

NC STATE UNIVERSITY

Invalidate and Update Protocols

Lecture 14
(Chapter 7, cont.)

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

CSC/ECE 506: Architecture of Parallel Computers

1

1

State-Transition Diagrams

- 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: $\langle \text{Any} \rangle \rightarrow I$
 - Intervention: $\{\text{Exclusive, Modified}\} \rightarrow \text{Shared}$

CSC/ECE 506: Architecture of Parallel Computers

4

4

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

CSC/ECE 506: Architecture of Parallel Computers

2

2

MSI: Processor-Initiated Transactions

CSC/ECE 506: Architecture of Parallel Computers

5

5

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 BusRd)
 - Flush: updates memory
- Actions**
 - Update state, perform bus transaction, flush value onto bus

CSC/ECE 506: Architecture of Parallel Computers

3

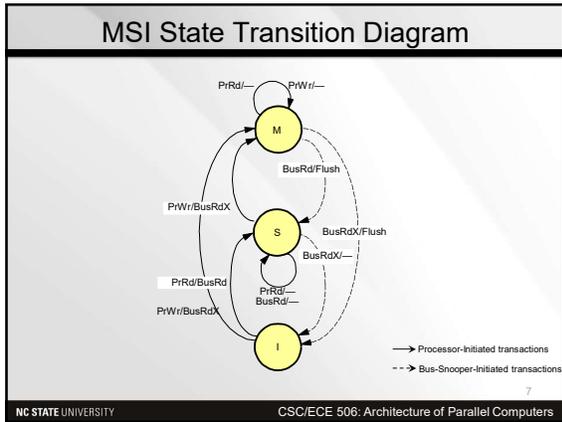
3

MSI: Bus-Initiated Transactions

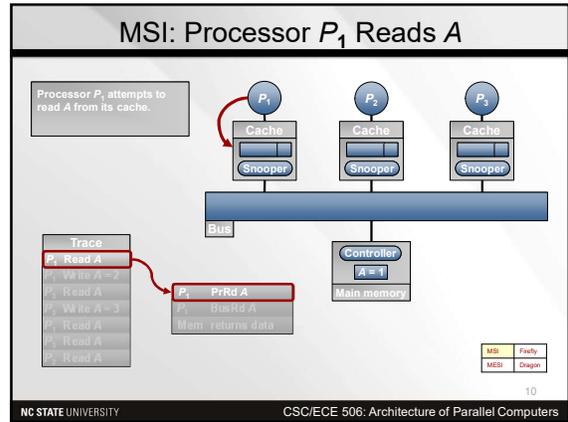
CSC/ECE 506: Architecture of Parallel Computers

6

6



7



10

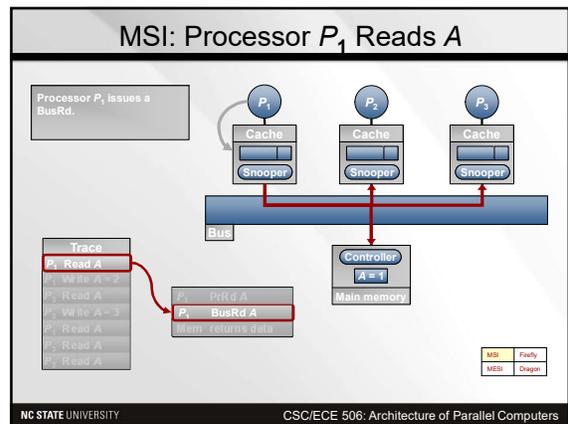
Lecture 15 Outline

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

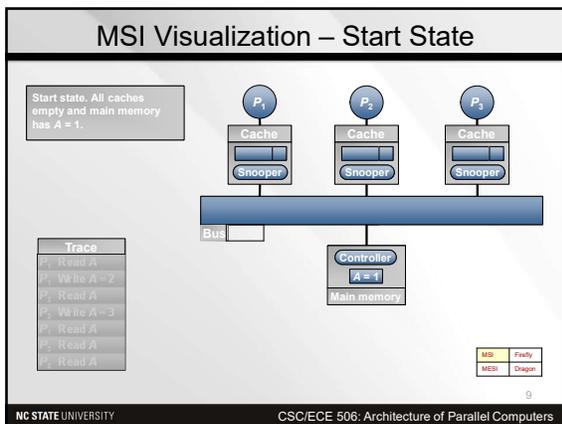
	Inval- idate	Update
3-state	MSI	Firefly
4-state	MESI	Dragon

NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

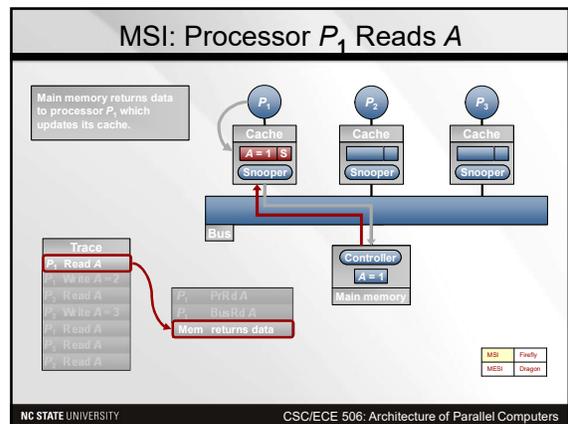
8



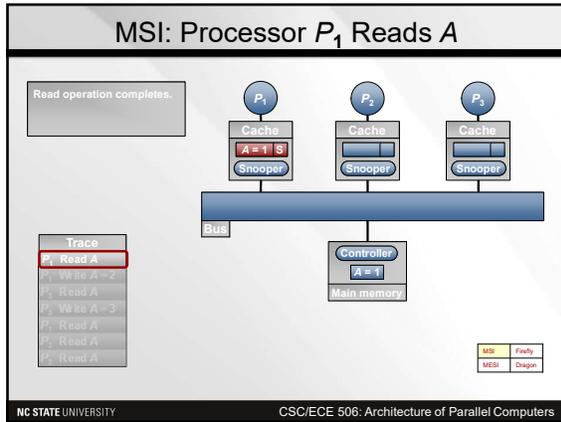
11



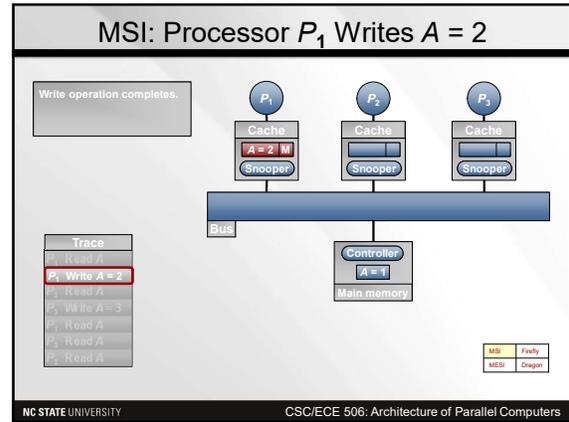
9



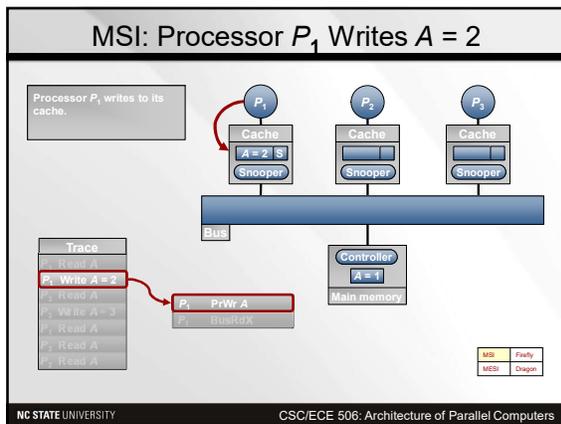
12



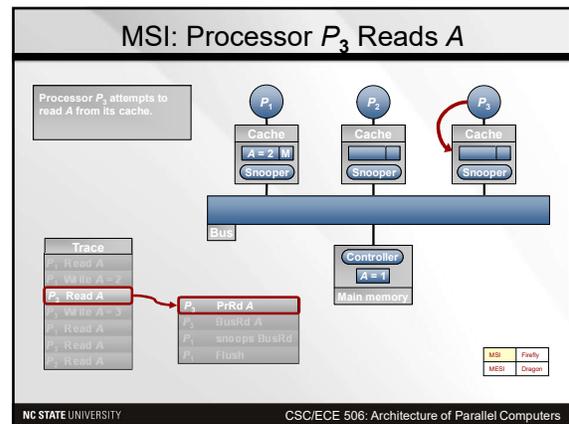
13



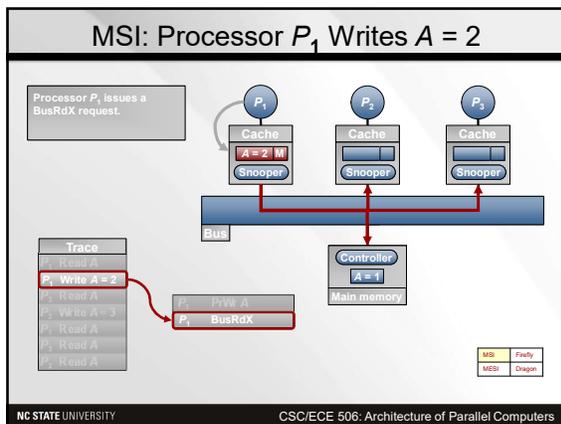
14



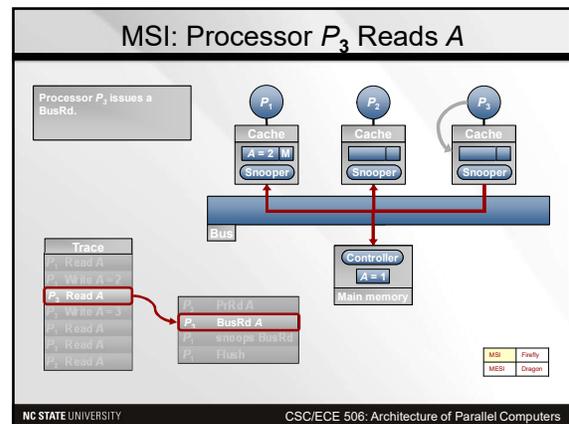
15



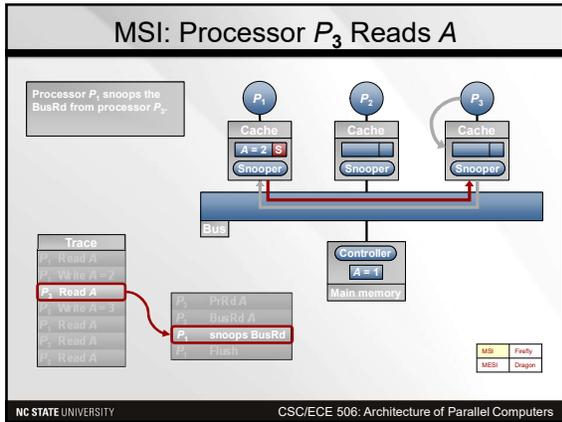
16



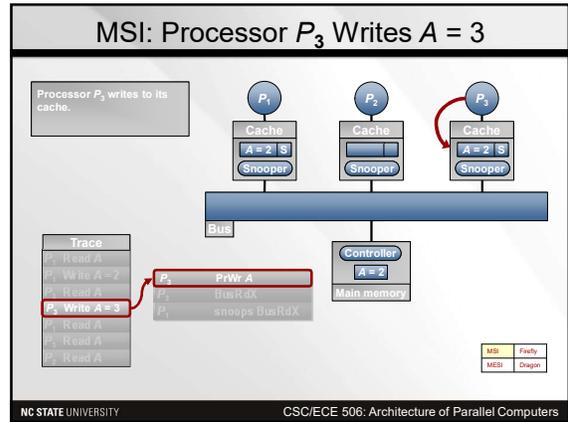
17



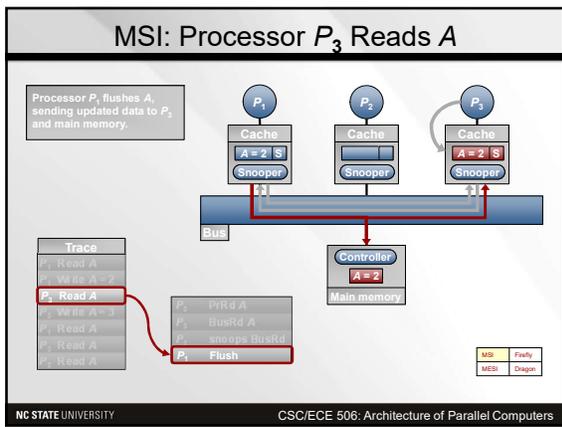
18



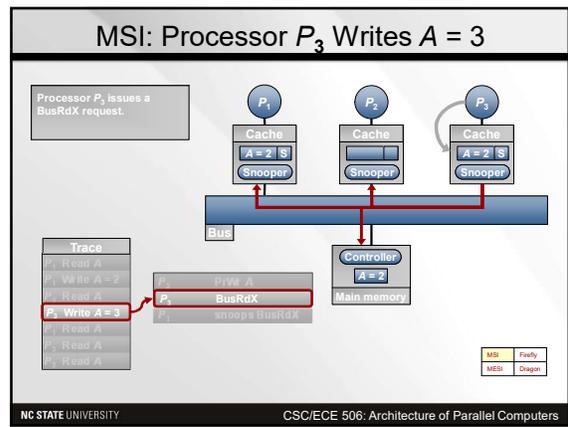
19



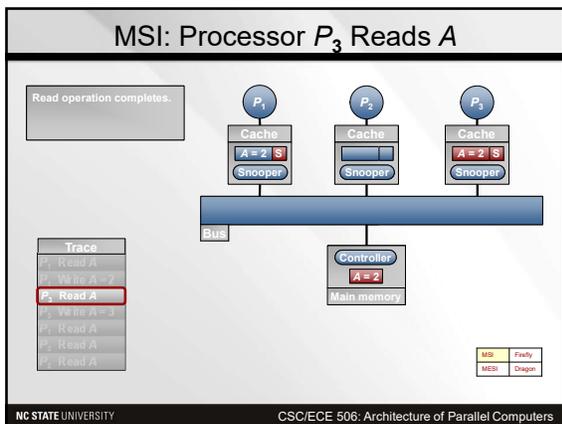
22



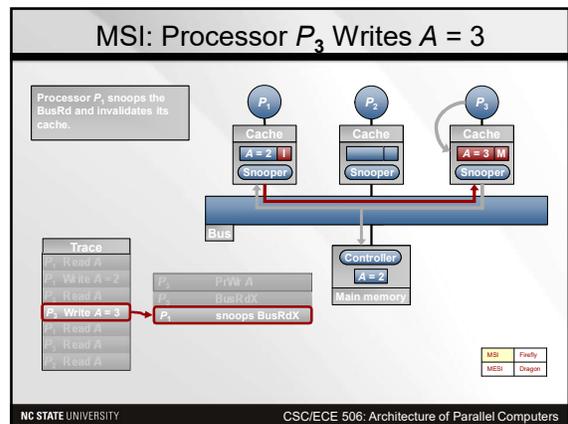
20



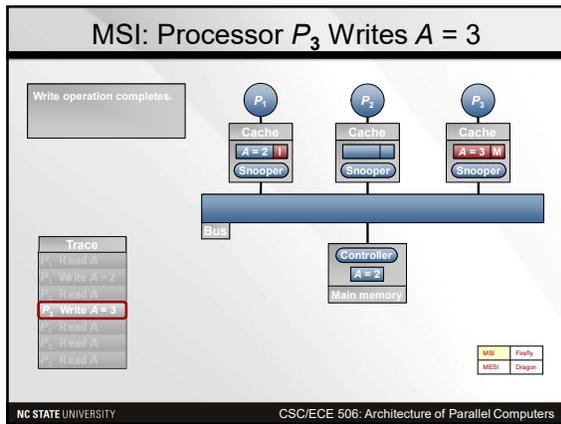
23



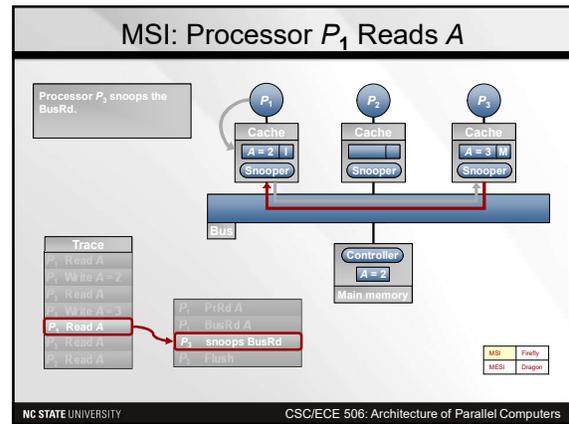
21



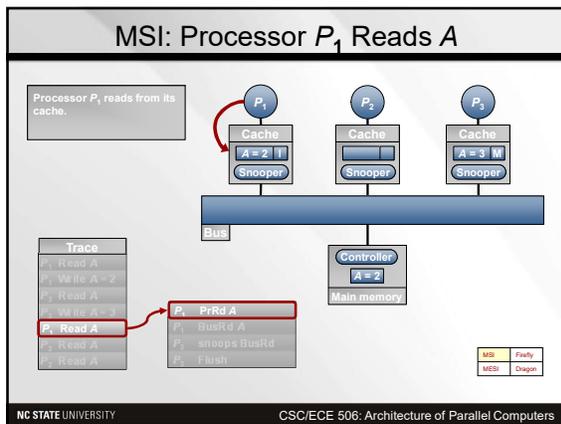
24



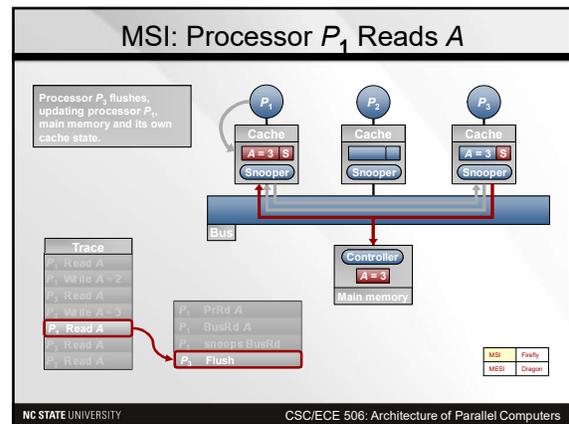
25



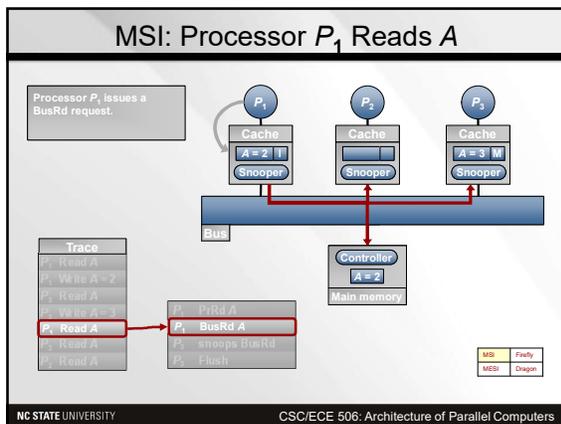
28



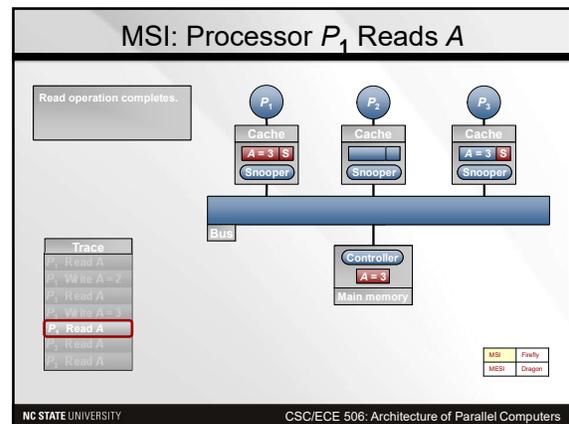
26



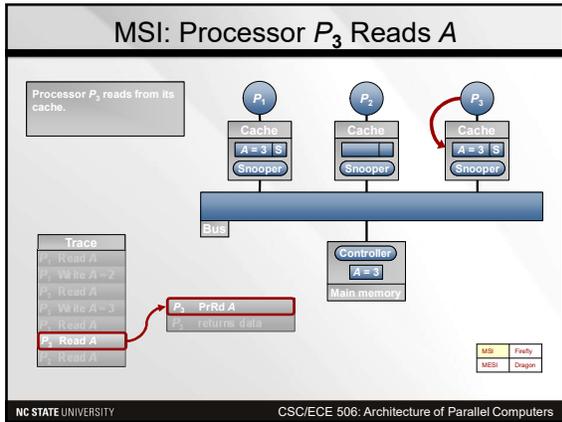
29



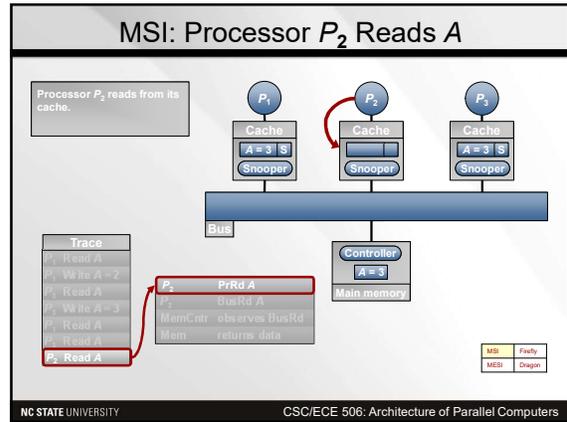
27



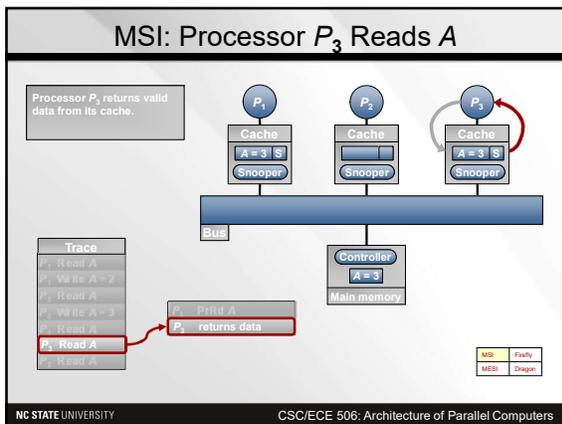
30



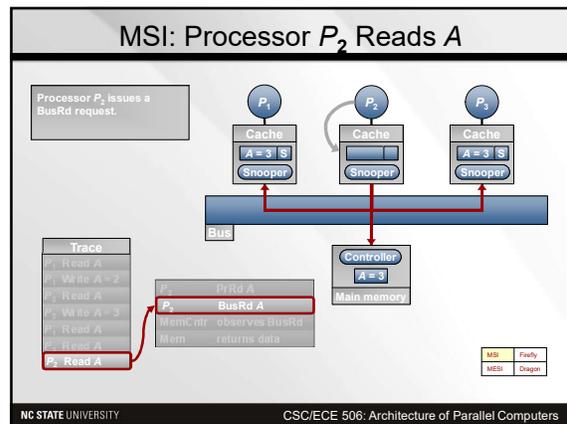
31



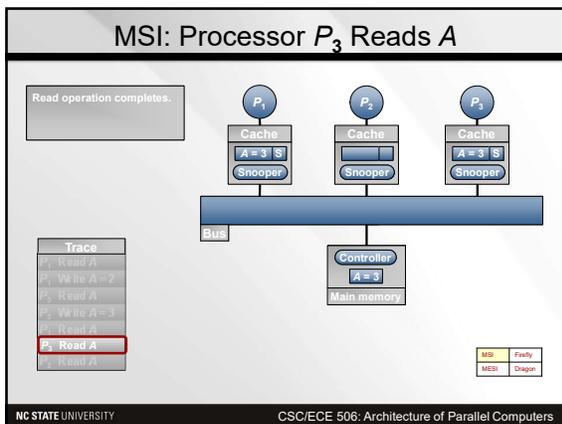
34



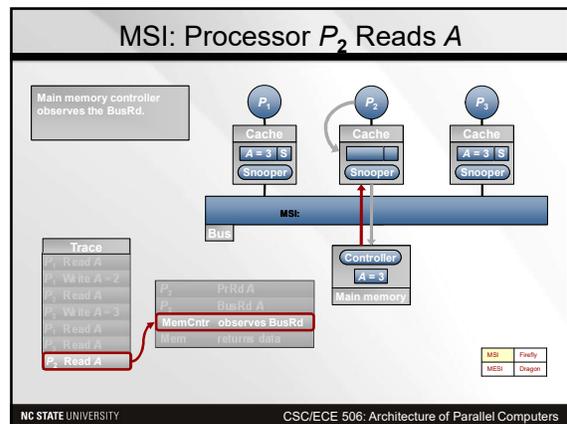
32



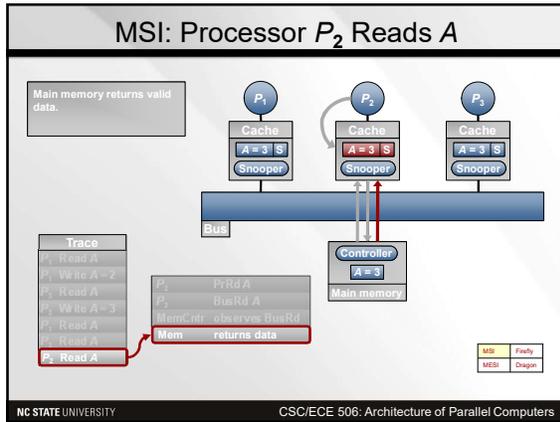
35



33



36



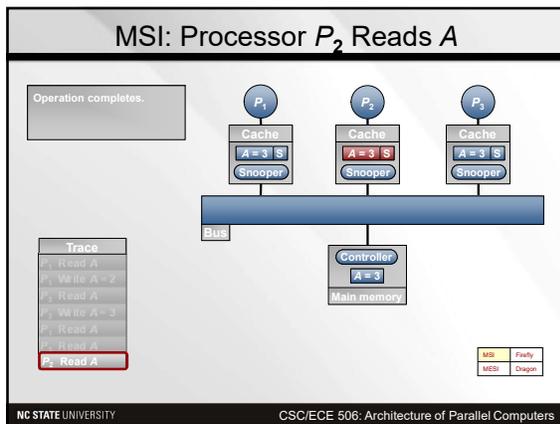
37

Notes on MSI Protocol

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

NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

40



38

Notes on MSI Protocol

- For $M \rightarrow I$, 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).

NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

41

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	M	-	-	BusRdX*	Mem
R3	S	-	S	BusRd/Flush	P1 cache
W3	I	-	M	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)

NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

39

Notes on MSI Protocol

- For $M \rightarrow I$, 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

NC STATE UNIVERSITY 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).

43

NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

43

MESI (4-state) Invalidation Protocol

- Here's a problem with the MSI protocol:
 - A {Rd, Wr} sequence causes two bus transactions
 - BusRd (I → S) followed by BusRdX or BusUpgr (S → 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)
- How does the protocol decide whether I → E or I → S?
 - Need to check whether someone else has a copy
 - "Shared" signal on bus: wired-or line asserted in response to BusRd



46

NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

46

Lecture 15 Outline

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

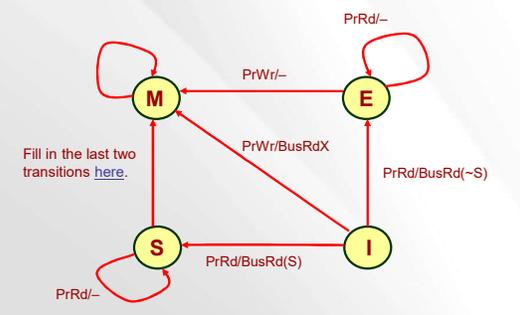
	Inval- idate	Update
3-state	MSI	Firefly
4-state	MESI	Dragon

44

NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

44

MESI: Processor-Initiated Transactions



47

NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

47

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?

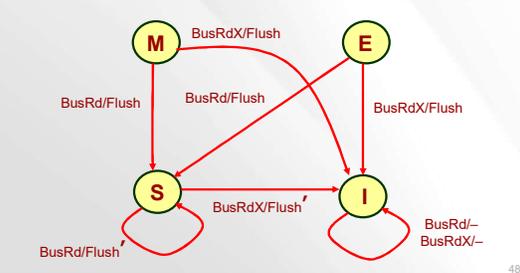
45

NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

45

MESI: Bus-Initiated Transactions

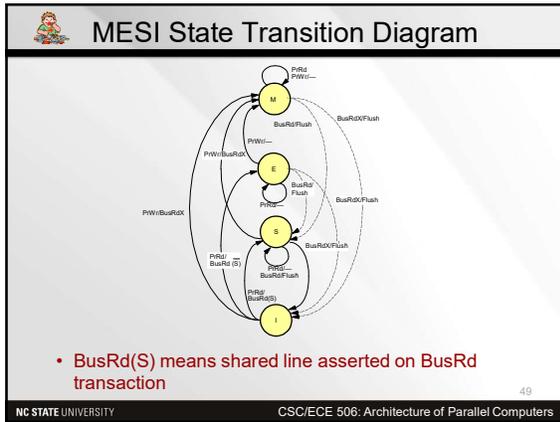
Flush' means flush only if cache-to-cache sharing is used; only the cache responsible for supplying the data will do a flush.



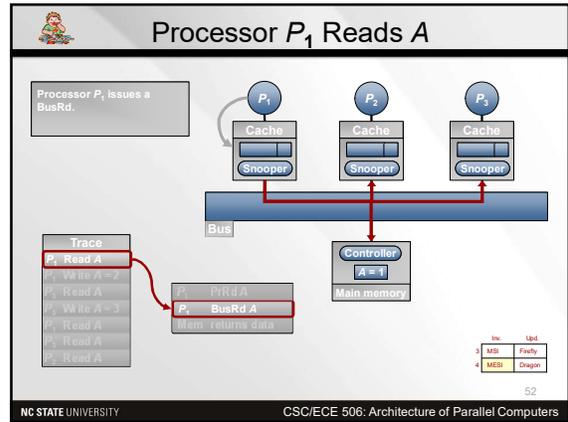
48

NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

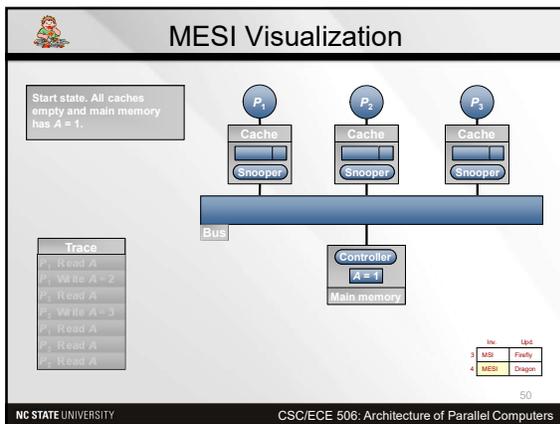
48



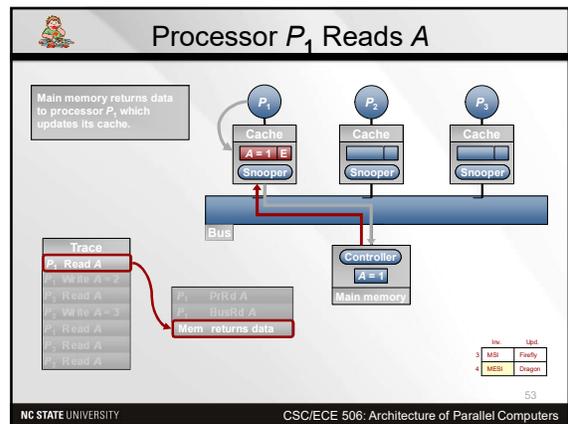
49



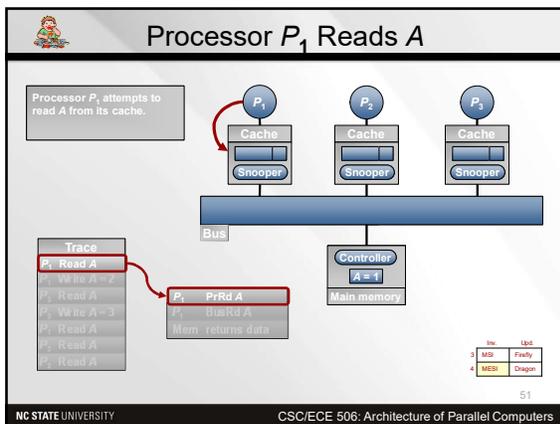
52



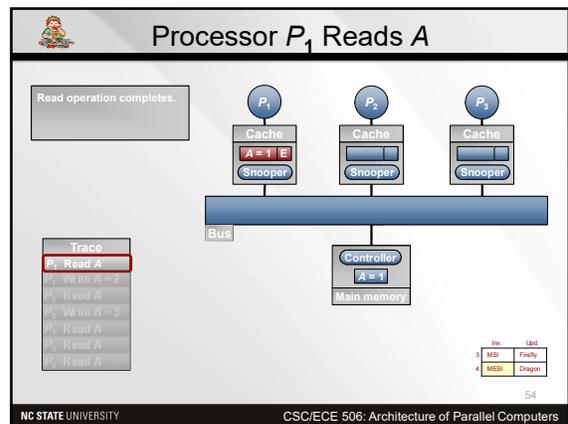
50



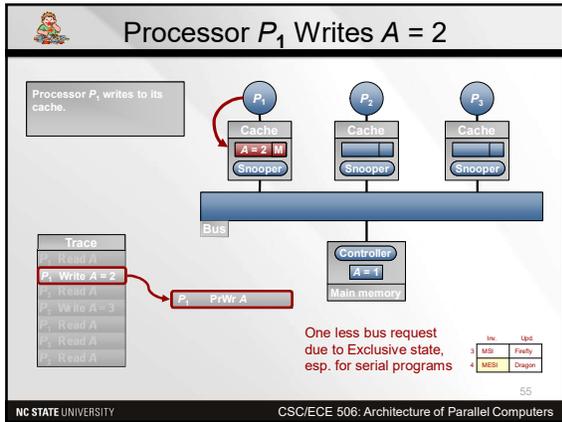
53



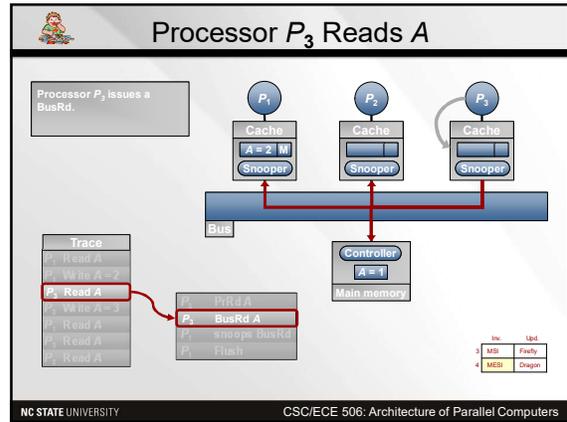
51



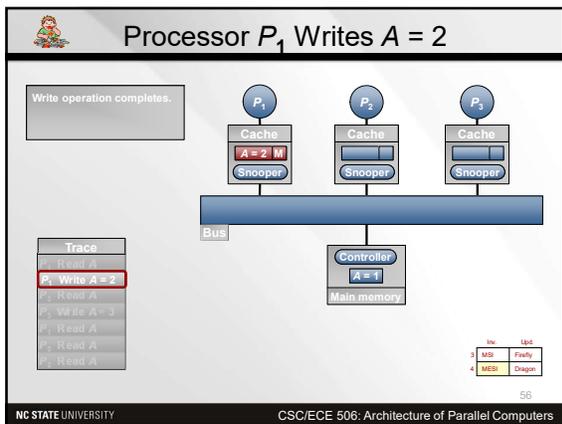
54



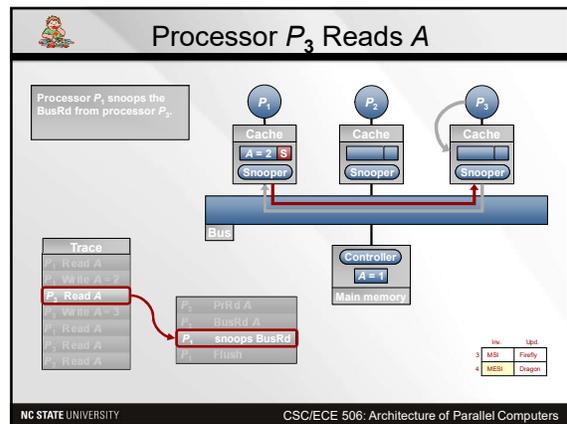
55



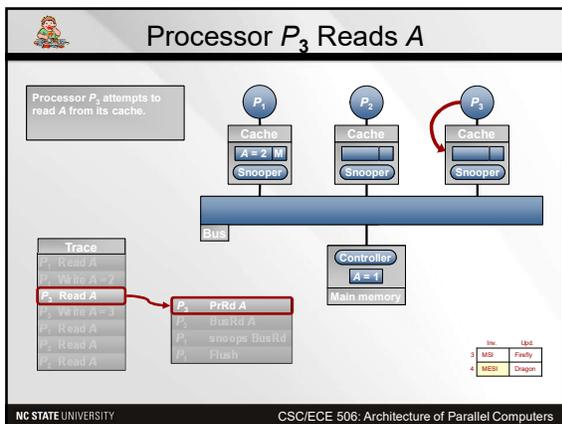
58



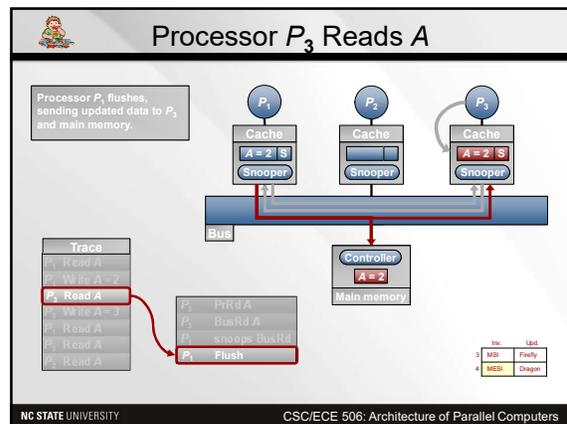
56



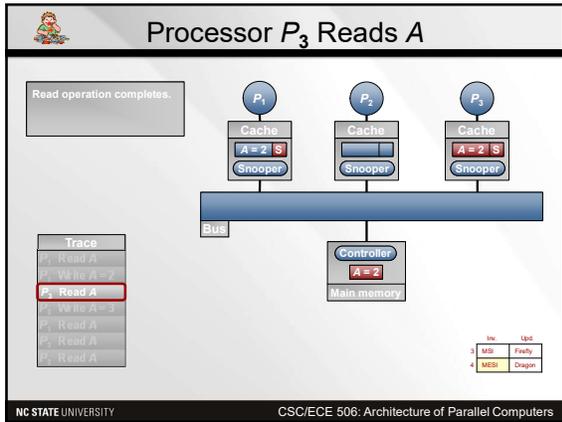
59



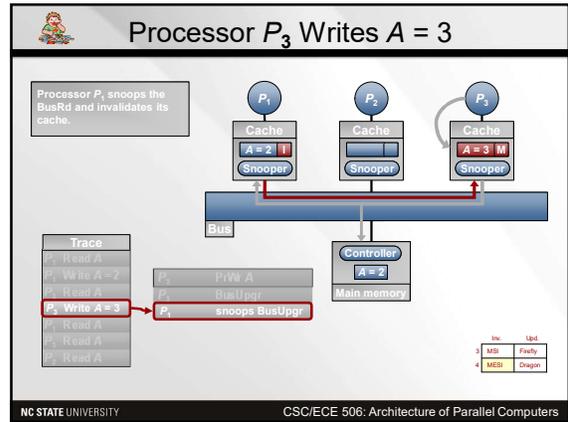
57



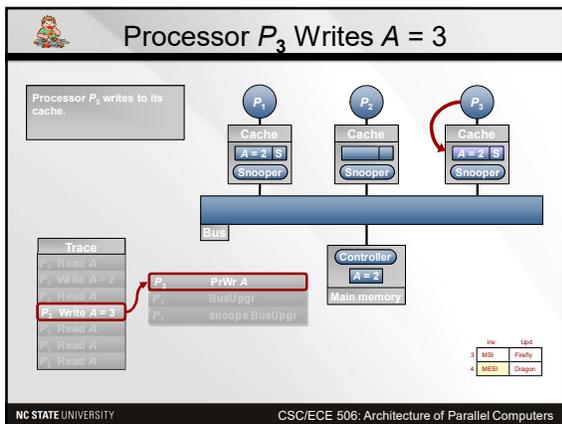
60



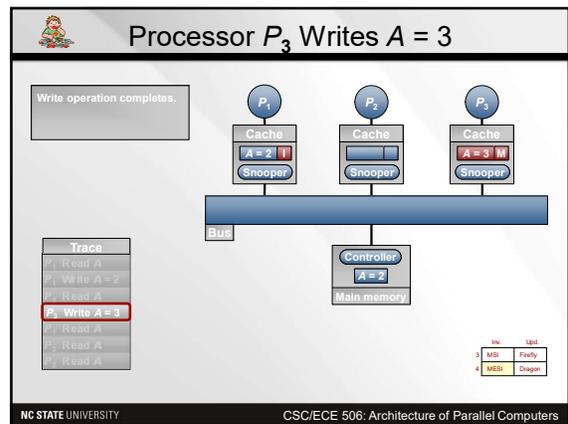
61



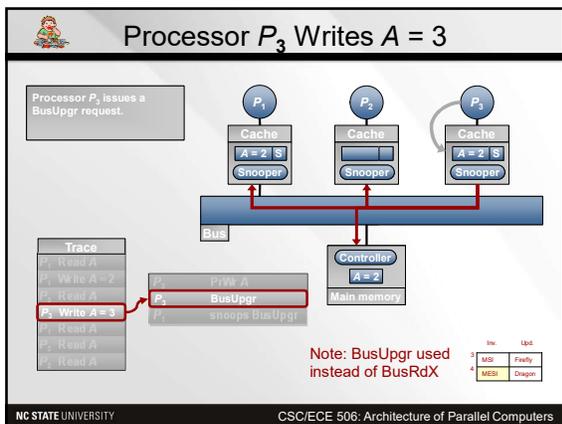
64



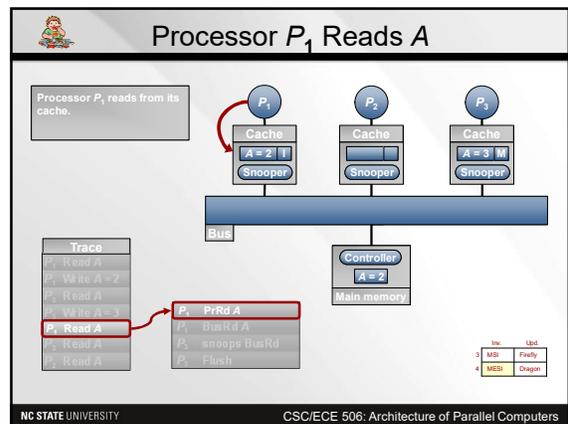
62



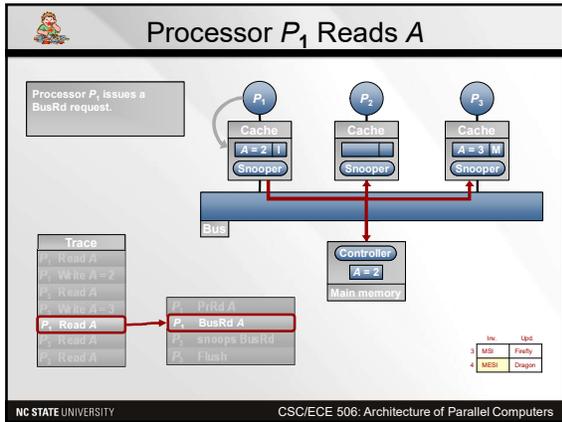
65



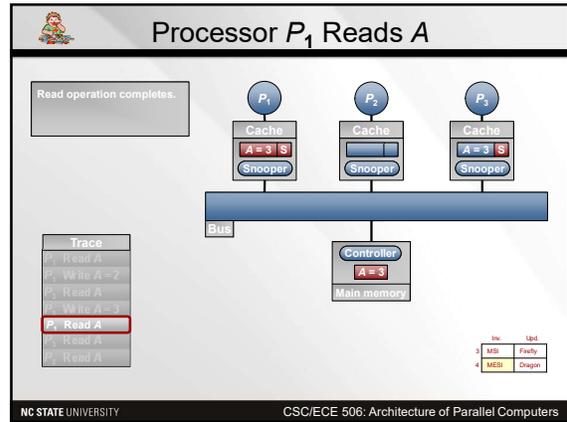
63



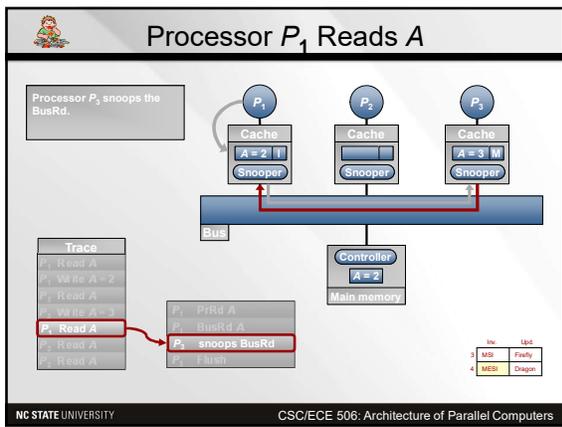
66



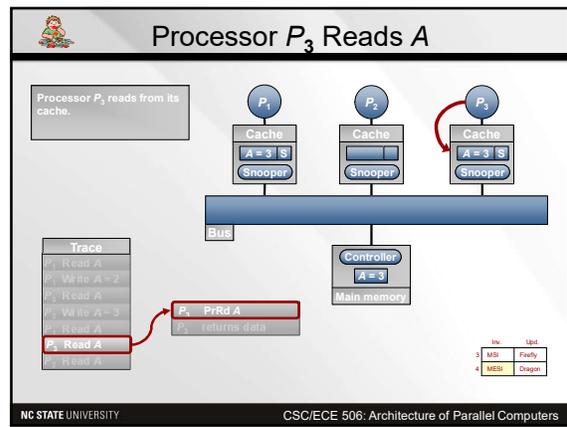
67



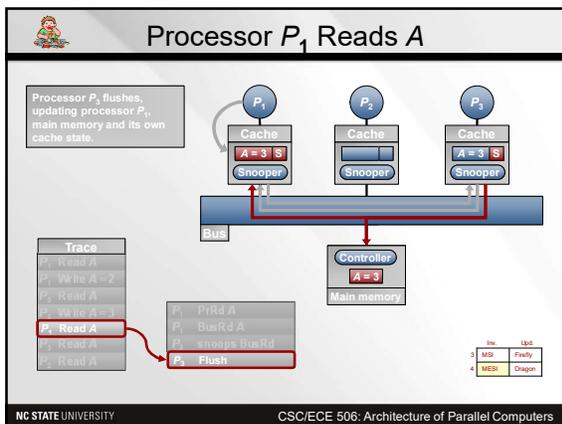
70



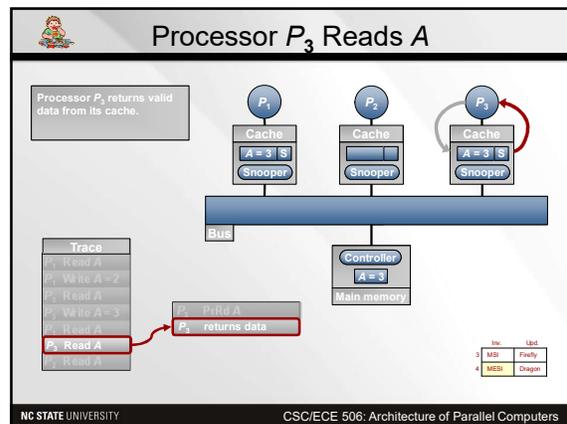
68



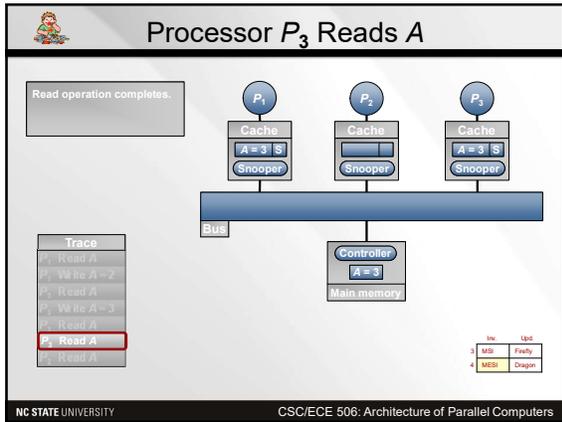
71



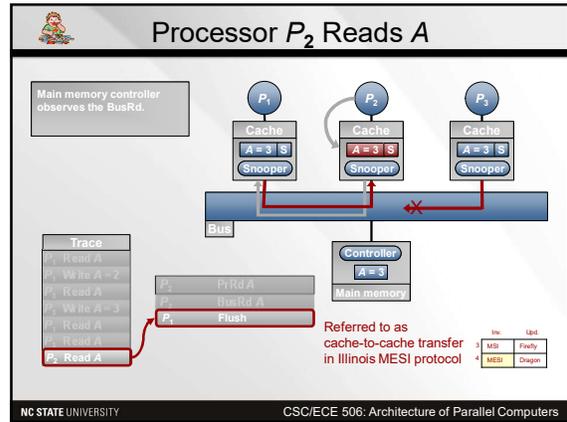
69



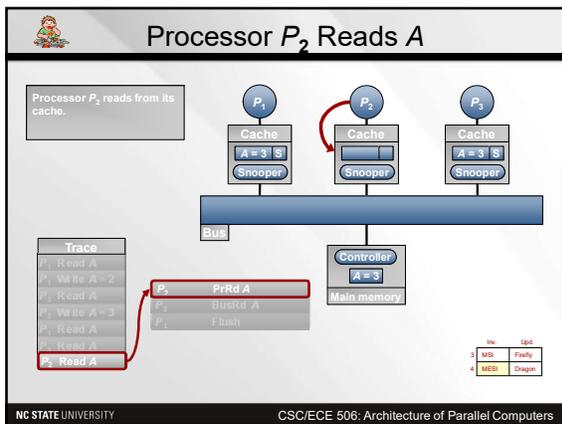
72



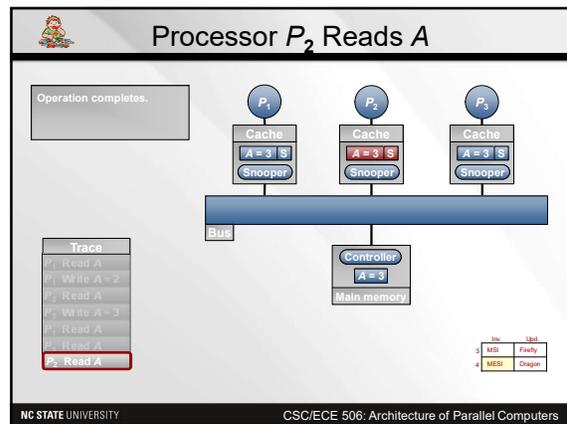
73



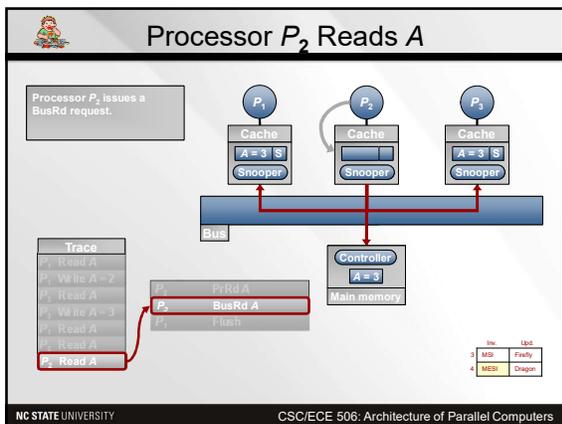
76



74



77



75

MESI Example (Cache-to-Cache Transfer)

Proc Action	State P1	State P2	State P3	Bus Action	Data From
R1	E	-	-	BusRd	Mem
W1	M	-	-	-	Own cache
R3	S	-	S	BusRd/Flush	P1 cache
W3	I	-	M	BusRdX	Mem
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 CSC/ECE 506: Architecture of Parallel Computers

78

Change from MSI (Cache-to-Cache Transfer)

Proc Action	State P1	State P2	State P3	Bus Action	Data From
R1	E	-	-	BusRd	Mem
W1	M	-	-	-	Own cache
R3	S	-	S	BusRd/Flush	P1 cache
W3	I	-	M	BusRdX	Mem
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 CSC/ECE 506: Architecture of Parallel Computers 79

79

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

Proc Action	State P1	State P2	State P3	Bus Action	Data From
R1	E	-	-	BusRd	Mem
W1	M	-	-	-	Own cache
R3	S	-	S	BusRd/Flush	P1 cache
W3	I	-	M	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 CSC/ECE 506: Architecture of Parallel Computers 82

82

Change from MSI (Cache-to-Cache Transfer)

Proc Action	State P1	State P2	State P3	Bus Action	Data From
R1	E	-	-	BusRd	Mem
W1	M	-	-	-	Own cache
R3	S	-	S	BusRd/Flush	P1 cache
W3	I	-	M	BusRdX	Mem
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 CSC/ECE 506: Architecture of Parallel Computers 80

80

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

Proc Action	State P1	State P2	State P3	Bus Action	Data From
R1	E	-	-	BusRd	Mem
W1	M	-	-	-	Own cache
R3	S	-	S	BusRd/Flush	P1 cache
W3	I	-	M	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 CSC/ECE 506: Architecture of Parallel Computers 83

83

Change from MSI (Cache-to-Cache Transfer)

Proc Action	State P1	State P2	State P3	Bus Action	Data From
R1	E	-	-	BusRd	Mem
W1	M	-	-	-	Own cache
R3	S	-	S	BusRd/Flush	P1 cache
W3	I	-	M	BusRdX	Mem
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 CSC/ECE 506: Architecture of Parallel Computers 81

81

- ### 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)
- NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers 84

84

Lecture 15 Outline

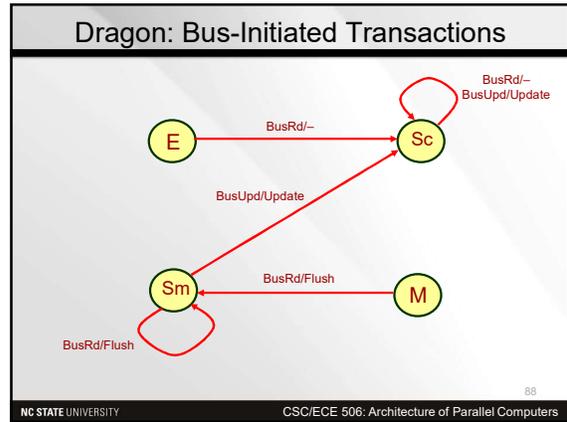
- MSI protocol
- MESI protocol
- Dragon protocol
- Firefly protocol

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

85

NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

85



88

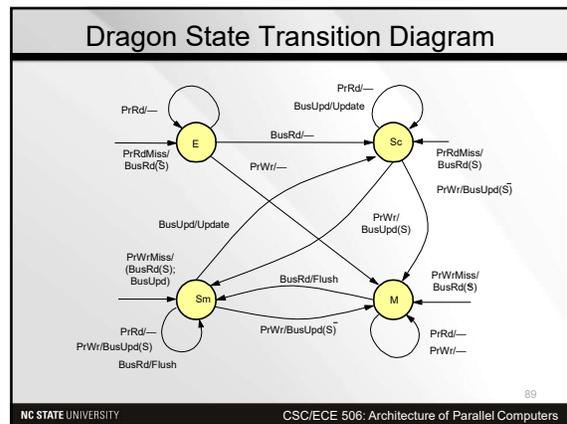
Dragon Writeback Update Protocol

- 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

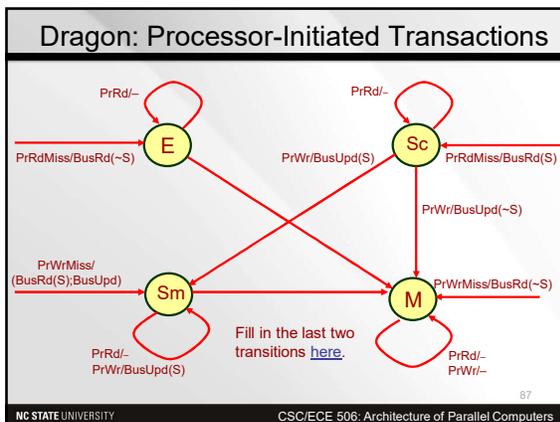
86

NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

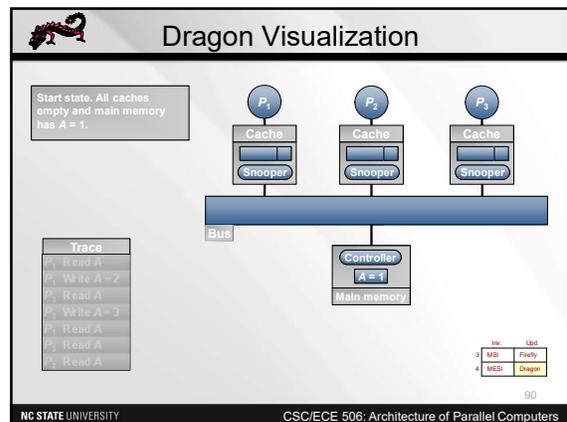
86



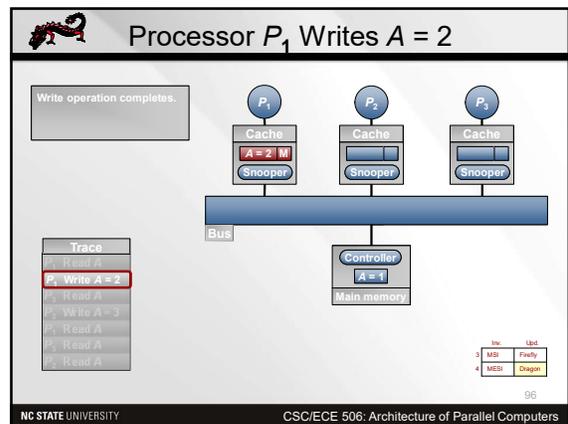
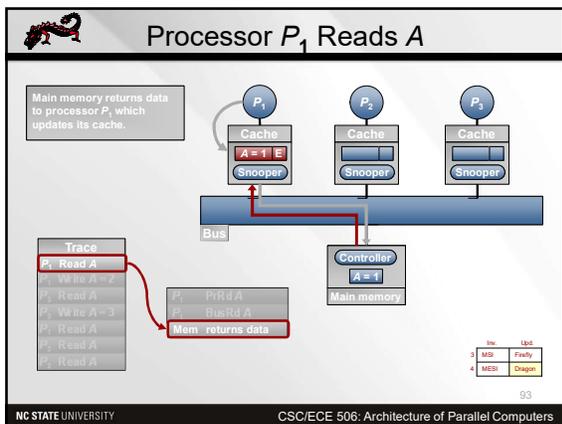
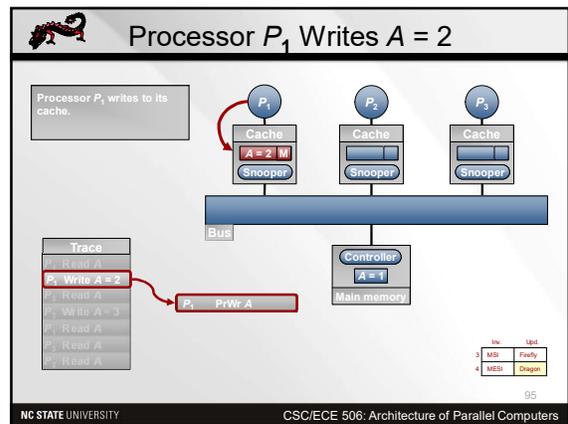
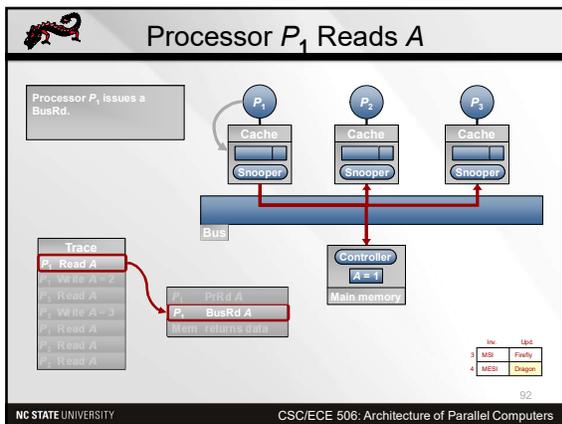
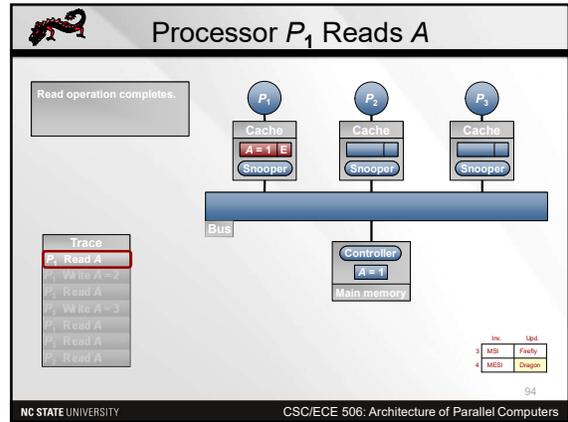
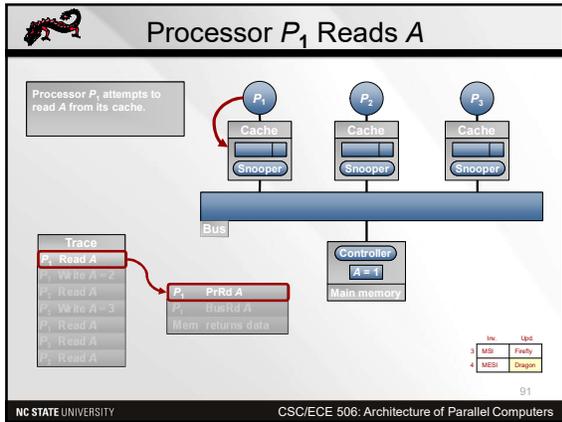
89

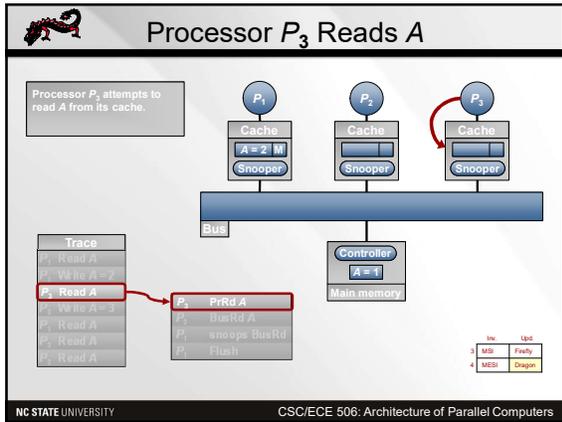


87

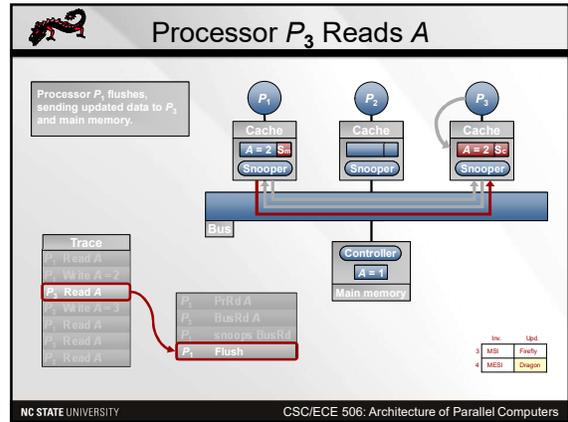


90

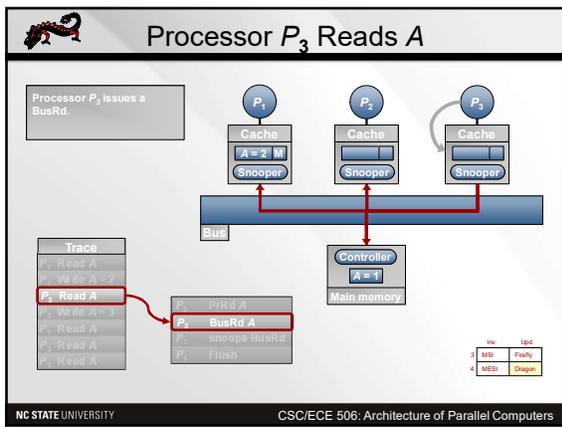




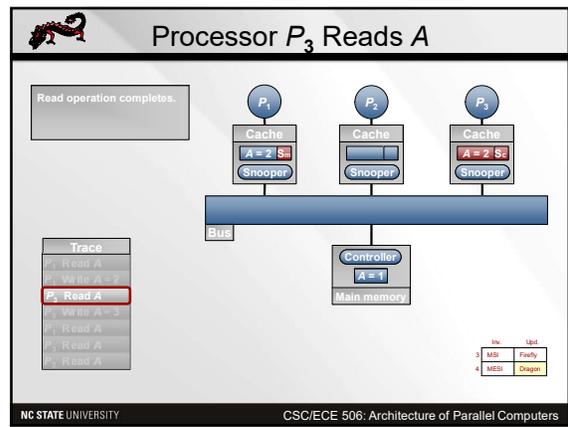
97



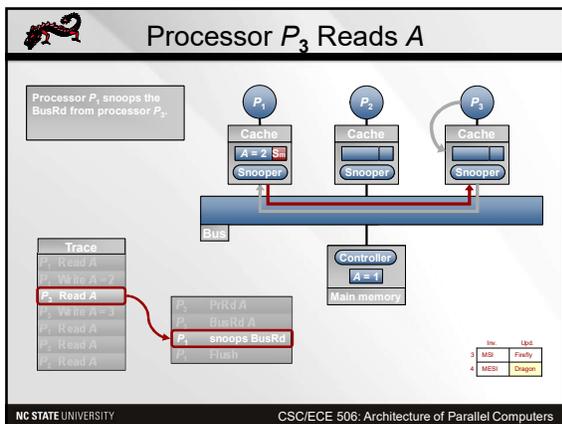
100



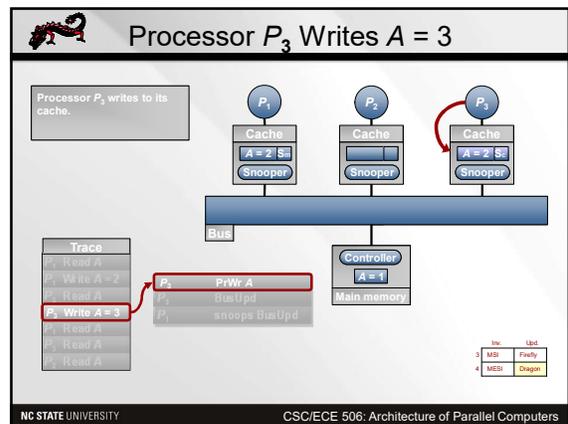
98



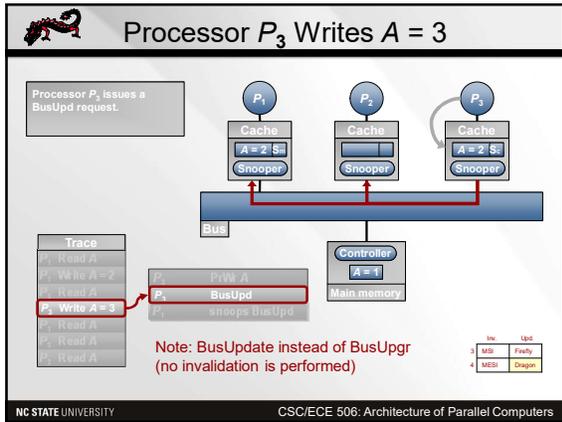
101



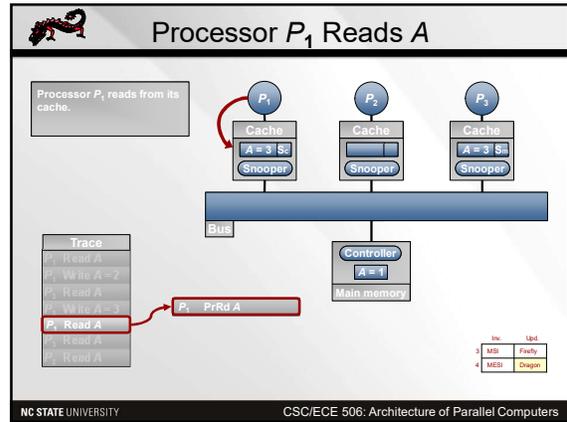
99



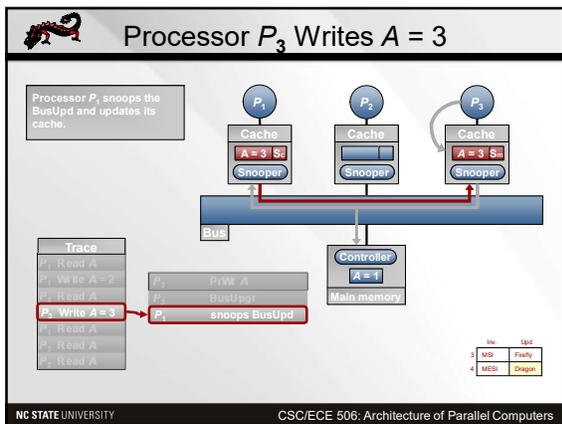
102



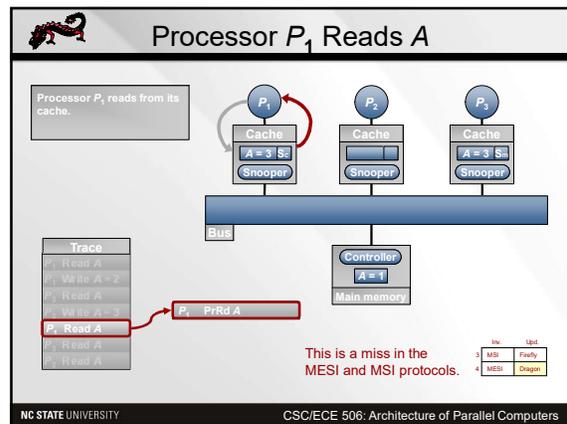
103



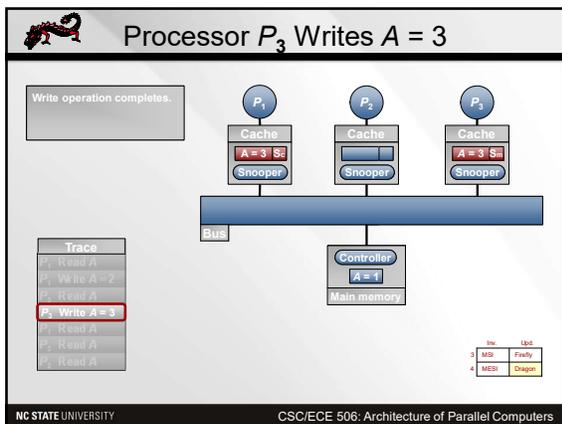
106



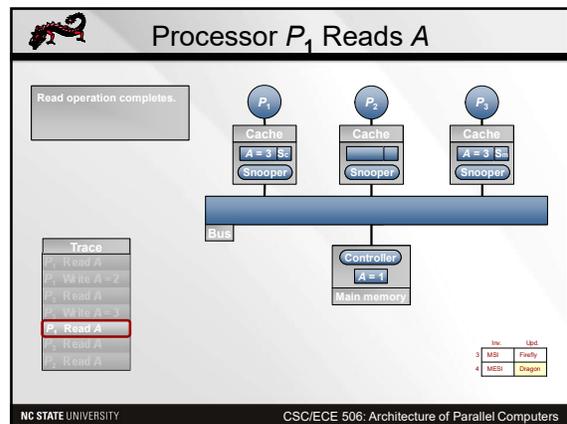
104



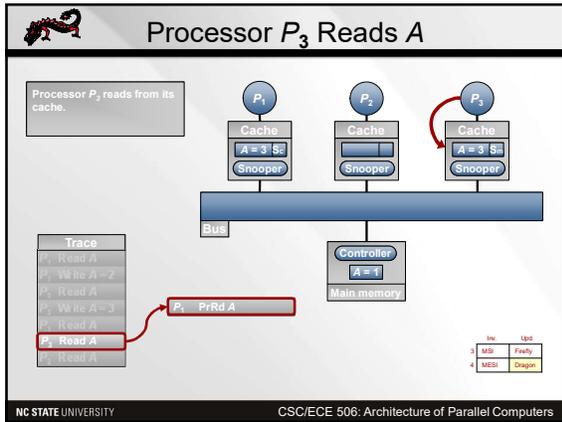
107



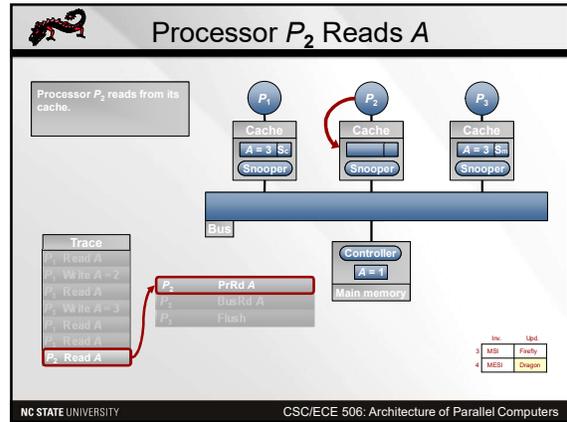
105



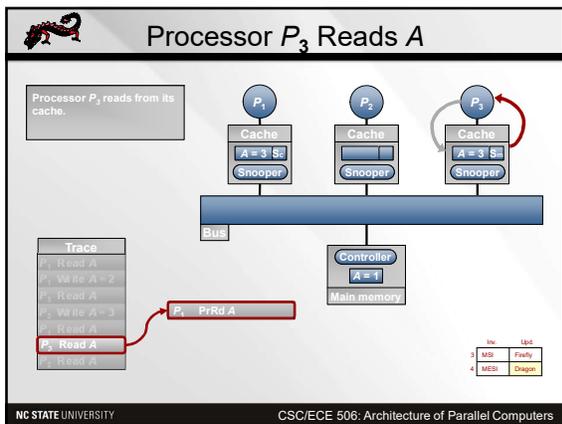
108



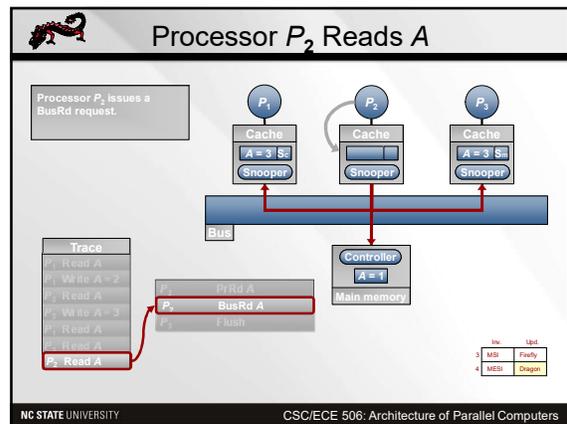
109



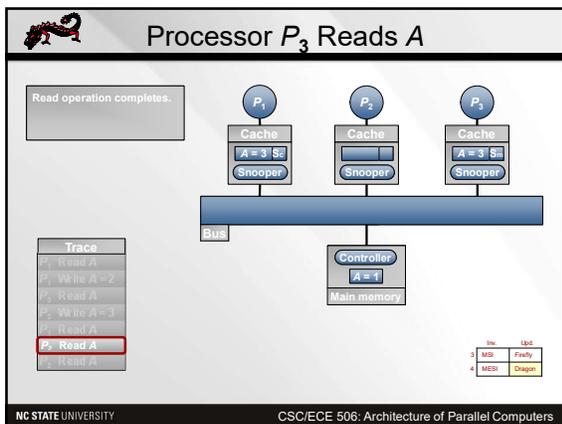
112



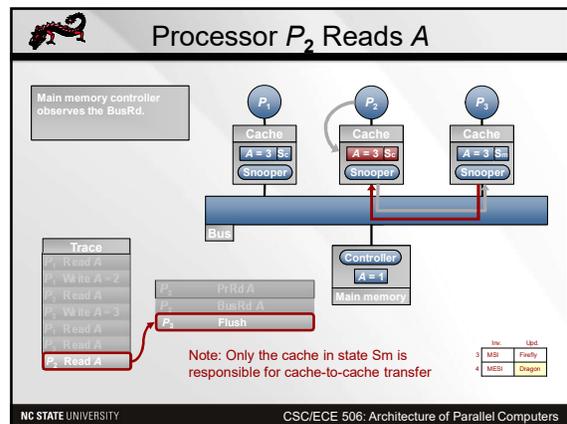
110



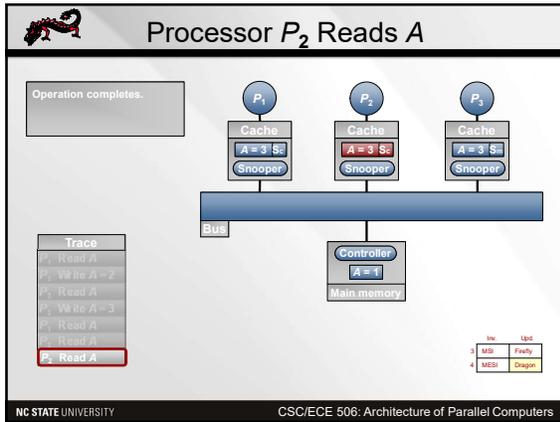
113



111



114



115

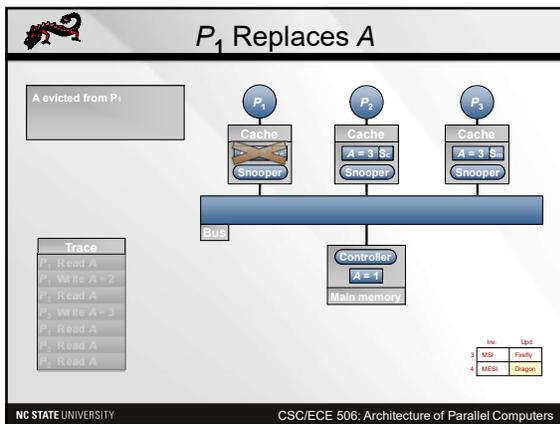
Dragon Example

Proc Action	State P1	State P2	State P3	Bus Action	Data from
R1	E	-	-	BusRd	Mem
W1	M	-	-	-	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

118 119
3 MSI Firefly
4 MESI Dragon

NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

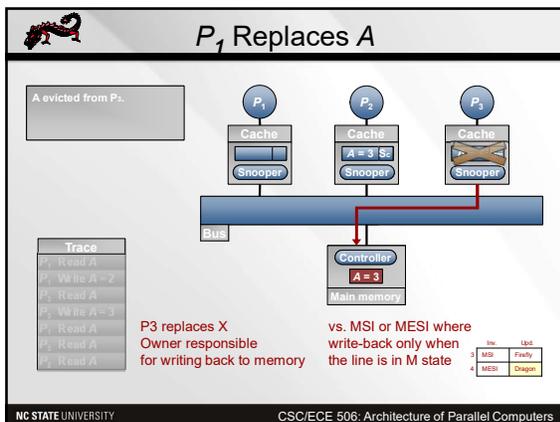
118



116

- ### 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.
- 118 119
3 MSI Firefly
4 MESI Dragon
- NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

119



117

Lecture 15 Outline

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

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

118 119
3 MSI Firefly
4 MESI Dragon

NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

120

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?



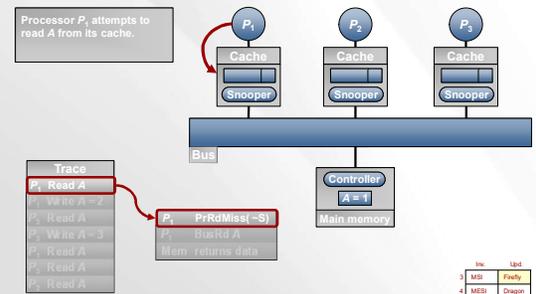
- The Firefly protocol, named after a multiprocessor workstation developed by DEC, is an example of such a protocol.

NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

121

Processor P_1 Reads A

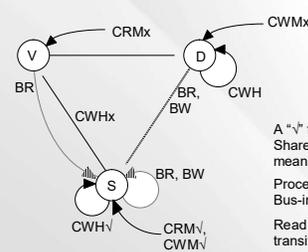
Processor P_1 attempts to read A from its cache.



NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

124

Firefly State-Transition Diagram



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

A "*" 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.

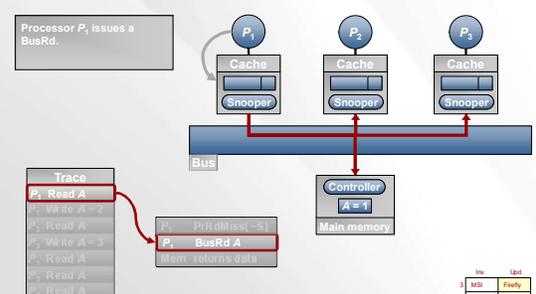
- Answer some questions about this diagram.

NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

122

Processor P_1 Reads A

Processor P_1 issues a BusRd.

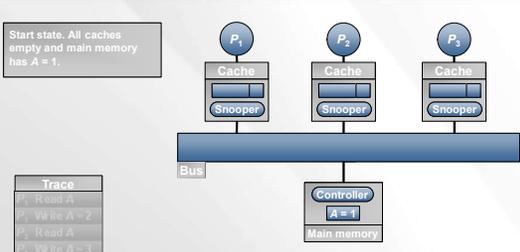


NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

125

Firefly Visualization

Start state. All caches empty and main memory has A = 1.

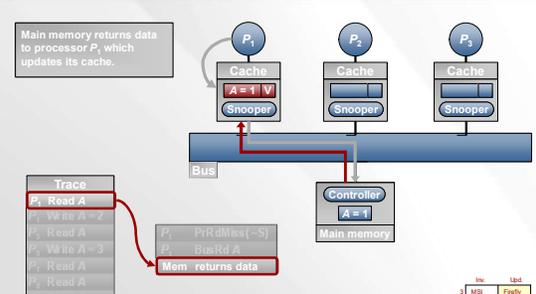


NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

123

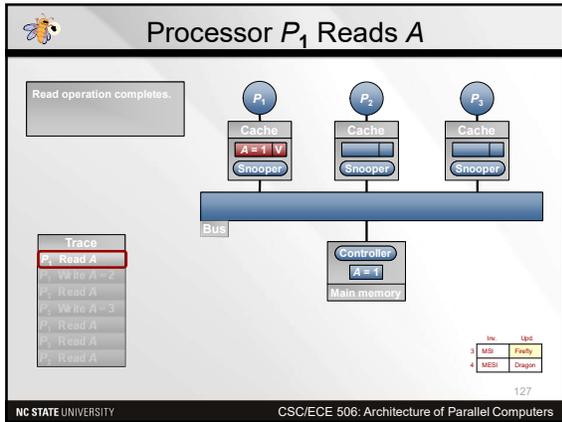
Processor P_1 Reads A

Main memory returns data to processor P_1 which updates its cache.

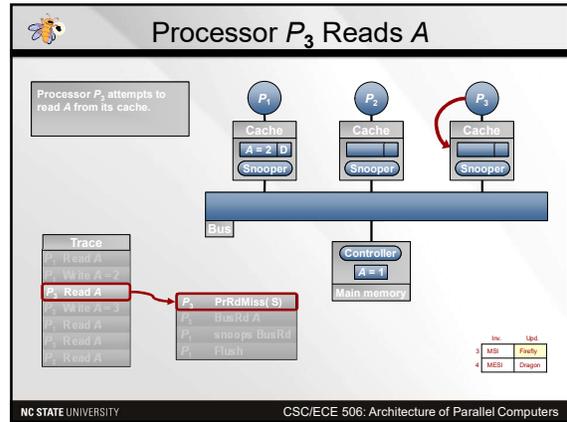


NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

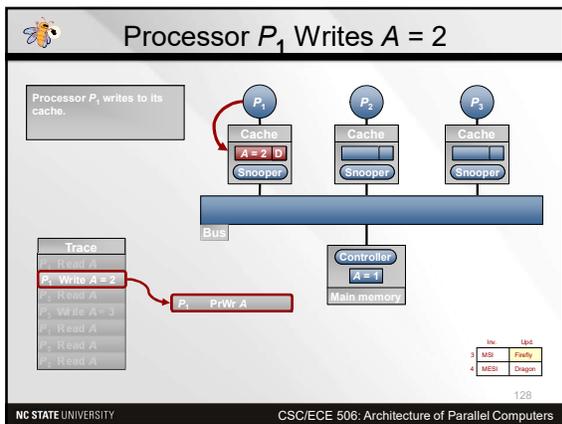
126



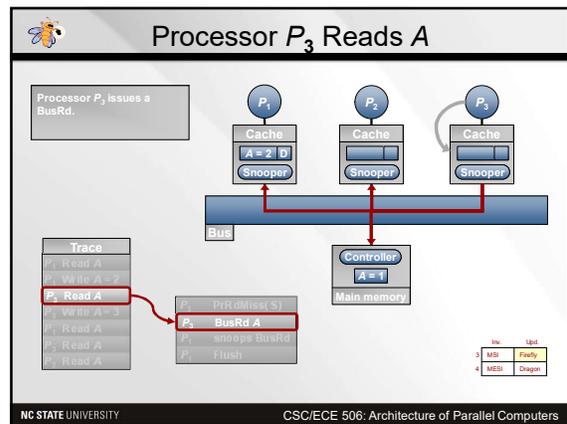
127



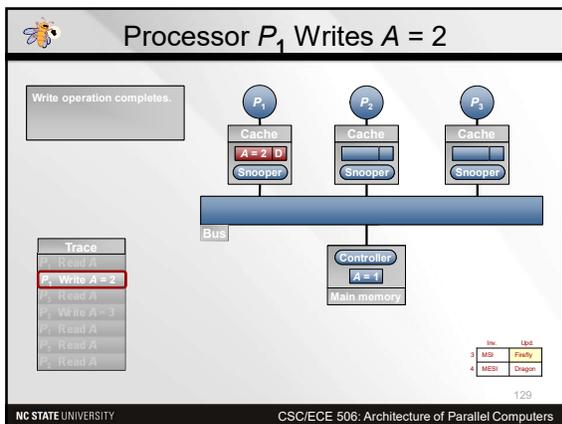
130



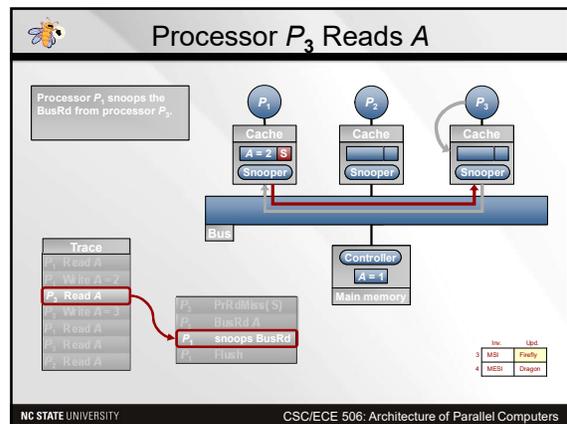
128



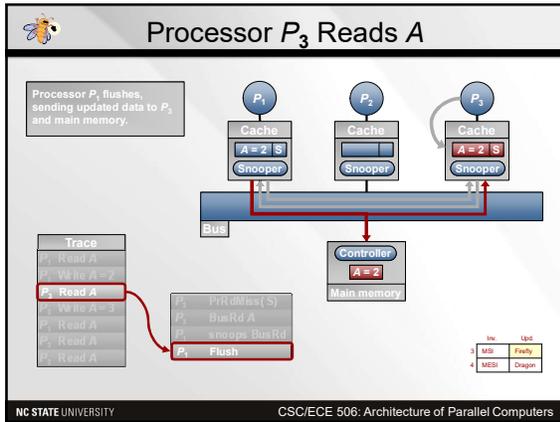
131



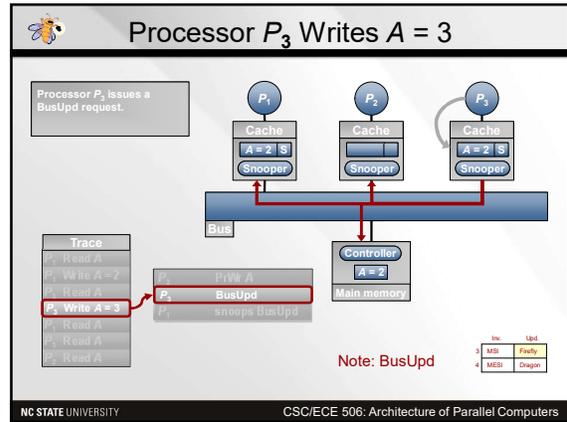
129



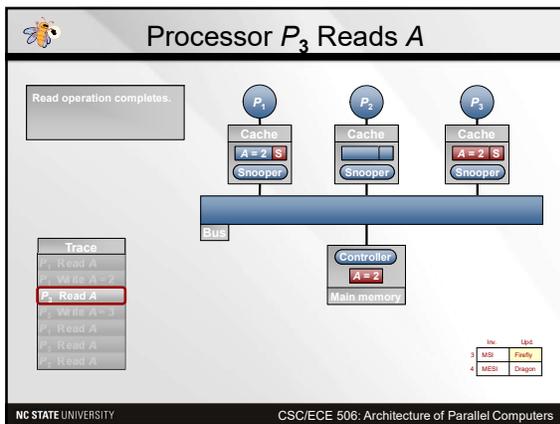
132



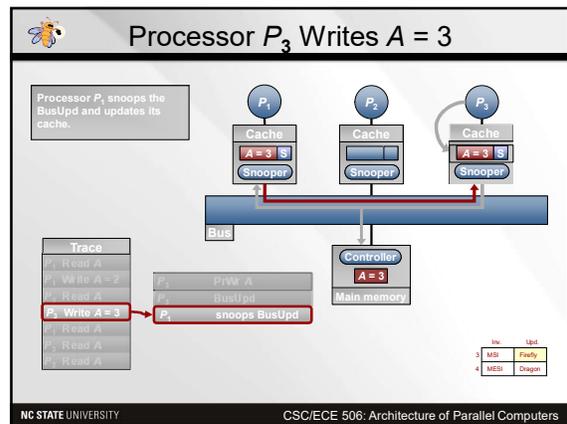
133



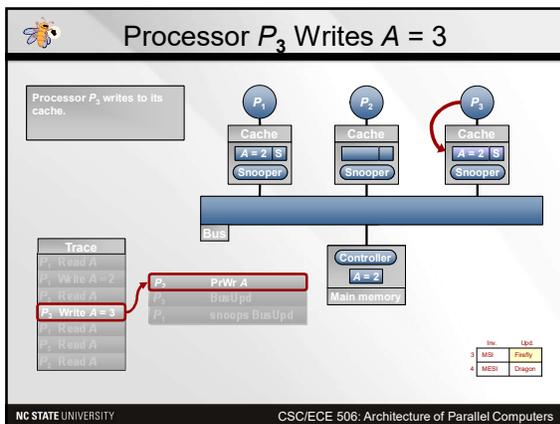
136



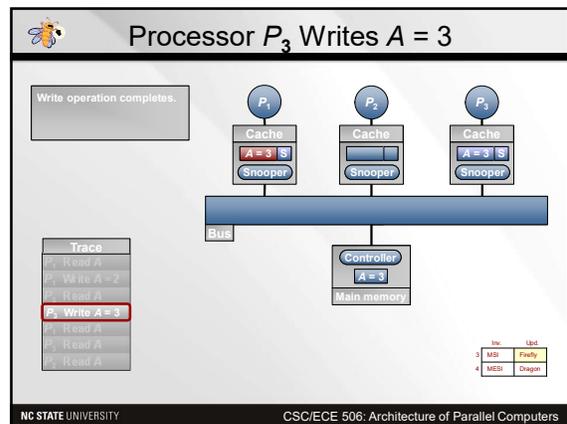
134



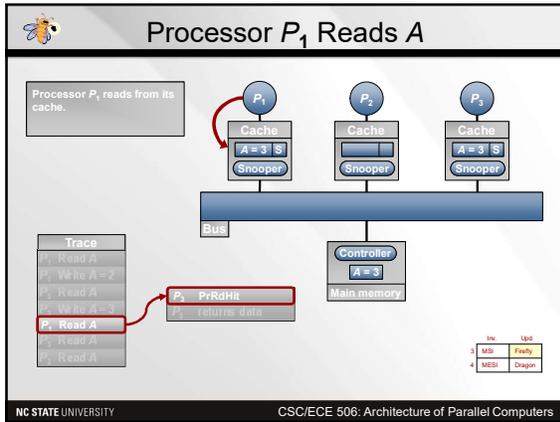
137



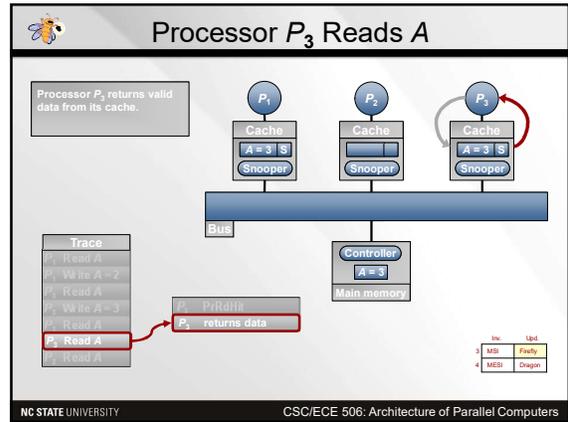
135



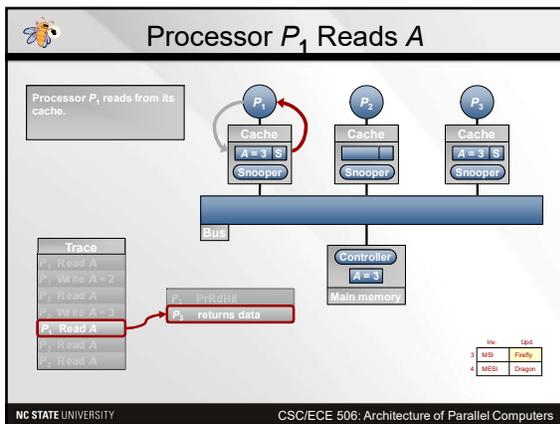
138



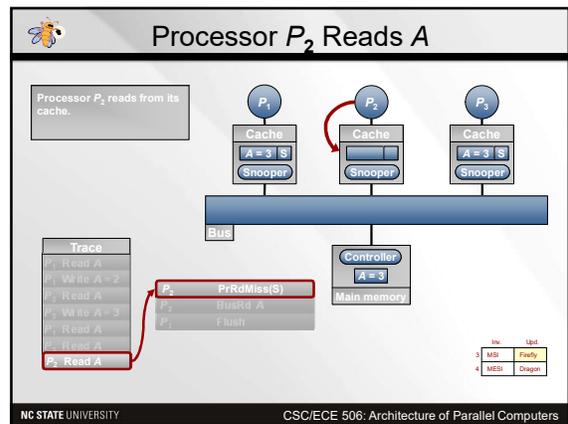
139



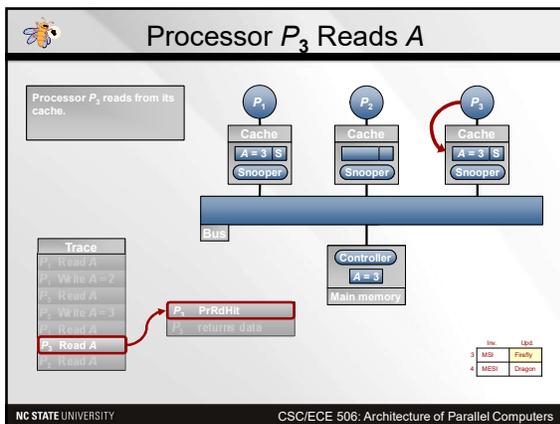
142



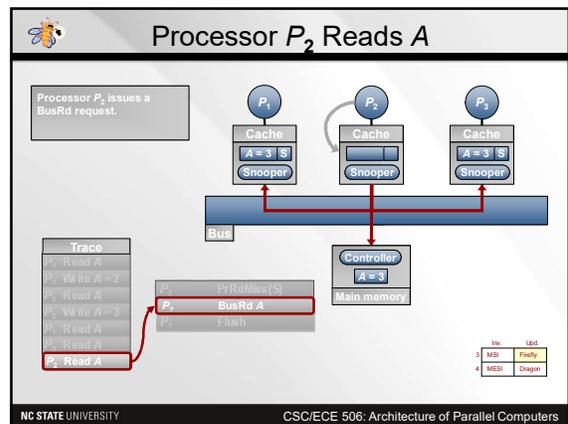
140



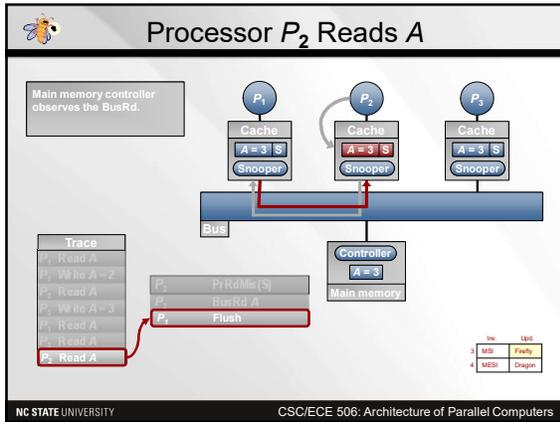
143



141



144



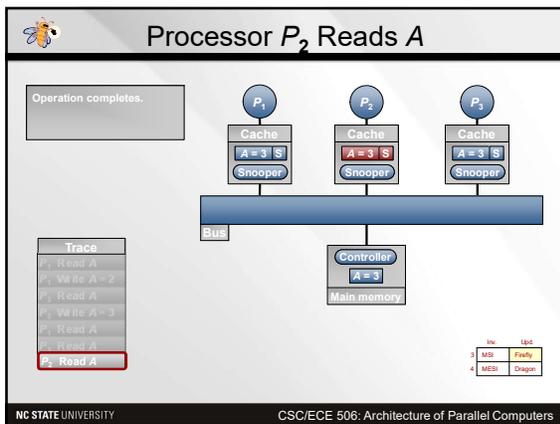
145

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

NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

148



146

- ### 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
- NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

149

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

NC STATE UNIVERSITY CSC/ECE 506: Architecture of Parallel Computers

147