-	Own cache
R2	S	S	S	BusRd/Flush'	P1/P3 Cache*

\* Data from memory if no cache-to-cache transfer, BusRd/ -

**NC STATE UNIVERSITY** 

### **Lower-Level Protocol Choices**

- Who supplies data on miss when not in M state: memory or cache?
- Original, Illinois MESI: cache
  - assumes cache is faster than memory (cache-to-cache transfer)
  - Not necessarily true
- Adds complexity
  - How does memory know it should supply data? (must wait for caches)
  - A selection algorithm is needed if multiple caches have valid data.
- Useful in a distributed-memory system
  - May be cheaper to obtain from nearby cache than distant memory
  - Especially when constructed out of SMP nodes (Stanford DASH)



# Lecture 15 Outline

- MSI protocol
- MESI protocol
- Dragon protocol
- Firefly protocol

	Inval- idate	Update
3-state	MSI	Firefly
4-state	MESI	Dragon

85

### **Dragon Writeback Update Protocol**

- Four states
  - Exclusive-clean (E): Memory and I have it
  - Shared clean (Sc): I, others, and maybe memory, but I'm not owner
  - Shared modified (Sm): I and others but not memory, and I'm the owner
    - Sm and Sc can coexist in different caches, with at most one Sm
  - **Modified** or dirty (M): I and, no one else
  - On replacement: Sc can silently drop, Sm has to flush
- No invalid state
  - If in cache, cannot be invalid
  - If not present in cache, can view as being in not-present or invalid state
- New processor events: PrRdMiss, PrWrMiss
  - Introduced to specify actions when block not present in cache
- New bus transaction: BusUpd
  - Broadcasts single word written on bus; updates other relevant caches



# **Dragon: Processor-Initiated Transactions**



NC STATE UNIVERSITY

Ē

### **Dragon: Bus-Initiated Transactions**



NC STATE UNIVERSITY

CSC/ECE 506: Architecture of Parallel Computers

88

### **Dragon State Transition Diagram**



CSC/ECE 506: Architecture of Parallel Computers

89

#### NC STATE UNIVERSITY



# **Dragon Visualization**



#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





92

#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY



### Processor $P_1$ Writes A = 2



#### NC STATE UNIVERSITY



### Processor $P_1$ Writes A = 2







### CSC/ECE 506: Architecture of Parallel Computers

#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY



**NC STATE UNIVERSITY** 

### Processor P<sub>3</sub> Reads A





### Processor $P_3$ Writes A = 3



### CSC/ECE 506: Architecture of Parallel Computers

#### NC STATE UNIVERSITY



### Processor $P_3$ Writes A = 3



#### NC STATE UNIVERSITY



### Processor $P_3$ Writes A = 3



CSC/ECE 506: Architecture of Parallel Computers

#### NC STATE UNIVERSITY



**NC STATE UNIVERSITY** 

### Processor $P_3$ Writes A = 3





**NC STATE UNIVERSITY** 

### Processor P<sub>1</sub> Reads A







#### NC STATE UNIVERSITY





### CSC/ECE 506: Architecture of Parallel Computers

#### NC STATE UNIVERSITY




### CSC/ECE 506: Architecture of Parallel Computers

#### NC STATE UNIVERSITY





### CSC/ECE 506: Architecture of Parallel Computers

#### NC STATE UNIVERSITY



**NC STATE UNIVERSITY** 

# Processor P<sub>3</sub> Reads A







#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY



**NC STATE UNIVERSITY** 

# Processor P<sub>2</sub> Reads A





**NC STATE UNIVERSITY** 

# $P_1$ Replaces A





# P<sub>1</sub> Replaces A



#### NC STATE UNIVERSITY



# **Dragon Example**

Proc Action	State P1	State P2	State P3	Bus Action	Data from
R1	E	-	-	BusRd	Mem
W1	М	-	—	_	Own cache
R3	Sm	_	Sc	BusRd/Flush	P1 cache
W3	Sc	-	Sm	BusUpd/Upd	Own cache
R1	Sc	_	Sm	-	Own cache
R3	Sc	_	Sm	_	Own cache
R2	Sc	Sc	Sm	BusRd/Flush	P3 cache

### **Lower-Level Protocol Choices**

- Can shared-modified state be eliminated?
  - If memory is updated too on BusUpd transactions (DEC Firefly)
  - Dragon protocol doesn't (assumes DRAM memory slow to update)
- Should replacement of an Sc block be broadcast?
  - Would allow last copy to go to Exclusive state and not generate updates
  - Replacement bus transaction isn't in critical path, but later update may be
- Shouldn't update local copy on write hit before controller gets bus
  - Can mess up serialization
- Coherence, consistency considerations much like write-through case
- In general, there are many subtle race conditions in protocols.



# Lecture 15 Outline

- MSI protocol
- MESI protocol
- Dragon protocol
- Firefly protocol

	Inval- idate	Update
3-state	MSI	Firefly
4-state	MESI	Dragon

# A Three-State Update Protocol

- Whenever a bus update is generated, suppose that main memory—as well as the caches updates its contents.
- Then which state don't we need?
- What's the advantage, then, of having the fourth state?

 The Firefly protocol, named after a multiprocessor workstation developed by DEC, is an example of such a protocol.

# **Firefly State-Transition Diagram**



Key:

CRM — CPU read miss CWM — CPU write miss CWH — CPU write hit BR — bus read BW — bus write

A " $\sqrt{}$ " following a transition means SharedLine was asserted. An "x" means it was not.

Processor-induced transitions — Bus-induced transitions

Read hits do not cause state transitions and are not shown.

• Answer some questions about this diagram.





123

#### NC STATE UNIVERSITY

Ę





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





126

#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY



# Processor $P_1$ Writes A = 2



#### NC STATE UNIVERSITY



### Processor $P_1$ Writes A = 2



#### 129

#### NC STATE UNIVERSITY





### CSC/ECE 506: Architecture of Parallel Computers

#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY



**NC STATE UNIVERSITY** 

# Processor P<sub>3</sub> Reads A





# Processor $P_3$ Writes A = 3



### CSC/ECE 506: Architecture of Parallel Computers

#### NC STATE UNIVERSITY



# Processor $P_3$ Writes A = 3



### NC STATE UNIVERSITY CSC/EC



# Processor $P_3$ Writes A = 3



### CSC/ECE 506: Architecture of Parallel Computers

#### NC STATE UNIVERSITY



**NC STATE UNIVERSITY** 

# Processor $P_3$ Writes A = 3







CSC/ECE 506: Architecture of Parallel Computers

#### NC STATE UNIVERSITY





### CSC/ECE 506: Architecture of Parallel Computers

#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY





#### NC STATE UNIVERSITY


## Processor P<sub>2</sub> Reads A



### CSC/ECE 506: Architecture of Parallel Computers

#### NC STATE UNIVERSITY



### Processor P<sub>2</sub> Reads A



### CSC/ECE 506: Architecture of Parallel Computers

#### NC STATE UNIVERSITY

# Firefly Example

Proc Action	State P1	State P2	State P3	Bus Action	Data From
R1	V	-	—	BusRd	Mem
W1	D	-	-	-	Own cache
R3	S	_	S	BusRd/Flush	P1 cache
W3	S	_	S	BusUpd	Own cache
R1	S	_	S	-	Own cache
R3	S	-	S	_	Own cache
R2	S	S	S	BusRd/Flush	P1 Cache

# Firefly Example

Proc Action	State P1	State P2	State P3	Bus Action	Data From
R1	V	-	-	BusRd	Mem
W1	D	-	-	-	Own cache
R3	S	-	S	BusRd/Flush	P1 cache
W3	S	_	S	BusUpd	Own cache
R1	S	-	S	-	Own cache
R3	S	_	S	_	Own cache
R2	S	S	S	BusRd/Flush	P1 Cache

CSC/ECE 506: Architecture of Parallel Computers

### Assessing Protocol Tradeoffs

- In the next lecture, we will look at results of comparing protocols by simulation.
- Methodology:
  - Use simulator; default 1MB, 4-way cache, 64-byte block, 16 processors. Some runs use 64K cache.
  - Focus on frequencies, not end performance for now
    - transcends architectural details, but not what we're really after
  - Use idealized memory performance model to avoid changes of reference interleaving across processors with machine parameters
    - Cheap simulation: no need to model contention

CSC/ECE 506: Architecture of Parallel Computers